



Outils et environnements pour l'amélioration incrémentale, la post-édition contributive et l'évaluation continue de systèmes de TA. Application à la TA français-chinois.

Lingxiao Wang

► To cite this version:

Lingxiao Wang. Outils et environnements pour l'amélioration incrémentale, la post-édition contributive et l'évaluation continue de systèmes de TA. Application à la TA français-chinois.. Traitement du texte et du document. Université Grenoble Alpes, 2015. Français. <NNT : 2015GREAM057>. <tel-01320566>

HAL Id: tel-01320566

<https://tel.archives-ouvertes.fr/tel-01320566>

Submitted on 24 May 2016

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

THÈSE

Pour obtenir le grade de

DOCTEUR DE L'UNIVERSITÉ GRENOBLE ALPES

Spécialité : **Informatique**

Arrêté ministériel :

Présentée par

Lingxiao WANG

Thèse dirigée par **Christian BOITET**
et codirigée par **Valérie BELLYNCK**

préparée au sein du **Laboratoire d'informatique de Grenoble**
dans l'**École Doctorale « Mathématiques, Sciences et Technologies de l'Information, informatique »**

Outils et environnements pour l'amélioration incrémentale, la post- édition contributive et l'évaluation continue de systèmes de TA. Application à la TA français-chinois.

Thèse soutenue publiquement le **14 décembre 2015**,
devant le jury composé de :

Prof. Catherine BERRUT

Professeur, Université Joseph Fourier, Présidente/Examinatrice

Prof. Xiaodong SHI

Professeur, Université de Xiamen, Rapporteur.

Prof. Eric WEHRLI

Professeur, Université de Genève, Rapporteur.

Dr. Hong-Thai NGUYEN

Ingénieur de recherche, DHATIM, Examineur.

Prof. Christian BOITET

Professeur, Université Joseph Fourier, Directeur.

MdC. Valérie BELLYNCK

Maître de conférence, Grenoble-INP, Co-Directrice.



UNIVERSITÉ DE GRENOBLE

N° attribué par la bibliothèque

/ / / / / / / / / / / / / / /

THÈSE

pour obtenir le grade de

DOCTEUR ÈS SCIENCES

délivré par l'**UNIVERSITÉ DE GRENOBLE**

Spécialité : “INFORMATIQUE”

Thèse préparée au laboratoire GETALP-LIG (CNRS-INPG-UJF) dans le cadre de
l'École Doctorale “Mathématiques, Sciences et Technologies de l'Information, l'informatique”

présentée et soutenue publiquement

par

Lingxiao WANG

Le 14 décembre 2015

**Outils et environnements pour l'amélioration incrémentale,
la post-édition contributive et l'évaluation continue de
systèmes de TA. Application à la TA français-chinois.**

JURY

Prof. Catherine BERRUT (LIG)

Prof. Xiaodong SHI (Xiamen)

Prof. Eric WEHRLI (Unige)

Dr. Hong-Thai NGUYEN

Prof. Christian BOITET

MdC. Valérie BELLYNCK

Présidente/Examineur

Rapporteur

Rapporteur

Examineur

Directeur de thèse

Codirectrice de thèse

Résumé

Cette thèse, effectuée dans le cadre d'une bourse CIFRE, et prolongeant un des aspects du projet ANR TRAQUIERO, aborde d'abord la production, l'extension et l'amélioration de corpus multilingues par traduction automatique (TA) et post-édition contributive (PE). Des améliorations fonctionnelles et techniques ont été apportées aux logiciels SECTRA_W et IMAG, et on a progressé vers une définition générique de la structure d'un corpus multilingue, multi-annoté et multimédia, pouvant contenir des documents classiques aussi bien que des *pseudo-documents* et des *méta-segments*. Cette partie a été validée par la création de bons corpus bilingues français-chinois, l'un d'eux résultant de la toute première application à la traduction littéraire.

Une seconde partie, initialement motivée par un besoin industriel, a consisté à construire des systèmes de TA de type Moses, spécialisés à des sous-langages, en français↔chinois, et à étudier la façon de les améliorer dans le cadre d'un usage en continu avec possibilité de PE. Dans le cadre d'un projet interne sur le site du LIG et d'un projet (TABE-FC) en coopération avec l'université de Xiamen, on a pu démontrer l'intérêt de l'apprentissage incrémental en TA statistique, sous certaines conditions, grâce à une expérience qui s'est étalée sur toute la thèse.

La troisième partie est consacrée à des contributions et mises à disposition de supports informatiques et de ressources. Les principales se placent dans le cadre du projet COST MUMIA de l'EU et résultent de l'exploitation de la collection CLEF-IP 2011 de 1,5 M de brevets partiellement multilingues. De grosses mémoires de traductions en ont été extraites (17,5 M segments), 3 systèmes de TA en ont été tirés, et un site Web de support à la RI multilingue sur les brevets a été construit. On décrit aussi la réalisation en cours de JIANDAN-EVAL, une plate-forme de construction, déploiement et évaluation de systèmes de TA.

Abstract

This thesis, conducted as part of a CIFRE grant, and extending one of the aspects of the ANR project TRAQUIERO, first addresses the production, extension and improvement of multilingual corpora by machine translation (MT) and contributory post-editing (PE). Functional and technical improvements have been made to the SECTRA and IMAG, and progress has been made toward a generic definition of the structure of a multilingual, annotated and multi-media corpus that may contain usual documents as well as *pseudo-documents* (such as Web pages) and *meta-segments*. This part has been validated by the creation of good French-Chinese bilingual corpora, one of them resulting from the first application to literary translation.

A second part, initially motivated by an industrial need, has consisted in building MT systems of Moses type, specialized to sub-languages, for french↔chinese, and to study how to improve them in the context of a continuous use with the possibility of PE. As part of an internal project on the LIG website and of a project (TABE-FC) in cooperation with Xiamen University, it has been possible to demonstrate the value of incremental learning in statistical MT, under certain conditions, through an experiment that spread over the whole thesis.

The third part of the thesis is devoted to contributing and making available computer tools and resources. The main ones are related to the COST project MUMIA of the EU and result from the exploitation of the CLEF-IP 2011 collection of 1.5 million partially multilingual patents. Large translation memories have been extracted from it (17.5 million segments), 3 MT systems have been produced (de-fr, en-fr, fr-de), and a website of support for multilingual IR on patents has been constructed. One also describes the on-going implementation of JIANDAN-EVAL, a platform for building, deploying and evaluating MT systems.

Remerciements

Il me sera très difficile de remercier tout le monde car c'est grâce à l'aide de nombreuses personnes que j'ai pu mener cette thèse à son terme.

J'aimerais tout d'abord remercier mon directeur de thèse, monsieur Christian BOITET, pour toute son aide. Je suis ravi d'avoir travaillé en sa compagnie car outre son appui scientifique, il a toujours été là pour me soutenir et me conseiller au cours de l'élaboration de cette thèse.

Je remercie également ma co-directrice de thèse, madame Valérie BELLYNCK, pour la gentillesse et la patience qu'elle a manifestées à mon égard durant cette thèse, pour tous les conseils et les programmes qu'elle a bien voulu m'envoyer.

J'adresse aussi mes remerciements à monsieur François BROWN DE COLSTOUN. Il m'a tout d'abord permis d'intégrer L&M en me proposant un sujet très intéressant et m'a laissé la liberté de l'orienter au cours du déroulement de ma thèse. C'est également grâce à sa collaboration avec le laboratoire LIG-GETALP que j'ai eu la chance de travailler dans L&M, ce qui c'est avéré une expérience très enrichissante. Je remercie aussi l'ANRT qui a cofinancé cette thèse.

Monsieur Xiaodong SHI et monsieur Eric WEHRLI m'ont fait l'honneur d'être rapporteurs de ma thèse. Leurs remarques m'ont permis d'envisager mon travail sous un autre angle. Pour tout cela je les remercie.

Mes remerciements vont également à madame Catherine BERRUT, monsieur Jean-Pierre CHEVALLET et monsieur Hong-Thai NGUYEN pour avoir accepté de participer à ce jury de thèse.

Je remercie ma chère épouse, ZHANG Ying, pour son soutien quotidien indéfectible et son enthousiasme contagieux à l'égard de mes travaux comme de la vie en général. Notre couple a grandi en même temps que mon projet scientifique, le premier servant de socle solide à l'épanouissement du second.

Mes remerciements s'adressent enfin à ma mère qui m'a toujours épaulé dans ce projet. Je sais que mon absence a été longue et j'espère pouvoir un jour rattraper le retard accumulé. Mon père, merci de m'avoir encouragé pendant ma vie. Je te dois beaucoup. Merci aussi à Xile, tu remplis ma vie de joie.

谨以此论文，献给我天堂里的父亲。

中文摘要

此论文获得了 CIFRE 奖学金，并且作为 ANR Traouiero 项目的延伸，首先通过机器翻译 (traduction automatique) 和后编辑 (post-édition) 的方法，创建，扩展和完善多语种语料库。针对 SECTra 和 iMAG 软件产品做了功能和技术上的改进，并且在结构上推进了对多语种，多元注释和多媒体语料库的广义定义，其中包含了常规文档，伪文档和元段的定义。这部分已通过建立良好质量的法汉双语语料库验证，该语料库中的一部分语料是第一次把后编辑应用到文学翻译中获得的。

第二部分，由最初工业领域对机器翻译的需求出发，致力于构建摩西机器翻译系统，针对法语↔汉语专业领域子语言，研究通过使用连续的后编辑方法改进机器翻译质量的可能性。在针对 LIG 网站的内部项目和与厦门大学合作的项目 (TABE-FC) 中，我们通过论文中所阐述的实验，证明了在一定条件下，增量训练对统计机器翻译系统的价值。

论文的第三部分是专门讲述在计算机支持和语言资源方面的贡献。主要贡献是在欧盟项目 COST MUMIA 中，对 CLEF-IP 2011 语料库中 150 万部分平行多语种专利语料库的处理结果。包括大翻译记忆的提取 (约 1750 万句对)，三个机器翻译系统的训练，以及构建支持多语言专利信息查找的网站。我们还阐述了正在构建中的 JianDan-eval 系统，一个用于构建，部署，评估机器翻译系统的平台。

Table des matières

Résumé	2
Abstract	2
Remerciements.....	3
中文摘要	4
Table des matières.....	5
Table des figures.....	8
Table des tableaux.....	10
Glossaire et abréviations.....	12
Introduction	14
Partie A Production, extension et amélioration de corpus multilingues par TA et PE contributive.....	17
Résumé	17
Chapitre I Amélioration d'aspects fonctionnels et techniques de SECTra et du logiciel iMAG pour les passerelles d'accès multilingue	17
I.1 SITUATION ET ETAT DE L'ART AU DEBUT DE LA THESE	17
I.1.1 <i>Modélisation et exploitation de corpus de traductions : SECTra_w</i>	17
I.1.2 <i>Accès multilingue à des sites Web : le logiciel iMAG</i>	27
I.1.3 <i>Travaux comparables et idées directrices pour le futur</i>	34
I.2 AMELIORATIONS DE SECTRA_W DANS LE CADRE DU PROJET TRAQUIERO	38
I.2.1 <i>Extension de fonctions existantes</i>	38
I.2.2 <i>Aspects de génie logiciel</i>	40
I.2.3 <i>Travail de spécification</i>	42
I.3 AMELIORATIONS DES iMAG DANS LE CADRE DU PROJET TRAQUIERO	43
I.3.1 <i>Paramétrisation</i>	43
I.3.2 <i>Travail de spécification</i>	44
Chapitre II Travail sur des aspects plus conceptuels et définition de nouvelles fonctionnalités	46
II.1 MODELISATION DE CORPUS DE TRADUCTIONS VARIES.....	46
II.1.1 <i>Variété des corpus de traductions et esquisse d'une méthode de formalisation</i>	46
II.1.2 <i>Esquisse d'une méthode de modélisation des corpus de traductions</i>	52
II.1.3 <i>Validation au niveau des métadonnées</i>	54
II.2 CONCEPTION DE NOUVELLES FONCTIONNALITES ET DEVELOPPEMENTS EN COURS.....	63
II.2.1 <i>Nouvelles fonctionnalités</i>	63
II.2.2 <i>Programmabilité</i>	65
II.2.3 <i>Modularité</i>	66
Chapitre III Variété des iMAG et de leurs usages : de l'accès multilingue à la création de bons corpus bilingues et à la traduction littéraire contributive de qualité	68
III.1 LISTE (AVEC LES MT ASSOCIEES).....	68
III.2 COMMENTAIRES SUR LES UTILISATIONS ACTUELLES	70
III.2.1 <i>Accès à des sites Web d'organismes ou de sociétés</i>	70
III.2.2 <i>Aide à la traduction de documents : rapports, parties de thèse, manuels</i> ...	70
III.2.3 <i>Accès multilingue à des documents pédagogiques : MACAU</i>	71
III.2.4 <i>Évaluation</i>	72
III.3 UTILISATIONS PLUS NOVATRICES : PRODUCTION DE BONS CORPUS PARALLELES, ET POST-EDITION DE TEXTES LITTERAIRES POUR L'AUTO-APPRENTISSAGE OU POUR LA TRADUCTION CONTRIBUTIVE.....	73
III.3.1 <i>Production de « corpus parallèles » de qualité</i>	73
III.3.2 <i>« Voyage au centre de la terre » de Jules Verne</i>	75
III.3.3 <i>« The Book of Me » de Powers</i>	76
III.3.4 <i>« IITB : Monastery, Sanctuary, Laboratory » de Rohit Manchanda</i>	77
Partie B Construction de systèmes de TA spécialisés à des sous-langages en français ↔ chinois	79
Chapitre IV Revue des systèmes TA français ↔ chinois en contexte industriel	80
IV.1 DEMANDE DE GROSSES SOCIETES	80
IV.2 ÉTAT DE L'ART DE LA TA DU CHINOIS	81
IV.2.1 <i>Historique</i>	81
IV.2.2 <i>Expérimentations</i>	83
IV.3 CONSTRUCTION DE SYSTEMES DE TA POUR LE CHINOIS BASES SUR MOSES EN CONTEXTE INDUSTRIEL	87
IV.3.1 <i>Choix du sous-langage et des couples à traiter</i>	87
IV.3.2 <i>Recherche infructueuse de corpus parallèles adaptés</i>	88

IV.3.3	<i>Production de corpus par PE de résultats de Google</i>	90
IV.3.4	<i>Construction de systèmes français→chinois</i>	91
IV.3.5	<i>Évaluations et perspectives</i>	93
Chapitre V	Construction de systèmes de TA pour le chinois avec Moses en contexte de recherche : le projet TABE-FC	96
V.1	BUTS DU PROJET TABE-FC.....	96
V.1.1	<i>Buts théoriques</i>	96
V.1.2	<i>Buts pratiques</i>	97
V.1.3	<i>Définition du projet</i>	97
V.2	CONSTITUTION DES CORPUS D'APPRENTISSAGE.....	98
V.2.1	<i>Recherche de sites et collecte de pages Web monolingues et bilingues</i>	98
V.2.2	<i>Nettoyage et filtrage</i>	98
V.2.3	<i>TA par GT, puis PE (production d'un corpus parallèle)</i>	100
V.3	CONSTRUCTION DE SYSTEMES DE TA.....	102
V.3.1	<i>Construction de systèmes Moses "ligne de base"</i>	102
V.3.2	<i>Avancement de l'expérimentation</i>	103
V.3.3	<i>Résultats provisoires</i>	103
Chapitre VI	Démonstration de l'intérêt de l'apprentissage incrémental en TA statistique	105
VI.1	CONTEXTE.....	105
VI.1.1	<i>Motivations</i>	105
VI.1.2	<i>Expérience sur le site du LIG</i>	106
VI.2	EXPERIMENTATION.....	107
VI.2.1	<i>Phase 1 (2-6/2013)</i>	107
VI.2.2	<i>Phase 2 (7-9/2013)</i>	110
VI.2.3	<i>Phase 3 (9-12/14 et 7-11/15)</i>	113
VI.3	ANALYSE DES RESULTATS.....	116
Partie C	Contribution d'outils et de ressources	117
Chapitre VII	Construction de systèmes de TA et support à la RI multilingue pour MUMIA	118
VII.1	CONTEXTE ET MOTIVATIONS.....	118
VII.1.1	<i>Description du projet MUMIA et du WG2</i>	118
VII.1.2	<i>PerFedPat et Khresmoi</i>	119
VII.1.3	<i>Objectif poursuivi</i>	120
VII.2	CONSTRUCTION DE MT ET DE STA A PARTIR DES CORPUS CLEF-IP 2011.....	120
VII.2.1	<i>Description du corpus CLEF-IP</i>	120
VII.2.2	<i>Extraction de MT à partir de CLEF-IP 2011</i>	122
VII.2.3	<i>Construction des systèmes de TA</i>	125
VII.3	EXPERIMENTATION ET ELARGISSEMENT A D'AUTRES LANGUES.....	125
VII.3.1	<i>Reconstruction de trois sites Web de brevets monolingues</i>	125
VII.3.2	<i>Accès multilingue en utilisant les systèmes de TA créés</i>	127
VII.3.3	<i>Accès multilingue utilisant d'autres systèmes et pour d'autres langues</i>	127
Chapitre VIII	Mise à disposition de ressources	130
VIII.1	CONTRIBUTION DE RESSOURCES STATIQUES SOUS FORME DE MT.....	130
VIII.1.1	<i>Formats choisis</i>	130
VIII.1.2	<i>Méthode de création</i>	132
VIII.1.3	<i>Résultats</i>	132
VIII.2	CONTRIBUTION SOUS FORME DE SYSTEMES DE TA.....	135
VIII.2.1	<i>Systèmes de TA téléchargeables</i>	135
VIII.2.2	<i>Systèmes de TA utilisables comme des services Web</i>	135
VIII.3	PASSERELLES iMAG VERS DES SITES WEB STATIQUES OU DYNAMIQUES.....	135
VIII.3.1	<i>Passerelles iMAG pour des sites statiques</i>	135
VIII.3.2	<i>Passerelles iMAG pour des sites dynamiques</i>	136
VIII.3.3	<i>Structure d'une contribution « dynamique » par iMAG</i>	137
VIII.3.4	<i>Remarques sur la création de certains des sites « contribués »</i>	138
Chapitre IX	Vers une plate-forme de construction, déploiement et évaluation de systèmes de TA: JianDan-eval	140
IX.1	CAHIER DES CHARGES, ARCHITECTURE, ET SPECIFICATIONS EXTERNES.....	140
IX.1.1	<i>Cahier des charges</i>	140
IX.1.2	<i>Spécifications externes</i>	141
IX.1.3	<i>Architecture logicielle</i>	145

IX.2	IMPLEMENTATION	146
IX.2.1	<i>Outils de base utilisés</i>	147
IX.2.2	<i>Composants utilisés</i>	147
IX.2.3	<i>Composants développés</i>	147
IX.3	EXEMPLE : CREATION D'UN SYSTEME DE TA.....	148
Conclusions et perspectives		150
Bibliographie		152
Table des définitions.....		157
Annexes	159	
ANNEXE 1 :	CORPUS DE LA CAMPAGNE D'EVALUATION DE TA DU PROJET TRANSAT	159
ANNEXE 2 :	PROTOCOLE D'EVALUATION POUR LE PROJET TRANSAT	161
ANNEXE 3 :	UN EXEMPLE DU CORPUS B@BEL	163
ANNEXE 4 :	CORPUS EOLSS	164
ANNEXE 5 :	UN EXEMPLE DE GRAPHE UNL AVEC CORRECTION	165
ANNEXE 6 :	DOCUMENT DE BREVET DU CORPUS CLEF-IP 2011 (EP-0000007-B2.XML)	167
ANNEXE 7 :	STRUCTURE DES DONNEES DE MINIDICTIONAIRES.....	170
ANNEXE 8 :	50 SEGMENTS EN « <i>VUE SECTra/POST-EDITION</i> »	171
ANNEXE 9 :	EXEMPLE DU CORPUS PARALLELE FRANÇAIS-CHINOIS CREE POUR L&M DANS LE DOMAINE DE L'ENERGIE.....	177
ANNEXE 10 :	EXEMPLE D'EXTRACTION DE BISEGMENTS A PARTIR DU CORPUS MULTIUN	204
ANNEXE 11 :	SCRIPT DE FILTRAGE DE CORPUS	206
ANNEXE 12 :	SOURCE DU PROGRAMME POUR CALCULER D_{MIX}	209
ANNEXE 13 :	100 BISEGMENTS ANGLAIS-FRANÇAIS EXTRAITS DE CLEF-IP 2011	212
ANNEXE 14 :	20 BISEGMENTS ANGLAIS-FRANÇAIS EN FORMAT TMX	220

Table des figures

Figure 1 : Architecture générale d'une iMAG pour un site élu.....	30
Figure 2 : Capture d'écran de l'iMAG LIG-LAB en chinois.....	31
Figure 3 : Architecture par agents SECTra_w, iMAG, PIVAX (Nguyen, 2009)	32
Figure 4 : Interface de « Translate Corpus »	39
Figure 5 : Interface de TRADOH.....	40
Figure 6 : Options de la fonction "export"	41
Figure 7 : Interface de sélection paramétrable dans SECTra_w	41
Figure 8 : Structure logique d'une base de données de corpus multilingues	43
Figure 9 : Fichier HTML et fichier compagnon .unl.....	50
Figure 10 : Document 2 traduit de l'anglais vers le français (<i>GROUND AND SOIL WATER CHARACTERISTICS</i>).....	50
Figure 11 : Exemple de structure et de description d'un dialogue du corpus ERIM	51
Figure 12 : Exemple du fichier <i>french.wpl</i> et <i>vietnamese.wpl</i>	51
Figure 13 : Capture d'écran de panneau de dictionnaires ajouté à SECTra_w.....	64
Figure 14 : Interface de SECTra_w intégrant les boutons « Delete », « Clean » et « Get ».....	64
Figure 15 : Traduction des segments sélectionnés et ajout à la MT.....	65
Figure 16 : Exemple de l'API « Call Tradoh »	65
Figure 17 : Post-édition d'un document français accédé en anglais (résumé de la thèse de Lingxiao WANG)	71
Figure 18 : Extraction of a "good" TM from a TM produced by "natural" post-edition.....	73
Figure 19 : Export of a « good » part of a TM	74
Figure 20 : Capture d'écran de iMAG français→chinois pour « Voyage au centre de la terre ».....	75
Figure 21 : Exemple de post-édition d'un chapitre de « Monastery, Sanctuary, Laboratory: 50 Years of IIT-Bombay » de Rohit Manchanda	78
Figure 22 : Architecture à 3 niveaux et 7 « missions » du projet TABE-FC (Chen, Wang et al., 2014)	97
Figure 23 : Exemple de page Web économique parallèle	98
Figure 24 : Exemple d'une page Web du site de "Bourse de Hong Kong" en format html.....	99
Figure 25 : Exemple de segments chinois-anglais extraits à partir de pages Web	100
Figure 26 : Capture d'écran de l'iMAG "Bourse de Paris" en chinois.....	101
Figure 27 : Comparaison de la traduction de GT et de la post-édition humaine.....	102
Figure 28 : Site du LIG vu en chinois à travers une iMAG	107
Figure 29 : Diminution de temps moyen de PE (par page standard) avec AI dans la phase 1 de l'expérience	109
Figure 30 : Capture d'écran de l'iMAG « Corpus par jour »	111
Figure 31 : Capture d'écran de Chamilo affichant le lien AXiMAG.....	112
Figure 32 : Diminution du temps moyen de PE (par page) avec AI dans la phase 3 de l'expérience ..	115
Figure 33 : Architecture de PerFedPat	119
Figure 34 : KHRESMOI.....	120
Figure 35 : Exemple de fichier XML Dans CLEF-IP	122
Figure 36 : Exemple de champ <claims> contenant 6 sous-champs <claim> dans <i>EP-0260000-B1.xml</i>	122
Figure 37 : Exemple d'un champ <invention-title> avec 3 attributs de langue différents et les contenus correspondants en 3 langues différentes.....	123
Figure 38 : Un champ <patent-document> avec attribut lang = "EN"	123
Figure 39 : Exemple de fichier XML monolingue	125
Figure 40 : Exemple de revendication dans le fichier EP0203923B1.xml.....	126
Figure 41 : Exemple de fichier HTML décoré	127
Figure 44 : PE en mode avancé, avec pseudo-trace montrant les différences entre les sorties de TA, la post-édition (utilisée comme référence), et la MT.	128
Figure 45 : Retraduction des segments du français vers le chinois pour DOC6 avec le système de TA français→chinois MosesLIG	129
Figure 46 : Exemple des données en format TXT (MT CLEF-IP anglais-français)	131

Figure 47 : Exemple des données en format TMX (MT CLEF-IP anglais-français).....	131
Figure 48 : Extraction d'une "bonne" MT de la MT produite par post-édition "naturelle"	132
Figure 49 : Segments post-édités pour la ressource énergie.....	134
Figure 50 : Exemple de contribution au format HTML (Chapitre 1 : Voyage au centre de la Terre) .	137
Figure 51 : Capture d'écran du site Web monolingue de CLEF-IP	138
Figure 52 : Capture d'écran de l'iMAG dédiée CLEF-IP	139
Figure 53 : Architecture initiale de gestion de travaux	146
Figure 54 : Tpetotal par page standard de différents systèmes de TA	149

Table des tableaux

Tableau 1 : Sites Web élus des iMAG dédiées disponibles en 2010.....	18
Tableau 2 : Données statistiques sur les segments post-édités dans SECTra_w depuis 2010	19
Tableau 3 : Liste des iMAG à MT dédiée construites depuis 2010.....	29
Tableau 4 : Exemples de sites Web de partage de corpus parallèles.....	34
Tableau 5 : Comparaison de l'organisations logiques, physiques, et interne de quelque corpus.....	52
Tableau 6 : Métadonnées du corpus BTEC (les segments extraits)	54
Tableau 7 : Métadonnées des données d'évaluation à la TRANSAT	55
Tableau 8 : Métadonnées du corpus UNESCO-B@bel.....	56
Tableau 9 : Métadonnées du corpus EOLSS au niveau de la macrostructure.....	57
Tableau 10 : Métadonnées d'un fichier HTML au niveau de la microstructure	58
Tableau 11 : Métadonnées d'un fichier UNL au niveau de la microstructure	58
Tableau 12 : Métadonnées d'un corpus EOLSS au niveau de la mésostructure	59
Tableau 13 : Métadonnées d'un corpus ERIM au niveau de la macrostructure.....	59
Tableau 14 : Métadonnées de la séance dans le corpus ERIM a au niveau de <i>la microstructure</i>	60
Tableau 15 : Métadonnées d'un corpus ERIM au niveau de la mésostructure	60
Tableau 16 : Métadonnées du corpus CLEF-IP 2011	61
Tableau 17 : Métadonnées d'un document de brevet.....	62
Tableau 18 : 10 paramètres de l'API de CREATDICO	63
Tableau 19 : Exemple de fichier de configuration de CREATDICO.....	63
Tableau 20 : Exemple d'un lien pour l'utilisation de l'API de CREATDICO.....	63
Tableau 21 : iMAG pour les sites Web de laboratoires et d'universités.....	68
Tableau 22 : iMAG pour les sites Web d'organismes et de sociétés	69
Tableau 23 : iMAG pour des projets et des expérimentations	69
Tableau 24 : Nombre de langue du projet MACAU (06/2013).....	72
Tableau 25 : Statistiques de documents dans MACAU (06/2013).....	72
Tableau 26 : Statistique sur 21 chapitres de « Voyage au centre de la terre »	75
Tableau 27 : Corpus source, cible traduite et cible corrigée	76
Tableau 28 : Statistique sur les données.....	83
Tableau 29 : Formule d'évaluation de l'automatisme et de la qualité d'un système de TA.....	84
Tableau 30 : Exemple de traduction de GT.....	84
Tableau 31 : Paramètres de configuration de Joshua	85
Tableau 32 : Comparaison d'exemples de traductions obtenues par TA et d'une référence.....	86
Tableau 33 : Exemple de résultat d'évaluation	87
Tableau 34 : Corpus collectés en cherchant des corpus pour le français→chinois.....	89
Tableau 35 : Exemples de bisegments français→chinois parmi les 9000 collectés ou produits	90
Tableau 36 : Comparaison des temps d'entraînement de Moses	91
Tableau 37 : Configuration de la machine	92
Tableau 38 : Scores BLEU pour différentes tailles du corpus d'entraînement.....	92
Tableau 39 : Statistiques sur le corpus MultiUN.....	92
Tableau 40 : Exemple de données de test.....	93
Tableau 41 : Statistiques des données de test.....	94
Tableau 42 : Score BLEU et exemples de sorties de systèmes de TA	95
Tableau 43 : Statistiques des pages Web collectées.....	99
Tableau 44 : Exemple de conversion des caractères chinois du traditionnel vers le simplifié	99
Tableau 45 : Statistiques sur la ressource économique et boursière	102
Tableau 46 : Statistiques sur les données d'entraînement de la phase 1	107
Tableau 47 : Évaluation du temps de post-édition (2-6/2013)	109
Tableau 48 : Évaluations basées sur des références (BLEU, NIST, TER).....	110
Tableau 49 : Statistiques de post-édition sur 21 articles français 4/7-13/9/2013	111
Tableau 50 : Statistiques de post-édition sur les supports de cours	112
Tableau 51 : Résultat de l'expérimentation (en français-chinois).....	113
Tableau 52 : Nombre de segments dans chaque MT.....	113
Tableau 53 : Statistiques sur les données pour l'AI (phase 3 de l'expérience).....	114

Tableau 54 : Évaluation du temps de post-édition (9-12/2014)	115
Tableau 55 : Données de test et scores BLEU	116
Tableau 56 : Nombre de segments extraits comme source et cible après l'alignement de segments dans les champs <i><title></i> et <i><claims></i>	124
Tableau 57 : Scores BLEU des systèmes de TA tirés de CLEF-IP	125
Tableau 58 : Segments post-édités dans SECTra_w à partir de 3 langues source	133
Tableau 59 : Segments parallèles obtenus à partir des MT (mêmes remarques)	133
Tableau 61 : Statistique des données pour la ressource énergie	134
Tableau 62 : Systèmes de TA téléchargeables	135
Tableau 63 : Systèmes de TA utilisables comme des services Web	135
Tableau 64 : Passerelles iMAG pour des sites statiques	136
Tableau 65 : iMAG pour des sites dynamiques.....	136
Tableau 66 : Résumé des données	148

Glossaire et abréviations

AI	Apprentissage incrémentale
API	Application Programming Interface
ARIANE-G5	Environnement de développement et d'exécution de systèmes (linguiciels) de traduction automatique
CLEF-IP	Conference and Labs of the Evaluation Forum, formerly known as Cross-Language Evaluation Forum - Intellectual Property
CLIPS	Communication Langagière et Interaction Personne Système
CNRS	Centre national de la recherche scientifique
DC	Dublin core
DSR	Digital Silk Road
EDF	Électricité de France
EOLSS	Encyclopedia of Life Support Systems
FC	Français-Chinois
GETA	Groupe d'Etude pour la Traduction Automatique
GETALP	Groupe d'Étude pour la Traduction/le Traitement Automatique des Langues et de la Parole
GI	Génie informatique
GT	Google Translate
HQ	Haute Qualité
IMAG	Passerelle interactive d'accès multilingue (interactive Multilingual Access Gateway)
ISCC	Institut des sciences de la communication
L&M	SAS Lingua et Machina
LIG	Laboratoire d'Informatique de Grenoble
LSPL	Langage Spécialisé pour la Programmation Linguistique
MACAU	Multilingual Access & Contributive Appropriation for Universities
MT	Mémoire de Traductions
MUMIA	Multilingual and multifaceted interactive information access
OMNIA	Outils et Méthodes Numériques pour l'Interrogation et l'Analyse des textes médiolatins
ONU	Organisation des Nations Unies
PAHO	Pan American Health Organization
PCT	Patent Cooperation Treaty
PE	Post-Édition
PIVAX	Base lexicale à pivot par acceptations (monolingues et interlingues)
QCM	Question à choix multiples
RI	Recherche d'Information
SAAS	Software as a Service
SECTRA_W	Système d'Exploitation de Corpus de Traductions sur le Web
SECTRA/EVAL	Première version de SECTra, développée pour le projet TRANSAT d'Orange Labs
SECTRA/TRAD	Deuxième version de SECTra, développée pour le projet EOLSS/UNL++
SECTRA/WEB	Troisième version de SECTra, développée pour le projet iMAG
SEGDOC	Segmentation de documents XML
TA	Traduction Automatique
TH	Traduction Humaine
TMX	Translation Memory eXchange

TRADOH	Un outil permet d'obtenir une traduction dans sa langue, par mise en œuvre automatique d'un ou plusieurs systèmes de TA disponibles en local ou à distance, avec composition éventuelle.
TRANSAT	Projet de TA de parole d'Orange 2004~2007
TXT	Texte brut
UNESCO/B@BEL	Partie du site Web de l'Unesco consacrée à la communication multilingue
UNL	Universal Networking Language (projet lancé par l'UNU (fin 1995) et langage "anglosémantique" d'hypergraphes associables aux énoncés en langue naturelle
XLIFF	XML Localisation Interchange File Format
XRCE	Xerox Research Center Europe

Introduction

Cette thèse a été effectuée dans l'équipe GETALP du LIG, et dans le cadre d'une bourse CIFRE avec Lingua et Machina, une jeune société qui vise à "prendre en charge la communication multilingue de l'entreprise". Le sujet initialement défini était centré sur l'amélioration de plusieurs aspects de génie logiciel du logiciel SECTRA_W/IMAG, réalisé par Cong Phap HUYNH dans le cadre de sa thèse (Huynh, 2010). Le point principal concernait la transformation de la partie SECTRA_W (*Système d'Exploitation de Corpus de Traductions* sur le Web) en un système programmable et extensible. On visait à pouvoir l'utiliser comme un "serveur corporal" gérant des suites de test ainsi que des corpus de développement pour le compte de systèmes de traduction automatique (TA) munis d'un environnement de développement complet, comme ARIANE-G5 et son successeur, en cours d'implémentation, ARIANE-Y. Un autre objectif était de pouvoir non seulement exploiter des corpus parallèles existants, pour les évaluer et/ou les améliorer par post-édition collaborative en ligne, mais aussi de pouvoir les étendre à de nouvelles langues par appel à des serveurs de TA, suivi de post-édition. Un dernier thème était la recherche d'une méthode de spécification formelle implémentable des « vrais » corpus multilingues, c'est-à-dire pas seulement des listes de "segments" multilingues, comme le BTEC (Boitet et al., 2007), qui ne sont en fait que de grandes « mémoires de traduction » (MT), même si on les appelle « corpus parallèles ». On souhaitait s'attaquer à la complexité des « vrais corpus », et passer à l'échelle, de façon, par exemple, à pouvoir traiter des corpus fortement structurés et de très grande taille comme ceux des brevets. Dans un corpus complexe, un document est formé d'un document maître (en XML, par exemple) ou d'une hiérarchie de tels documents, accompagné d'une collection de fichiers « satellites » (images, vidéos...) et éventuellement d'annotations contenues dans des fichiers « compagnons ».

On visait aussi à résoudre un certain nombre de problèmes liés à l'interaction entre SECTRA_W et le logiciel iMAG (*interactive Multilingual Access Gateways*) qui utilise SECTRA_W comme un « dorsal » et permet d'accéder à des sites Web « élus » dans un grand nombre de langues, avec possibilité d'améliorer les « prétraductions » produites par des serveurs de TA en les corrigeant (« post-éditant ») directement sur la page Web, ou dans l'interface de SECTRA_W dédiée à la post-édition.

Ces objectifs ont évolué à cause des besoins de l'entreprise, qui désirait d'abord construire des systèmes de TA « maison » français↔chinois en utilisant la boîte à outils MOSES (Koehn et al., 2007). Pour cela, il faut disposer de grands corpus parallèles de bonne qualité, dans le bon sens, et représentatifs des sous-langages des clients potentiels, en l'occurrence EDF, RENAULT, etc. Cela a d'abord mené à l'étude, l'expérimentation et l'évaluation d'aligneurs divers, de segmenteurs du chinois, et de divers systèmes de TA existants. Dans un deuxième temps, le travail s'est plus orienté vers la TA proprement dite. Il s'est d'abord agi de construire un environnement de préparation et d'exploitation de systèmes MOSES, intégré aux outils de L&M (LIBELLEX, MYRIAM) ou utilisés par L&M (METRICC, XELDA...). L'obstacle majeur à surmonter était l'absence de corpus parallèles français-chinois. Nous avons alors construit un corpus de 9000 segments, d'abord par post-édition de résultats de GOOGLE TRANSLATE (GT), puis par post-édition de résultats d'une première version d'un système MOSES-L&M-FC. À peu près à la même période, L&M nous a demandé d'étudier la nouvelle possibilité offerte par MOSES de faire de l'apprentissage incrémental (AI). Nous l'avons fait, ainsi que quelques essais préliminaires, puis L&M m'a dirigé sur autre chose, jugeant l'approche peu prometteuse. Il est vrai que les améliorations constatées étaient faibles. Pourtant, les gains de temps étaient considérables (environ 1h pour l'AI sur quelques dizaines ou centaines de nouveaux segments au lieu de 20h pour un réapprentissage complet). Notre intuition était que, dans le cas de sous-langages, on devait arriver, en faisant quelques dizaines d'itérations d'AI, à produire des prétraductions

meilleures que celles de GT, BING, SYSTRAN ou NIUTRANS, au moins en ce qui concerne la qualité d'usage pour la tâche de post-édition, et peut-être aussi pour la tâche de compréhension. Encouragé par mes directeurs de thèse, j'ai alors orienté ma recherche dans cette direction, pendant plusieurs mois. Au terme d'une expérience sur le sous-langage du site Web du LIG¹, nous avons pu publier (à COLING 2012) des résultats encourageants : notre courbe de "temps de post-édition" descendait assez régulièrement, et, au bout d'une vingtaine d'itérations avec un réapprentissage total au milieu, n'était plus qu'un petit peu au-dessus de la "ligne de base" correspondant à la PE des résultats de GT (environ 10 mn/page). Depuis, j'ai préparé et mené une troisième expérience, en améliorant l'automatisation du processus et des mesures associées, et j'ai pu démontrer que, au moins dans le cas d'un sous-langage comme celui du site du LIG, et du français-chinois, la combinaison de l'apprentissage incrémental par périodes, et spécialisé à un sous-langage, pouvait donner des résultats nettement meilleurs que ceux des systèmes de TA généralistes.

Ma recherche s'est ensuite trouvée orientée vers le passage à l'échelle, la TA de brevets, et la construction de grandes ressources de bonne qualité, dans le cadre du projet COST "MUMIA" de l'UE. Mon directeur de thèse était en effet VP de MUMIA et en charge du WG2. Dans le WG2 auquel j'ai participé, il s'agissait d'étudier et de prototyper des "infrastructures" matérielles et logicielles pour la recherche d'information dans un cadre multilingue, multimodal et multi-facette. J'ai été amené à traiter la collection CLEF-2012 (la même que CLEF-2011), constituée à partir de 1,5 millions de brevets partiellement traduits par des professionnels. Une bonne proportion des "segments", initialement rédigés en français, allemand ou anglais, a été traduite dans une deux des autres langues. Il ne s'agit pas de collections parallèles : il y a un fichier par brevet, en XML, où chaque segment contient sa version originale et éventuellement une ou des versions dans d'autres langues. À partir de cette collection, j'ai détecté la langue source de chaque segment, et construit 3 mémoires de traductions (très bonnes par construction), une pour chaque langue source. Je les ai aussi utilisées comme base pour l'apprentissage de 3 systèmes Moses (allemand→français, français→allemand, français→anglais). Enfin, avec l'aide d'un étudiant de M1 en TER (Huanan SUN), j'ai construit 3 collections de brevets monolingues, chacun étant seulement dans sa langue source. Grâce à 3 iMAG, il est possible d'y accéder dans diverses langues. Les traductions sont évidemment très bonnes pour les langues initiales (en tout cas, sur les parties du corpus réservées pour les tests), mais il est aussi possible d'y accéder dans d'autres langues, par exemple en chinois, d'améliorer les résultats de TA par PE, et de recycler les "bonnes traductions" pour construire un système spécialisé s'améliorant au fur et à mesure de l'usage.

Durant la dernière partie de ma thèse, je suis revenu au thème de la TA français-chinois, en participant à deux projets, MACAU-OFI et TABE-FC. MACAU-OFI est un projet défini par R. Kalitvianski et Ch. Boitet en 2012, visant à mettre à disposition des étudiants étrangers des supports de cours dans leur langue, en utilisant une passerelle iMAG dédiée, et en demandant aux étudiants eux-mêmes de « post-éditer ». Durant l'été 2013, j'ai ainsi participé à l'encadrement de deux stages d'été d'étudiants chinois (en master informatique à l'UJF), qui ont post-édité environ 520 pages standard (130K mots) dans le domaine des outils formels pour l'informatique. En 2013-2014, j'ai aussi participé de façon très active à la définition et au début de la réalisation du projet TABE-FC monté avec l'université de Xiamen, dans le cadre d'une année sabbatique passée à Grenoble par le Dr Y. CHEN. Il s'agit de construire des systèmes de TA permettant à des Chinois d'avoir un accès en chinois de bonne qualité (et surtout bien plus fiable et fidèle que les systèmes généralistes) aux "brèves"² des bourses

¹ <http://liglab.fr>

² « *flash reports* » en anglais

francophones³, et plus généralement aux sites Web économiques en français, et inversement pour des Français désireux d'intervenir sur les bourses de Shanghai, Shenzhen et Hong Kong. Cela m'a fait revenir aux thèmes plus liés au génie logiciel.

Au total, mon apport se situe dans quatre domaines principaux : (1) le génie logiciel des systèmes non seulement d'exploitation, mais aussi maintenant de création et de gestion de « vrais » corpus multilingues, (2) la TA, avec des contributions portant sur l'apprentissage incrémental, la TA français-chinois, ainsi que les environnements de construction et de déploiement de systèmes de TA de type MOSES ou similaire, (3) la mise à disposition de ressources (mémoires de traductions, systèmes de TA associés), et (4) la spécification et l'implémentation en cours d'une infrastructure pour l'évaluation, la plate-forme JIANDAN-EVAL, qui permettra les évaluations classiques, ainsi que l'évaluation comparative et « en usage » de systèmes de TA de toutes les architectures existantes.

La première partie de ce mémoire concerne la production, l'extension et l'amélioration de corpus multilingues par traduction automatique (TA) et post-édition contributive (PE). Des améliorations fonctionnelles et techniques ont été apportées aux logiciels SECTRA_W et IMAG produits lors des thèses de C.P. HUYNH et H.T. NGUYEN. Nous avons progressé vers une définition générique de la structure d'un corpus multilingue, multi-annoté et multimédia, pouvant contenir des documents classiques aussi bien que des *pseudo-documents* (comme des pages Web) et des *méta-segments*. Cette partie a été validée par la création de bons corpus bilingues français-chinois, l'un d'eux résultant de la toute première application à la traduction littéraire (un roman de Jules Verne), projet personnel mené pour progresser en français.

La seconde partie est centrée sur nos travaux en TA proprement dite. Initialement motivée par un besoin industriel, cette partie de notre recherche a consisté à étudier comment construire des systèmes de TA de type Moses, spécialisés à des sous-langages, en français↔chinois, et à étudier la façon de les améliorer dans le cadre d'un usage en continu avec possibilité de post-édition (PE) contributive en ligne. Dans le cadre d'un projet interne sur le site du LIG et d'un projet (TABE-FC⁴) en coopération avec l'université de Xiamen, nous avons pu démontrer l'intérêt de l'apprentissage incrémental en TA statistique, sous certaines conditions, grâce à une expérience qui s'est étalée sur toute la thèse.

Dans la troisième partie de ce mémoire, nous présentons nos contributions en termes de mise à disposition de supports informatiques et de ressources. Les principales se placent dans le cadre du projet COST MUMIA de l'EU et résultent de l'exploitation de la collection CLEF-2011 de 1,5 M brevets partiellement multilingues. De grosses mémoires de traductions en ont été extraites (17,5 M segments), trois systèmes de TA en ont été tirés (allemand→français, anglais→français, français→allemand), et un site Web de support à la RI multilingue sur les brevets a été construit. Avant de conclure, nous terminons en décrivant aussi la spécification et la réalisation en cours de JIANDAN-EVAL, une plate-forme de construction, déploiement et évaluation de systèmes de TA.

³ Un tel système, ALTFLASH (Uchino et al., 2001), a été déployé pour le Nikkei à partir de 2001.

⁴ TA pour les sites boursiers et économiques appliquée au français-chinois. En anglais: MTSE-FC.

Partie A Production, extension et amélioration de corpus multilingues par TA et PE contributive

Résumé

La partie A présente l'amélioration d'aspects fonctionnels et techniques de SECTRA_W et du logiciel iMAG pour les passerelles d'accès multilingue. Cette partie comporte aussi des aspects plus conceptuels, et la définition de nouvelles fonctionnalités. Nous montrons enfin la variété des iMAG et de leurs usages, de l'accès multilingue à la création de bons corpus bilingues et à la traduction littéraire contributive de qualité.

Chapitre I Amélioration d'aspects fonctionnels et techniques de SECTra et du logiciel iMAG pour les passerelles d'accès multilingue

Ce chapitre présente la situation et l'état de l'art au début de la thèse, puis les améliorations de SECTRA_W/iMAG étudiées et réalisées dans le cadre du projet TRAQUIERO.

I.1 Situation et état de l'art au début de la thèse

I.1.1 Modélisation et exploitation de corpus de traductions : SECTra_w

I.1.1.1 Présentation générale

I.1.1.1.1 Motivations et bref historique

En 2005-2007, notre équipe, alors le GETA du CLIPS, a participé au projet TRANSAT dans le cadre d'un contrat avec FRANCE TELECOM R&D. Dans ce projet, nous avons utilisé la toute première version de SECTra, SECTRA/EVAL, pour organiser une campagne d'évaluation de l'utilité potentielle des systèmes de traduction automatique commerciaux dans le domaine de la traduction de la parole. Le domaine privilégié dans cette étude était l'assistance à un touriste en situation difficile. Une autre partie de l'étude portait sur le domaine de la restauration.

De février 2008 à octobre 2008, la deuxième version de SECTra, SECTRA/TRAD, a été développée puis utilisée pour le projet EOLSS/UNESCOL réalisé dans le cadre d'un contrat entre l'Association Champollion et la fondation UNDL⁵. Dans ce projet, SECTRA a assuré le support et la gestion de la traduction de bonne qualité de 25 articles de l'encyclopédie EOLSS (*Encyclopedia of Life Support Systems*), soit environ 220K mots (880 pages standard), ou 13676 segments. En particulier, SECTRA_W a fourni pour ce projet un environnement collaboratif en ligne pour la post-édition humaine appliquée aux résultats de systèmes de TA et de déconvertisseurs UNL.

Depuis 2009, la troisième version, SECTRA_W, SECTRA/WEB, développée fin 2008, est utilisée comme "dorsal" pour l'accès multilingue de bonne qualité à des sites Web "élus". La conception et le développement de SECTra ont fait l'objet de la thèse de Cong Phap HUYNH (Huynh, 2010). À ce moment, mi-2010, nous avons alors construit des passerelles "iMAG" pour une trentaine de sites Web, dont celui du LIG. Chaque iMAG que nous avons construite

⁵ The UNDL FOUNDATION is a private Swiss law Foundation with head office in Geneva, Switzerland, legally representing the United Nations Organization in protecting its property rights pertaining to the UNL language and system, and legally representing the EOLSS Publishers and the UNESCO Joint Committee in translating the EOLSS with the use of the UNL technology.

apparaît comme un Wiki permettant l'accès multilingue à un site Web élu et la contribution à l'amélioration des traductions des segments (phrases et titres) de ses pages.

Essentiellement, une iMAG fournit une interface interactive permettant aux utilisateurs de voir et de post-éditer le site Web élu en plusieurs langues. En arrière-plan, tous les processus de gestion des segments "multilingualisés" du site Web sont réalisés par SECTRA_W, par SEGDOC (système de segmentation), et par TRADOH (intergiciel d'appel à des systèmes de TA).

1.1.1.1.2 Apports de la thèse de C.P. HUYNH

La thèse de C.P. HUYNH a apporté des réponses théoriques et pratiques à trois grands défis. Le premier défi consistait à offrir un support informatique unifié à l'évaluation des systèmes de TA. Le deuxième défi concernait le support contributif et collaboratif au travail humain sur des corpus variés en contexte multilingue. Le troisième défi était la construction d'un support informatique à l'exploitation de corpus de traductions dans des applications novatrices comme l'accès multilingue à des sites Web et la recherche d'information en contexte multilingue et multimédia (OMNIA). Plusieurs notions nouvelles ont été précisées (comme *segment multilingualisé et contextualisé*, *pseudo-document*, *métadocument*, etc.), et plusieurs principes généraux (*proactivité*, *délégation*, etc.) ont été introduits. C.P. HUYNH⁶ a dégagé six problèmes associés à chacun de ces trois défis, à dominante conceptuelle (par exemple, définition étendue d'un « contexte » de segment), algorithmique (par exemple, programmabilité du traitement des corpus), et programmatoire (par exemple, traitement de masses de données).

1.1.1.1.3 Évolutions depuis début 2010

Jusqu'au milieu de l'année 2010, 30 iMAG avaient été définies. Chaque iMAG est associée à une mémoire de traductions (MT) gérée par SECTRA_W. Une iMAG de démonstration partage en général une MT avec d'autres iMAG. Les MT partagées les plus utilisées sont DEMO, DEMO1, DEMO2. Voici quelques données présentées dans la thèse de Phap (Tableau 1).

Tableau 1 : Sites Web élus des iMAG dédiées disponibles en 2010

LIG laboratory	Digital Silk Road
Danang city	Da Nang University of Technology
TOL	ISCC
Unesco/B@bel	Systran
La Métro	Forum Lyon
Mica	Campus France
GETALP	Floralis
TechniLang	Ordinaide
Homerica	aikicorenc
MT 25 Years On	Essilor
Michelin	Winsoft
LeMonde.fr	UNDL-foundation
GETALP presentation	UNESCO_Babel
XD-consulting	ARDI Rhône

Entre 2010 et 2012, avant le début de cette thèse, nous avons participé (lors d'un stage en alternance de M2P-GI puis d'un stage prédoctoral effectué dans le cadre du projet ANR TRAQUIERO) à l'amélioration SECTRA_W, qui en est à la version 2. Grâce au configurateur

⁶ C'est son nom d'usage dans l'équipe.

d'IMAG réalisé par H. T. NGUYEN lors de son post-doctorat (projet IMAG/LAMETRO en 2010 pour l'Expo de Shanghai en octobre 2010), le GETALP a créé 80 IMAG de démonstration, partageant 3 mémoires de traductions principales. Il y a 8 langues source, et chaque site Web peut être accédé en plus de 10 langues cible (en fait, dans toutes les langues traitées par GT, système qui en traite le plus) et 5 autres MT.

Début 2015, environ 45% des segments visités via une IMAG (plus de 370 000), et donc « pré-traduits » automatiquement, avaient été post-édités par des contributeurs. La post-édition a concerné majoritairement les paires anglais→chinois, anglais→français, et français→chinois, comme le montre le Tableau 2.

Tableau 2 : Données statistiques sur les segments post-édités dans SECTra_w depuis 2010

Paire de langues (L1→L2)	Bisegments	Mots source L1 (Pages standard)	Mots cible L2 (Pages standard)	Taille de L1	Taille de L2
anglais → français	121 074	2 542 731 (10 170 p.)	2 613 351 (10 453 p.)	10,1Mo	10,4 Mo
anglais → chinois	208 106	4 370 530 (17 482 p.)	6 063 942 (151 159 p.)	19,1 Mo	17,6 Mo
français → anglais	29 079	627 661 (2510 p.)	610 098 (2 440 p.)	4 Mo	3,9 Mo
français → chinois	10 890	228 703 (914 p.)	317 322 (793 p.)	1,5 Mo	1,25 Mo
chinois → anglais	2 013	58 656 (146 p.)	42 275 (169 p.)	240 Mo	263 Mo
chinois → français	10 062	291 192 (727 p.)	211 185 (844 p.)	874 Mo	1 Mo

1.1.1.2 Résumé des avancées et des limitations au niveau conceptuel

1.1.1.2.1 Avancées

Structure des données. Les données sont organisées autour des MT. L'administrateur peut créer une MT et la « dédier » à un corpus (comme les 25 articles de EOLSS) ou à un pseudo-corpus (comme les pages Web accessibles par un certain ensemble d'url), ou bien il la déclare comme « partageable » par plusieurs IMAG. Dans la partie SECTRA_w/IMAG, il n'y a pas de corpus bien identifié au sens classique, sauf pour les MT dédiées. En effet, dans le cas d'une MT partagée par plusieurs IMAG, ce sont toutes les pages Web (des "pseudo-documents" et pas des documents, d'ailleurs) visitées depuis ces IMAG qui constituent "le corpus", qui devient de fait l'ensemble des pages Web dont les segments ont été mis dans cette MT. Il y a bien une avancée, puisqu'on peut traiter des pseudo-documents et partager des MT, mais elle s'accompagne d'une confusion, sur laquelle nous reviendrons plus loin.

Au moment de la post-édition d'une page Web sous SECTRA_w, on édite en fait la partie de la MT constituée de l'ensemble (sans répétition) des segments apparaissant dans l'instance de ce pseudo-document. Cet ensemble est présenté en une suite de "pages logiques", chacune contenant un nombre maximum de segments, paramétrable. La valeur par défaut est actuellement 20, car une page standard contient en moyenne 20 phrases de 12 à 13 mots. Mais cela peut varier selon le type de document. Si un segment apparaît plusieurs fois dans une page Web, la suite des pages logiques ne correspond donc pas à la suite des segments dans la page.

C'est un point apparemment positif, puisqu'un segment répété ne sera post-édité qu'une fois. Cependant, il serait préférable de présenter la suite des segments dans l'ordre du document, sans éliminer les segments répétés, mais en présentant toutes leurs post-éditions dans leurs occurrences successives. En effet, deux occurrences d'un même segment dans une page Web,

et plus généralement dans un document, peuvent devoir être traduites différemment. Par exemple, dans le site Web du LIG, le segment « Recherche » doit être traduit par « *Search* » s'il apparaît dans un onglet de navigation, et par « *Research* » si c'est un item du plan du site.

Un point très important est que chaque segment source a un identifiant unique qui est utilisé pour le lier avec les éléments qu'on considère comme ses annotations (des traductions automatiques avec leur origine, des post-éditions avec leurs scores et leur dernier contributeur, et on pourrait avoir une liste de versions, un minidictionnaire, un graphe UNL, etc.).

Appel de systèmes de TA. SECTRA_W permet de paramétrer la liste des systèmes de TA appelés (comme GT, SYSTRAN, et REVERSO) pour produire des traductions candidates. Pour cela, il utilise l'intergiciel TRADOH (originellement créé en 2002-2003 par Hung VO-TRUNG (Vo Trung, 2004) dans le cadre de sa thèse, puis revu par nous en 2013.

Support de l'évaluation de systèmes de TA. SECTRA_W supporte l'évaluation de systèmes de TA depuis sa première version, mais cette partie est restée indépendante du reste. En particulier, la partie concernant la PE contributive ne permet pas de faire de l'évaluation subjective classique. On se contente d'un « score de qualité » défini à partir du profil du dernier contributeur à la PE d'un segment, et modifiable par les contributeurs inscrits et connectés.

La partie dédiée à l'évaluation reste très intéressante. Décrivons-la brièvement. SECTRA/EVAL intègre dès octobre 2007 des outils d'évaluation subjective⁷ et objective⁸. On peut effectuer les deux types les plus courants d'évaluation subjective (adéquation, fluidité) en utilisant plusieurs juges, pour un ou plusieurs systèmes de TA à la fois. L'administrateur d'une campagne d'évaluation peut interdire ou permettre à plusieurs juges d'effectuer l'évaluation subjective sur la même unité de données (segment). Les libellés et le nombre de réponses possibles sont paramétrables. Au total, on peut définir de nombreuses configurations de campagnes d'évaluation. Il y a bien sûr aussi les outils d'import, d'export, et de calcul de résultats.

SECTRA_W/TRAD fournit de plus des interfaces graphiques pour faciliter la sélection de données (tout ou partie d'un corpus), et le lancement des programmes de calcul des mesures d'évaluation objective sur les données sélectionnées.

« **Base corporale** ». SECTRA_W est une sorte de « base corporale⁹ », Il fournit des services Web pour l'exploitation de corpus multilingues, l'amélioration des traductions, et l'extension « en largeur » à d'autres langues, par appel à des serveurs de TA. Depuis 2008, SECTRA_W permet la post-édition collaborative de pages Web à l'aide d'un ou plusieurs systèmes de TA, et leur évaluation par la distance de post-édition, qui est assez bien corrélée au temps passé à la post-édition, ce qui donne une bonne mesure de « qualité d'usage », pour la post-édition tout au moins.

SECTRA_W, et d'ailleurs aucun système d'évaluation de TA, ne donne de moyen d'évaluer la qualité d'usage quand l'utilisateur est une personne ne connaissant pas la langue source et cherchant à comprendre une page Web. Dans le cas où il s'agit d'un système de e-commerce,

⁷ C'est-à-dire, faisant appel à des jugements humains, comme l'évaluation d'adéquation, de fluidité, de fidélité, de grammaticalité, de précision terminologique, etc.

⁸ Il s'agit non seulement des mesures de type BLEU, NIST ou METEOR, fondées sur des calculs de similarité avec des traductions de référence, et donnant des résultats sur tout un corpus, mais pas sur tel ou tel bisegment, mais aussi de mesures liées à la tâche, comme la "distance mixte de post-édition" introduite par Boitet et Pineau en 2004 (Pineau, 2004) Pineau, M. (2004). Comparaison de résultats de traduction (automatique ou non). Université Joseph Fourier, Rapport de TER..

⁹ Terme proposé par Ch. Boitet pour faire pendant à celui de "base lexicale".

on devrait pouvoir le faire en mesurant le nombre d'achats par rapport au nombre et à la durée des visites, mais ce n'est qu'une hypothèse. Dans d'autres cas, par exemple l'accès à des manuels scientifiques, on pourrait évaluer la compréhension "via la traduction" par une technique de QCM¹⁰.

Dans cette thèse, nous nous limitons à l'évaluation de la qualité d'usage quand les utilisateurs font de la post-édition.

Notions. Lors de la conception de SECTRA_W, C. P. HUYNH a clarifié et bien défini plusieurs notions, comme segment monolingue multilingualisé, pseudo-document, métasegment, super-segment, infra-segment, etc.

Nous les reprenons ici, et en introduisons quelques-unes dont la nécessité est apparue durant le projet TRAQUIERO.

1.1.1.2.2 Définitions de diverses notions

Les définitions nouvelles sont signalées par une étoile.

a. Phrase et groupe

Définition 1. Une *phrase* (*sentence* en anglais) est l'unité élémentaire d'un énoncé, formée de plusieurs mots ou groupes de mots, et qui présente un sens complet. (*TheFreeDictionary*¹¹).

Définition 2*. Un *syntagme* ou *groupe* (*phrase* en anglais) est un constituant possible d'une phrase. En général, un *titre* est un groupe nominal.

b. Segment, fragment, élément non textuel (hors-texte)

Définition 3. Un *segment* est l'unité de traduction de base des traducteurs humains. Il s'agit d'une phrase, d'un titre, ou d'un terme dans une nomenclature.

Définition 4*. Un *fragment* (*chunk* en anglais) est une partie d'un segment, qui peut être un *groupe syntaxique* (*phrase* en anglais) ou un simple n-gramme.

Définition 5*. Un segment peut contenir des *éléments non textuels*, ou *hors-texte*, comme des images, des formules, ou des balises, qui ont un rôle linguistique et une valeur non linguistique.

Ainsi, une image ou une expression mathématique fonctionne usuellement comme un groupe nominal, une relation mathématique fonctionne comme un groupe nominal ou un groupe verbal, et une balise peut seulement contrôler la présentation (gras, souligné) ou avoir aussi un rôle de ponctuation (élément de liste à puces par exemple).

On appelle aussi « segment » (en TMX ou en XLIFF) un segment source accompagné d'une ou plusieurs traductions. Il faut donc raffiner la définition précédente.

¹⁰ Question à choix multiples

¹¹ <http://fr.thefreedictionary.com/phrase>

Dans le format TMX¹² (*Translation Memory eXchange*), un segment multilingue est représenté par un élément `<tu>`.

Par exemple :

```
<tu>
  <tuv xml:lang="en"><seg>How are you?</seg></tuv>
  <tuv xml:lang="fr"><seg>Comment vas tu?</seg></tuv>
  <tuv xml:lang="vi"><seg>Chú có khỏe không?</seg></tuv>
</tu>
```

Si $N = 2, 3, \dots$, on parle de segments bilingues, trilingues, etc. On appellera donc *segment monolingue* un segment multilingue réduit à un seul segment (source).

Définition 6. Un *segment monolingue* est un segment dont le contenu textuel est en une seule langue.

Définition 7*. Un *segment multilingue* est un segment dont le contenu textuel est dans plusieurs langues, chaque version étant considérée comme contenant exactement la même information, exprimée de façon correcte.

C'est par exemple le cas des segments contenus dans des documents de brevet (voir VII.2.1).

Définition 8. Un *segment monolingue multilingualisé (annoté)* est un objet contenant un segment « source » primaire, une ou plusieurs traductions (automatiques ou humaines ou automatiques post-éditées) pour une ou plusieurs langues, et des annotations, en général des objets, comme des arbres linguistiques, des graphes UNL, des résultats d'évaluation(s), et des références aux contributions ayant produit chaque objet non primaire.

Définition 9*. Un *segment multilingue multilingualisé (annoté)* est un objet contenant un segment multilingue, dans N langues « sources », et, dans M autres langues, une ou plusieurs traductions (automatiques ou humaines ou automatiques post-éditées), ainsi que des annotations, comme celles d'un segment monolingue multilingualisé et annoté.

Définition 10*. Le *chemin traductionnel* d'une annotation, en particulier d'une traduction ou d'une post-édition, est la suite des opérations l'ayant produite, ainsi que les intervenants humains impliqués, et les éventuels objets auxiliaires utilisés.

Par exemple, une traduction français-chinois peut avoir été produite par le chemin traductionnel :

```
TA (Systran, fr_en); PE(Zhong, Sectra); TA(Neon, en_zh); PE(Wang, Sectra);
```

c. *Métasegment, document, métadocument*

Définition 11*. Un *métasegment* est un segment comportant une ou plusieurs *variables*, éventuellement typées (nombre, date, balise faible...).

Définition 12*. Un *document* est un ensemble formé par un support et une information, celle-ci enregistrée de manière persistante. Nous nous intéressons aux documents textuels, qui contiennent des "segments" textuels.

Définition 13*. Un *métadocument* est un document pouvant contenir des métasegments.

Définition 14*. Un *pseudo-document* est défini par une référence (nom de fichier, url, uri) à un document qui peut varier au cours du temps.

Par exemple, une page Web peut changer un peu, pas du tout ou totalement d'un appel au suivant.

¹² http://www.lisa.org/fileadmin/standards/tmx1.4/tmx.htm#Misc_GroupTU

- Un petit changement peut être dû au fait que le document est en fait un métadocument, et qu'au moins une valeur d'une variable d'un métasegment a changé (par exemple, la date, ou le nombre de visites).
- Un changement total peut être dû au fait que l'instance précédente a été remplacée par un document tout à fait différent. C'est par exemple le cas d'une page Web contenant chaque jour des nouvelles, ou un nouvel éditorial (comme le "World Web" de l'Unesco).

d. *Contexte*

Définition 15. *Contexte* : Le contexte $m-n$ d'un segment source par rapport à un document, ou plus généralement à une instance d'un pseudo-document, est défini par :

- la liste des m segments (de même langue) qui le précèdent.
- la liste des n segments (de même langue) qui le suivent.
- l'instanciation des variables, s'il y en a, dans ces $m+n+1$ segments.

Le contexte $m-n$ d'un segment (cible) dans une version résultat de TA ou de PE, dans un segment monolingue, ou multilingue multilingualisé, est défini par

- la liste des m segments (de même langue et de même version) qui le précèdent.
- la liste des n segments (de même langue et de même version) qui le suivent.
- l'instanciation des variables, s'il y en a, dans ces $m+n+1$ segments pour la même version.

Deux résultats de TA appartiennent à la même version s'ils ont été produits par la même instance du même système de TA.

Deux post-éditions appartiennent à la même version si elles apparaissent dans la même version d'un document (ou pseudo-document) post-édité.

Par conséquent, le contexte d'une PE peut dépendre du choix des PE des segments précédents et suivants. Exemple : 2 réviseurs choisissent des PE différentes dans la MT et produisent 2 versions (tout à fait correctes) du même document.

1.1.1.2.3 *Limitations*

Absence d'une définition claire des "corpus" et des "corpus de traductions". Cela est dû au fait que SECTRA_W d'abord été développé pour ce qu'on appelle des corpus parallèles en TA statistique, c'est-à-dire des listes de segments se correspondant dans une ou plusieurs langues, sur le modèle du BTEC (*Basic Travel Expression Corpus*) d'ATR (Boitet, Boguslavskij et al., 2007). Dans un second temps, on a traité un "vrai" corpus (monolingue), formé de 25 articles de l'encyclopédie EOLSS, présentés dans autant de fichiers *.aspx* et de fichiers compagnons *.unl*, auxquels on a attaché les traductions en français de chaque segment. Enfin, on a traité des corpus en quelque sorte virtuels, formés des pseudo-documents correspondant aux pages Web de tel ou tel site Web.

Définition 16*. Un *corpus* est un ensemble usuellement fermé de documents homogènes du point de vue de leur structure, de leur(s) langue(s), de leur genre et de leur domaine.

Par exemple, on peut parler du corpus (monolingue) des articles du Monde sur la culture de 2000 à 2002. Le caractère fermé est relatif, car on considère que certains corpus croissent par adjonction de nouveaux ensembles. Par exemple, le corpus bilingue HANSARD des débats du Parlement canadien peut être considéré comme une suite de corpus annuels, le dernier étant en construction.

Définition 17*. Un *corpus de traductions* est un corpus au sens précédent, contenant les traductions de tout ou partie de ses segments, dans une ou plusieurs langues.

Absence de traitement des corpus à proprement parler. SECTRA_W permet l'exploitation de corpus existants, l'évaluation des traductions et l'élargissement à plus de langues, mais n'est pas assez générique. Par exemple, il ne permet pas la création de nouveaux corpus (au sens classique), ni l'annotation selon des critères librement définissables par un linguiste. Ce progrès était prévu dans la thèse de Phap, mais n'a pas encore été réalisé.

Décalage de segments. SECTRA_W appelle les systèmes de TA pour produire des pré-traductions, et le fait maintenant via l'intergiciel TRADOH, mais il lui manque un sous-système pour synchroniser et gérer les segments source et les résultats de traduction. Parfois, un système de TA renvoie plusieurs résultats de traduction pour un même segment, et alors SECTRA_W les aligne avec d'autres, ce qui détruit la cohérence des bisegments, au moins jusqu'à ce que ces autres segments soient (re)traduits et que leurs résultats soient bien alignés.

Confusion de notions. En SECTRA_W, il y a plusieurs confusions entre des notions très différentes, ce qui bloque certaines avancées souhaitables. Par exemple, il y a confusion entre "mémoire de traductions" et « projet ». Quand l'utilisateur crée un nouveau « projet », en effet, le système crée une nouvelle MT. Mais un projet devrait pouvoir concerner plusieurs MT à la fois, ou une seule MT, existante, si par exemple on veut y ajouter des annotations telles que des arbres linguistiques, ou des caractérisations d'erreur. Or, cela ne justifie pas la création de nouvelles MT.

Il y a aussi dans l'actuel SECTRA_W une confusion entre « corpus » et « mémoire de traductions ». Cela est dû, historiquement, à l'évolution vers le traitement de pages Web. Auparavant, dans le projet EOLSS par exemple, on avait un vrai « corpus », c'est-à-dire une collection organisée de documents, avec leurs fichiers satellites (images, etc.) et leurs fichiers compagnons (pour EOLSS, des fichiers de graphes UNL). Après segmentation, on introduisait les segments dans la MT, et les documents n'étaient plus supposés changer. On pouvait donc facilement noter dans la base de données toutes les occurrences d'un segment¹³.

Avec les pages Web, on a dû traiter des *pseudo-documents*, c'est à dire des documents qui ont le même nom (une *url*), mais qui peuvent changer d'un appel à l'autre. On a essayé d'adapter la méthode précédente, mais la solution retenue est source d'erreurs et augmente considérablement le temps de recherche d'un segment. Elle consiste en effet à associer une petite MT à chaque pseudo-document, liant ainsi corpus et MT. Quand on traite une page Web, on regarde dans sa MT s'il a déjà été vu, et si oui on prend la meilleure PE qu'on y trouve. Mais ce segment a pu être déjà post-édité dans un autre pseudo-document (une autre page Web), et du coup le système peut présenter une sortie de TA au lieu d'une PE déjà faite. De plus, si le système ne trouve pas le segment dans la MT du pseudo-document, il le cherche dans tous les autres documents du même "corpus", ce qui peut faire beaucoup, et cela d'autant plus que, si la MT est partagée, il y a des pseudo-documents liés à d'autres IMAG qui sont examinés, alors que, conceptuellement, ils sont dans un corpus différent.

La solution serait de reprendre toutes ces MT, de les fusionner avec la MT principale, et d'associer simplement à chaque pseudo-document la liste des segments déjà vus dans une des instances de ce pseudo-document, et, pour chaque instance, de construire un « document squelette » constitué du code html et de références aux segments présents dans cette instance.

¹³ TM d'IBM et SDL TRADOS le font, ce qui permet de reconstituer la suite des segments des documents déjà traités, mais LIBELLEX ne le fait que depuis 2013.

1.1.1.3 Résumé des avancées et des limitations au niveau algorithmique (conception informatique)

1.1.1.3.1 Avancées

Post-édition "sans couture". SECTRA_W communique avec la partie IMAG, qui peut être considérée comme un "frontal" permettant la post-édition directement sur la page Web, ce qu'on appelle la post-édition "sans couture". Cela veut dire qu'on ne change pas de contexte entre lecture et correction d'une page Web. De plus, quand l'utilisateur fait de la post-édition dans la page Web, il peut visualiser la page Web source en parallèle. Un autre aspect de cette communication est qu'on peut passer de la post-édition d'un segment sur la page Web à sa post-édition dans la page (logique) de post-édition "avancée" de SECTRA contenant ce segment, puis revenir à la page Web, à l'endroit où on était, une fois qu'on a post-édité quelques segments en mode "avancé". En pratique, les post-éditeurs experts utilisent les deux modes. Le mode avancé leur permet de post-éditer très vite de nombreux segments, et le mode sans couture leur permet de voir l'ensemble de la page Web, et de corriger de petites erreurs qui sautent aux yeux alors qu'on les repère mal dans les petites zones de texte (*textareas*) de l'interface de PE de SECTRA_W.

Possibilité de partager des MT. Le fait qu'une MT puisse être partagée par une ou plusieurs IMAG est un avantage si les sites Web "élus" sont comparables, et plus précisément si leurs sous-langages diffèrent peu. C'est par exemple le cas de sites de tourisme, ou de sites de villes, ou de sites de cours d'informatique en ligne. En effet, des segments identiques (parfois longs) apparaissent alors, et la post-édition d'un site peut bénéficier à un autre.

Gestion des utilisateurs. SECTRA_W utilise le même type de gestion des utilisateurs que XWiki, avec deux niveaux : individu et groupe, un individu pouvant appartenir à plusieurs groupes. Il serait préférable de séparer SECTRA_W du logiciel IMAG et de créer un service unique faisant le lien entre ces gestionnaires d'utilisateurs, par exemple fondé sur un annuaire du genre LDAP, de façon à éviter de demander à un contributeur de s'authentifier plusieurs fois alors qu'il a tous les droits nécessaires.

SECTRA_W associe à chaque utilisateur un profil général incluant des informations sur ses compétences et ses droits en général (par exemple, un utilisateur peut accéder à la MT DEMO1, mais pas aux autres MT). Il faudrait aussi associer à chaque utilisateur un raffinement de son profil, pour chaque projet. Par exemple, quelqu'un pourrait être administrateur mais pas post-éditeur dans un projet, et l'inverse dans un autre.

1.1.1.3.2 Limitations

Structure de la base de données. Quand l'administrateur crée un « projet » dans SECTRA_W, il crée une nouvelle « mémoire de traductions », qui donne lieu à la création d'une nouvelle base de données (BD) en MySQL. Quand l'utilisateur crée une IMAG associée à cette MT, SECTRA_W crée 3 tables dans cette BD. La première table contient tous les segments source. La deuxième contient tous les segments cible traduits par un ou plusieurs système de TA. La troisième contient les post-éditions des segments post-édités en au moins une langue cible.

Par exemple, le site Web du LIG est en français. Ce site contient une centaine de pages Web, et son IMAG, donc aussi SECTRA_W, prévoit plus de 30 langues cible. Donc, à une page Web du site LIG, SECTRA_W associe une table source, une ou plusieurs tables contenant les TA, et plus de 30 tables contenant les segments post-édités. Cette structure de BD est difficile à manipuler, et augmente la complexité de la recherche de l'information. Cette limitation provient d'un mauvais choix au niveau du schéma conceptuel, et par suite du schéma logique des MT.

Passage à l'échelle. SECTRA_W, dans son état de 2011, était donc basé sur un schéma conceptuel de la structure des données, qui a mené à un schéma logique multipliant les tables, et ne permettant pas de traiter efficacement de grosses quantités de données. Nous savions donc qu'il faudrait intervenir à ce niveau pour pouvoir réellement passer à l'échelle.

1.1.1.4 Résumé des avancées et des limitations au niveau de l'implémentation

1.1.1.4.1 Avancées

Import de corpus pour l'évaluation. SECTRA_W permet d'importer divers corpus multilingues disponibles et de faciliter la sélection des corpus source pour organiser des campagnes d'évaluation. SECTRA_W accepte des fichiers d'entrée en format texte (.txt). Au niveau physique, chaque corpus¹⁴ peut être constitué par un fichier source contenant des segments source, un ou plusieurs fichiers de traductions "candidates" (produites par différents systèmes de TA ou vers différentes langues) à évaluer, et un ou plusieurs fichiers de traductions de référence. Cependant, un corpus d'évaluation peut ne contenir que des segments source, car SECTRA_W permet de produire des traductions candidates par appel de TA, et des traductions de référence par post-édition en ligne. Au niveau interne, les segments dans les fichiers (source, candidats, et références) sont alignés.

Convivialité et souplesse de l'interface. L'interface de SECTRA_W permet de changer la taille des colonnes, et de cacher/montrer des colonnes en utilisant la souris ou un tableau de configuration. Les utilisateurs peuvent choisir le nombre de segments affichés à chaque fois sur l'interface. Leurs travaux sont enregistrés automatiquement et les données déjà manipulées sont distinguées des autres au niveau de la couleur du fond.

Mode de lecture et mode avancé. SECTRA_W/IMAG offre deux interfaces de post-édition. Le premier est intégré au mode de lecture. L'utilisateur peut faire la post-édition sur le document, et voir très facilement l'effet de la post-édition dans le contexte de lecture. Le deuxième est un mode « avancé ». L'utilisateur peut se concentrer sur le segment. Les segments sont bien présentés dans les pseudo-documents.

API. XWiki permet de publier des services, et ainsi d'outiller leur utilisation par des logiciels tiers. Cela est fait dans SECTRA_W pour les fonctions de recherche exacte d'un texte dans un corpus (*TM exact search by API*) et pour mettre à jour la traduction mémorisée pour un texte dans un corpus (*Update TM by API*).

1.1.1.4.2 Limitations

XWiki et serveur de déploiement. SECTRA_W a été développé en utilisant la plate-forme XWiki (version 1.3.1), la technologie AJAX, et les langages de script JAVASCRIPT, GROOVY et VELOCITY. Depuis 2010, SECTRA_W n'a pas suivi les mises à jour de XWiki. Passer à une version plus récente de XWiki comme la version 7 serait actuellement extrêmement lourd et risqué, car le langage VELOCITY (d'Apache Foundation) qui permet de traiter les BD MYSQL a lui-même changé. C'est pour cela que SECTRA_W est déployé sur TOMCAT 5, correspondant à l'ancien VELOCITY, et qu'on ne peut pas passer à TOMCAT 6 ou 7 pour l'instant. Enfin, une autre limitation est que XWiki n'est pas très compatible avec certains navigateurs tels que IE, car il utilise des fonctions évoluées (comme AJAX) qui ne sont pas encore standardisées sur les différents navigateurs.

¹⁴ Ce terme est utilisé par abus de langage dans le contexte des campagnes d'évaluation. Il ne s'agit de rien d'autre que d'une suite de tests formée d'un **ensemble** de paires (segment source, segment traduit), et pas d'un ensemble de **documents**, un document n'étant jamais réductible à l'ensemble de ses phrases.

Export. SECTRA_w offre une fonction capable d'exporter le résultat de la campagne TRANSAT, mais n'a pas de fonction permettant de réaliser l'export sur les MT génériques, comme la grande MT LAMETRO. Il n'y a pas non plus de moyen de choisir les segments post-édités.

Absence de l'aide lexicale proactive annoncée. Dans sa thèse, C.P. HUYNH a présenté la conception d'une aide lexicale proactive intégrable à SECTRA_w, sous la forme d'un "minidictionnaire" précalculé et associé à chaque segment. De son côté, H.T. NGUYEN avait réalisé une fonction permettant de produire un minidictionnaire à partir d'une liste de mots, par lemmatisation, recherche dans une BD lexicale PIVAX, et formatage en HTML. Mais la jonction des deux n'a pas pu être réalisée avant leur départ du laboratoire.

I.1.2 Accès multilingue à des sites Web : le logiciel iMAG

I.1.2.1 Présentation générale

I.1.2.1.1 Bref historique et motivations

Le concept d'iMAG a été proposé par Ch. Boitet et V. Bellynck en 2005 (Boitet et al., 2005), étudié en 2006 et 2007 par les brats de projets d'étudiants, et a atteint l'état de prototype opérationnel grâce à C.P. HUYNH en novembre 2008 (Boitet et al., 2008), avec une première démonstration sur le site Web du laboratoire LIG. Il a été adapté au site Web DSR (Digital Silk Road) en avril 2009 par C.P. HUYNH, puis à plus de 80 autres sites Web grâce au configurateur réalisé par H.T. NGUYEN. Ces iMAG servent à des démonstrations et à des expérimentations. Citons ici (Boitet et al., 2010).

Une iMAG est une passerelle interactive d'accès multilingue (interactive Multilingual Access Gateway), ressemblant beaucoup à Google Translate, à première vue : on donne une URL (site Web de départ) et une langue d'accès, et on navigue ensuite dans cette langue d'accès. Lorsque le curseur passe sur un segment (le plus souvent une phrase ou un titre), une palette montre le segment source et propose de contribuer en corrigeant le segment cible, en fait en post-éditant un résultat de traduction automatique (TA). Avec Google Translate, la page ne change pas après la contribution, et si une autre page contient le même segment, sa traduction est toujours le résultat de TA grossière, pas la version polie post-éditée. La boîte à outils de traduction Google Translation Toolkit permet de traduire par TA et ensuite de post-éditer en ligne des pages Web complètes tirées de sites tels que Wikipédia, mais, de nouveau, les segments corrigés n'apparaissent pas quand on regarde plus tard la page de Wikipédia dans la langue d'accès.

En revanche, une instance d'iMAG (dite iMAG-S) est dédiée à un site Web élu (S), ou plutôt au sous-langage défini par une ou plusieurs URL et leur contenu textuel. C'est un bon outil pour rendre S accessible dans beaucoup de langues, immédiatement et sans responsabilité éditoriale. Les visiteurs de S ainsi que des post-éditeurs et des modérateurs payés ou non contribuent à l'amélioration continue et incrémentale des segments textuels les plus importants, et éventuellement de tous.

Une instance iMAG-S contient une MT. Les segments sont pré-traduits non pas par un système de TA unique, mais par un ensemble (sélectionnable) de systèmes de TA gratuits. Systran et Google sont principalement utilisés aujourd'hui, mais des systèmes spécialisés, développés à partir de la MT post-éditée, vont bientôt être utilisés. Les visiteurs de S ainsi que des post-éditeurs et des modérateurs payés ou non contribuent à l'amélioration continue et incrémentale des segments textuels les plus importants, et éventuellement de tous.

Les plates-formes contributives puissantes SECTra_w sont utilisées pour supporter les MT. Les pages traduites sont construites avec les meilleures traductions des segments disponibles au moment de l'accès. Pendant la lecture d'une page traduite, il est possible non seulement de contribuer au segment sous le curseur, mais aussi de passer de façon transparente sous l'environnement de post-édition en ligne de

SECTra_w, muni de bonnes fonctions de filtrage et de recherche-remplacement, et ensuite de revenir dans le contexte de lecture.

1.1.2.1.2 Lancement de maquetages

Deux « maquettes » d'IMAG ont été étudiées avant novembre 2008, par Mohammad Daoud dans son M2R (Daoud, 2007), et par Carlos Ramisch dans son stage ENSIMAG (Ramisch, 2008).

Le premier n'a pas pu produire la maquette prévue, car il avait voulu écrire lui-même un segmenteur de pages Web (en "segments textuels", c'est à dire en phrases ou titres), alors qu'on l'avait prévenu qu'il s'agissait d'un problème difficile et qu'on lui avait demandé d'utiliser indirectement GT comme segmenteur.

Dans le cadre d'un projet du LIG appelé "IMAG/LIG", financé par 5 ou 6 mois de bourses CNRS du LIG, C. Ramisch a voulu lui aussi construire un prototype complet à partir de zéro. Il avait prévu d'utiliser LINGPIPE¹⁵ pour la segmentation, mais il aurait fallu l'adapter aux pages Web. Il n'a pas eu le temps d'implémenter ce qu'il avait spécifié, car c'était un trop gros travail pour le temps disponible. Malheureusement, une fois son projet soutenu, il n'a pas pu "passer la main", étant pris à la fois par les cours de M2R et un stage à XRCE.

1.1.2.1.3 Prototype opérationnel de C. P. HUYNH

Durant l'été 2008, dans le cadre du même projet IMAG/LIG, C. P. HUYNH et H. T. NGUYEN ont travaillé sur le site Web du LIG, pour segmenter ses pages et aligner les segments déjà traduits en anglais. Le meilleur segmenteur qu'ils avaient trouvé était justement celui de GT. L'idée était qu'ils fournissent au logiciel IMAG en construction une MT initiale pour le site du LIG, en appelant GT sur les segments non traduits. Comme ce logiciel attendu n'était pas disponible, Phap a proposé de le réaliser très vite, par extension de SECTRA, qui était déjà assez mûr, ayant été utilisé par des dizaines de contributeurs dans le cadre du projet EOLSS-UNDL.

En à peine 3 semaines, juste avant de partir pour un stage prédoctoral international au Japon (NII), il a produit un premier prototype opérationnel, et une première IMAG, l'IMAG-LIG, qui a pu être démontrée à la direction du LIG (B. Plateau) vers le 20/11/2008.

Le fonctionnement est le suivant. La page consultée est accédée via la passerelle IMAG, qui réécrit les url de façon à ce que les pages accédées par la suite lors de la navigation passent aussi par elle. La page à traduire est envoyée à GT, on récupère le résultat, puis on détecte les segments marqués par GT. On les cherche alors dans la MT, et on les ajoute s'ils n'y sont pas. Dans les deux cas, on met aussi dans la MT associée à l'IMAG les (pré)traductions produites par GT. La page traduite est reconstruite à partir du code HTML de la page source, des meilleures traductions des segments (trouvées dans SECTRA_W), et de code JavaScript spécifique.

Au tout début, l'IMAG-LIG produite par le projet IMAG/LIG contenait la traduction d'environ environ 2000 segments du site Web du LIG. Elle en contient maintenant près de 15000.

1.1.2.1.4 Évolutions depuis début 2010

Début 2010, une IMAG dédiée a été créée pour la Métro de façon à permettre l'accès à son site Web en chinois et en anglais à l'occasion de l'exposition universelle de Shanghai qui eut lieu en octobre 2010. Ce fut la première expérience en grandeur réelle.

¹⁵ <http://alias-i.com/lingpipe/>

Ensuite, plusieurs sites IMAG à mémoire de traductions dédiée ont été créés. La table suivante donne les URL des sites principaux.

Tableau 3 : Liste des IMAG à MT dédiée construites depuis 2010

Nom abrégé	Accès au site élu : remplacer dans http://service.aximag.fr/xwiki/bin/view/imag/home home par :	Page d'accueil de chaque site :
Expo de Shanghai	lametro	http://www.lametro.fr
Projet Traouiero	traouiero	http://getalp.imag.fr/traouiero/
Projet MACAU	macau	http://tools.aximag.fr/macau/chamilo-macau/index.php
AMIES	AMIES	http://www.agence-maths-entreprises.fr
Grenoble INP	GINP	http://www.grenoble-inp.fr/accueil/
Powers	Powers	http://www-clips.imag.fr/geod/User/laurent.besacier/TRANSLATION-EXP/The_Book_of_Me.html
Michelin	Michelin	http://www.michelin.fr
Bio Clean	Bio-Clean	http://projet-bioclean.3beesonline.com/fr/accueil/
Projet Akenou	akenou	http://getalp.imag.fr/akenou/

De nombreuses IMAG ont aussi été créées sur des MT partagées. En particulier, une a été créée dans le cadre d'un projet de l'ISCC (institut des sciences de la communication), un institut du CNRS dirigé par Dominique WOLTON.

1.1.2.2 Conception générale avec implémentation "intégrée"

Voici un peu plus de détails sur l'implémentation du logiciel IMAG, "intégré" à SECTRA. Lors du projet ANR TRAQUIERO, V. Bellyneck a prévu de la rendre plus modulaire, grâce à l'introduction d'un "relais IMAG", de l'utilisation de TRADOH pour appeler les systèmes de TA, et de l'utilisation de SEGDOC pour effectuer la segmentation-normalisation, mais cette transformation n'est pas encore terminée.

Voici le schéma de l'architecture logicielle actuelle.

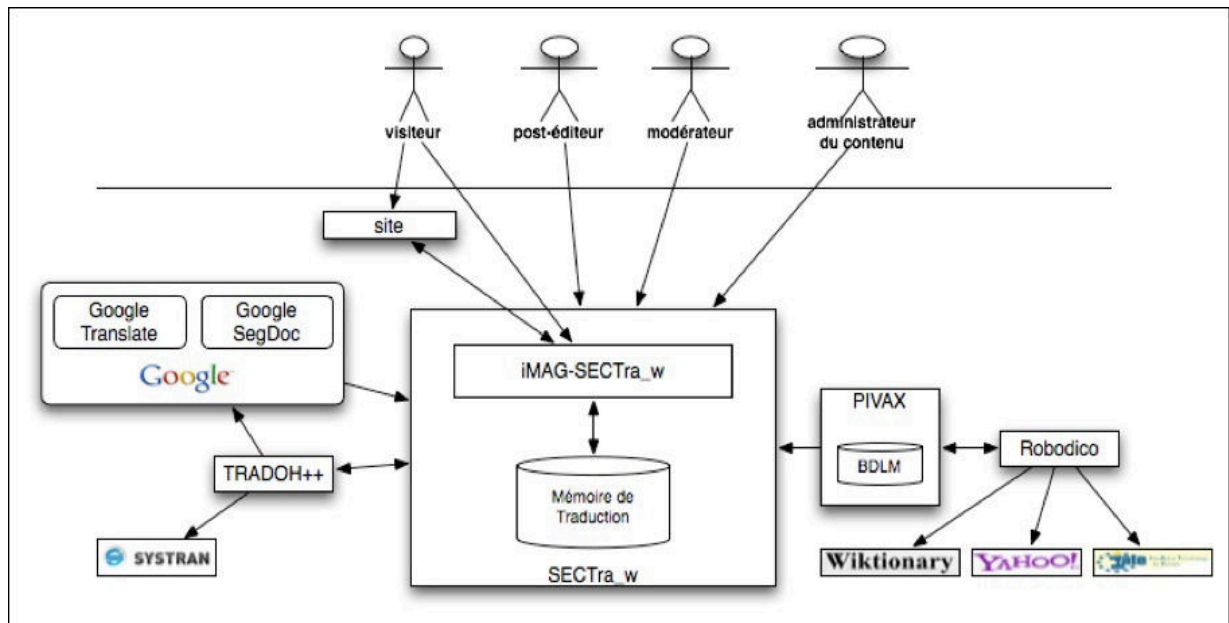


Figure 1 : Architecture générale d'une iMAG pour un site élu

Lorsque le visiteur demande à accéder au site élu S dans une autre langue que la langue source, la page d'accueil de S est envoyée à GT de façon à récupérer le découpage en segments de texte, et les traductions de chacun des segments par cet outil, s'il supporte les langues désirées en source et cible. Ultérieurement, lors de la navigation interne au site S, chaque page de S est traitée par le même procédé.

D'autres outils de traduction automatique comme REVERSO, SYSTRAN ou BING peuvent être utilisés pour prétraduire chaque segment. Les segments traduits sont mémorisés comme les correspondants (dans la langue de visite) des segments d'origine.

Dans la page Web reconstruite par l'iMAG, le texte de chaque segment en langue source est remplacé par le texte correspondant en langue cible, et le code nécessaire à l'insertion d'une bulle contextuelle (surgissant au survol par le curseur de chaque segment) est ajouté avec le texte dans la langue source, le couple (segment en langue source, segment en langue cible) constituant une correspondance de segments pour le système SECTRA_W.



Figure 2 : Capture d'écran de l'iMAG LIG-LAB en chinois

La bulle contextuelle est appelée « bulle-iMAG ». Elle montre le texte du segment en langue source, et présente un bouton dont l'action associée transforme la bulle en une bulle de contribution : elle contient alors un formulaire dont une zone de saisie de texte (« *textarea* ») permet de modifier à la volée la traduction proposée. Le visiteur peut alors contribuer à l'amélioration de la traduction de ce segment de la page Web courante selon son droit d'édition et le mode de publication du site. Une fois que la contribution est envoyée, elle est stockée dans une table associée à la page en question, identifiée par son URL comme un pseudo-document, et il y a une duplication dans une table globale des segments source-cible pour chaque site. Donc, dans une page, s'il y a un même segment (par exemple un item de menu, ou une phase type) qui est déjà traduit dans une autre page, la traduction est proposée avec la mention « *Mémoire de traductions* ».¹⁶

Dans l'architecture complète prévue, on a intégré l'aide lexicale proactive, ce qui explique pourquoi PIVAX apparaît dans le schéma suivant, tiré de la thèse de H.T. NGUYEN. On construit des instances spécifiques de composants SECTRA_W (pour la MT), de PIVAX (pour la base lexicale) et d'une iMAG dédiée. Il y a une connexion entre ces instances locales et une instance centrale de données afin d'initialiser les données pour un nouveau site.

Dans cette architecture, SECTRA_W, PIVAX, et les iMAG communiquent les uns avec les autres pour effectuer des tâches. SECTRA_W communique avec PIVAX pour demander des minidictionnaires, et inversement PIVAX communique avec SECTRA_W pour extraire des unités lexicales à partir des corpus.

Les iMAG communiquent avec SECTRA_W pour demander des traductions, soit à partir de MT, soit à partir de systèmes de TA, et inversement SECTRA_W communique avec les iMAG pour demander des segments source, des fichiers squelette, et des dictionnaires de hors-texte.

¹⁶ Remarque : tous les segments d'un site élu sont actuellement supposés écrits dans une même langue.

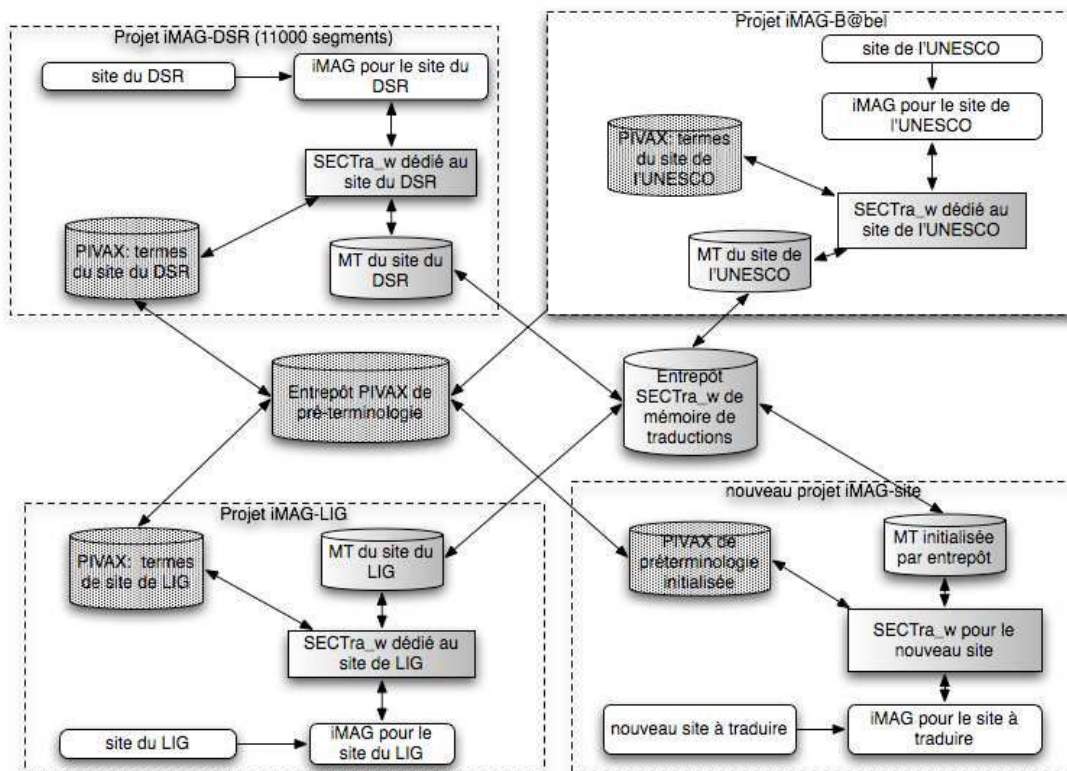


Figure 3 : Architecture par agents SECTra_w, iMAG, PIVAX (Nguyen, 2009)

1.1.2.3 Résumé des avancées et des limitations du logiciel iMAG

1.1.2.3.1 Avancées

Introduction d'une possibilité de modération. Il faut dans certains cas pouvoir éviter que des contributeurs anonymes n'améliorent pas les traductions, mais au contraire les remplacent par des textes injurieux, diffamatoires ou autres. C'est un risque non négligeable si le site élu est, par exemple, lié à des activités politiques. C'est pourquoi, à l'occasion du travail sur l'iMAG de la Métro (le "grand Grenoble"), H. T. NGUYEN a introduit une possibilité de modération, inspirée de celle de Wikipédia. En mode "modéré", un contributeur peut modifier une traduction comme il veut, et il la verra à l'écran, mais elle ne sera acceptée définitivement dans la MT qu'après avoir été validée par un modérateur ayant les droits suffisants.

Visualisation de la fiabilité supposée des segments en langue d'accès. La case à cocher « *fiabilité* » permet de voir, pour chaque segment, s'il est le résultat brut d'une TA, s'il a été produit par un contributeur anonyme (ou enregistré mais pas connecté), ou bien par un contributeur connecté. Cette visualisation ne modifie pas la couleur du fond, ou la couleur des fontes. On se borne à encadrer chaque segment par deux parenthèses spéciales, respectivement rouges, orange et vertes dans les trois cas ci-dessus.

Association d'un niveau de fiabilité supposée et d'un score de qualité à chaque segment traduit. Le terme "niveau de fiabilité" n'est pas très bon, mais nous n'en avons pas trouvé de meilleur. C'est en fait l'origine d'une traduction stockée dans la MT. Il est exprimé par 1 à 5 étoiles : ☆ pour une traduction mot-à-mot (pidgin), ☆☆ pour une traduction automatique, ☆☆☆ pour un humain connaissant les deux langues, ☆☆☆☆ pour un traducteur professionnel, et ☆☆☆☆☆ pour un traducteur certifié par le site Web associé à l'iMAG. Ce niveau est de même nature que les 3 couleurs des parenthèses spéciales, et simplement un peu plus précis. Pour les couleurs, on a voulu se limiter à 3, pour des raisons ergonomiques.

Le *score de qualité* est une note entre 0 et 20. Chaque contributeur inscrit dans SECTRA_W a un score par défaut pour chaque couple de langues sur lequel il peut intervenir, par exemple 12/20 pour "moyen", ou 14/20 pour "bon". Il peut lui-même modifier ce score pour noter ce qu'il pense de la qualité d'un segment qu'il a post-édité (ou qu'un autre a post-édité, mais on le fait plus rarement). Par exemple, on mettra 9/20 si on a conscience de ne pas avoir trouvé un équivalent convenable, et 16/20 si on est sûr d'avoir produit une très bonne traduction.

Possibilité de post-éditer directement sur la page Web. L'utilisateur peut faire de la post-édition sur une page Web en langue cible, directement dans le "contexte de lecture". Ses post-éditions sont présentées immédiatement sur la page Web¹⁷, et le résultat est sauvegardé dans la MT associée. Beaucoup de post-éditeurs utilisent en pratique les deux modes de post-édition. Dans le mode avancé, on va très vite, mais on ne voit pas le contexte global de la page, ni la présentation. Quand on revient sur la page, on voit des coquilles, des pronoms incorrects, etc., et on les corrige directement sur la page.

Visualisation parallèle d'une page en cible et source. Quand on utilise une IMAG, on peut post-éditer sur la page Web cible, en visualisant en parallèle la page Web source.

1.1.2.3.2 Limitations

Le logiciel IMAG actuel présente trois groupes de défauts du point de vue fonctionnel et architectural.

- Les dialogues de l'interface de tous les composants devraient être multilingues. Pour l'instant, ils sont seulement en anglais.
- Des fonctionnalités critiques ne sont pas implémentées, ou pas comme il faudrait pour obtenir les résultats attendus.
 - **Gestion de la mémoire de traductions.** Il faudrait éliminer la « *sous-MT* » associée à chaque pseudo-document.
 - **Contextes linguistiques.** Il faudrait les prendre en compte.
 - **Consultations et contributions lexicales.** Il faudrait intégrer une aide lexicale non seulement au niveau de l'interface de PE de SECTRA, mais aussi dans la "palette IMAG".
 - **Visualisation de l'apport de la PE.** Il faudrait pouvoir voir non seulement la page source, mais aussi la page cible en mode « trace » (TA brute initiale et vue intuitive des modifications dues à la PE).
 - **Mesure des temps de PE.** La communication des temps de PE et des scores de qualité entre l'interface IMAG et SECTRA_W n'était pas satisfaisante.
 - **Liaison entre les deux interfaces de PE.** Il faudrait introduire une liaison directe entre un segment sur la page Web et le segment correspondant dans l'interface de SECTRA_W, et dans les deux sens. Pour l'instant, on est limité à l'accès à une « *page logique* » de SECTRA à partir d'un segment.
- Le logiciel IMAG actuel n'est toujours pas protégé contre des attaques où des robots font des requêtes qui sont autorisées pour tous les visiteurs. Cela conduit à une surcharge du système et donc à des dénis de service.

L'architecture actuelle repose essentiellement sur l'exploitation d'un seul logiciel dans lequel tout est géré, ce qui n'est pas compatible avec les besoins en modularité.

¹⁷ Il faut rafraîchir la page pour que les parenthèses encadrant ce segment réapparaissent.

I.1.3 Travaux comparables et idées directrices pour le futur

I.1.3.1 Travaux comparables sur la création et l'exploitation des corpus multilingues

Pour améliorer notre système, nous avons essayé de étudier systématiquement les travaux comparables pour trouver les points intéressants sur la manipulation des corpus multilingues. Nous avons trouvé 2 types de systèmes, (1) les systèmes de stockage et d'accès, et (2) les systèmes de collecte directe et d'exploitation de corpus multilingues.

I.1.3.1.1 Systèmes de stockage et d'accès

Systèmes de stockage basé sur un service Web. Outre le stockage, ces systèmes permettent la navigation, le téléchargement, et la visualisation des données.

Certains fournissent aussi quelques fonctions utiles, comme la recherche dans un sous-ensemble de textes ou de phrases d'une certaine catégorie, ou la production de concordances. Par exemple, le site Web OPUS¹⁸ (Tiedemann, 2012) contient une grande quantité de corpus multilingues. L'utilisateur peut télécharger les corpus à partir des liens du site Web. Dans le Tableau 4, nous présentons les sites Web principaux permettant de partager des corpus parallèles.

Tableau 4 : Exemples de sites Web de partage de corpus parallèles

Nom du corpus	Lien
Europarl	http://statmt.org/europarl/
JRC-Acquis	https://ec.europa.eu/jrc/en/language-technologies/jrc-acquis
Parallel corpora at PELCRA	http://pelcra.pl/new/
UM-Corpus (Tian et al., 2014)	http://nlp2ct.cis.umac.mo/um-corpus/
The Bible (Resnik et al., 1999)	http://www.umiacs.umd.edu/~resnik/parallel/bible.html

Systèmes de plus haut niveau pour les brevets. Un système de gestion de documents de brevets est un autre type de système de gestion de corpus. Généralement, un tel système contient une série de fonctionnalités permettant de manipuler les brevets, comme la recherche, la comparaison, l'extraction, etc.

Par exemple, le système PATENTSCOPE (Pouliquen and Mazenc, 2011) permet d'accéder aux demandes internationales selon le Traité de coopération en matière de brevets (PCT¹⁹ pour PATENT COOPERATION TREATY) en texte intégral, le jour même de leur publication, ainsi qu'aux documents de brevet des offices de brevets nationaux ou régionaux participants.

Systèmes plus génériques. Ce type de système est beaucoup utilisé pour gérer et partager plusieurs types de données. Par exemple, META-SHARE (Federmann et al., 2012) est utilisé pour partager des corpus. Ce type de système a but pour gérer les publications, l'ensemble des données, des fichiers multimédia, etc. Il permet en particulier d'accéder à des corpus et de les télécharger.

I.1.3.1.2 Systèmes de collecte directe et d'exploitation

Quelques systèmes contiennent des outils ou des langages de scripts utilisables pour collecter et manipuler quelques types de corpus, et en particulier des corpus multilingues (comparables ou parallèles).

Voici trois exemples, ERIM-COLLECTE, de notre labo, BITEXTOR, et PACO².

¹⁸ <http://opus.lingfil.uu.se>

¹⁹ http://fr.wikipedia.org/wiki/Traité_de_coopération_sur_les_brevets

ERIM-Collecte : ERIM-COLLECTE (Fafiotte, 2004) a été construit dans le cadre du projet ERIM (Environnement Réseau pour l'Interprétariat Multimodal), visant à créer un environnement en réseau pour l'aide à la communication orale multilingue, dans lequel on peut collecter des corpus de dialogues parlés spontanés bilingues et multilingues. ERIM-COLLECTE permet l'enregistrement systématique des actes et données de l'interaction, pour tous les participants (deux locuteurs ou plus, un interprète ou plus). L'enregistrement est fait localement lors de la conversation. En fin de dialogue, les descripteurs et fichiers produits localement sont transmis à un serveur de collecte, où ils sont regroupés et structurés.

Bitextor : BITEXTOR (Esplá-Gomis, 2009) est un outil d'extraction de bi-textes à partir de pages Web. Tout d'abord, il télécharge les pages Web visées. Ensuite, il normalise les fichiers téléchargés. Enfin, il compare les fichiers en utilisant divers critères pour extraire les phrases parallèles, que nous appellerons plus précisément "bisegments". Dans la terminologie actuelle, un "corpus parallèle" est une collection de bisegments, et bien plus rarement une collection de vrais documents multilingues alignés au niveau des segments (phrases et titres).

PaCo² : PACO² (*Parallel Corpora Collector*) (San Vicente and Manterola, 2012) est un outil de collecte de corpus parallèles à partir du Web. Il met en œuvre les techniques les plus modernes (à l'état de l'art) permettant de trouver des bi-textes dans un domaine Web qui en contient. De plus, il peut trouver de tels domaines Web automatiquement.

1.1.3.1.3 *Systèmes d'extraction indirecte*

Les corpus comparables sont considérés comme des ressources importantes pour collecter des corpus parallèles, surtout quand on n'en trouve pas directement. Ici, le terme "corpus comparable" dénote une collection de "vrais" documents, par exemple des pages Web, ou des documents en WORD ou en PDF, qui sont de typologies comparables (par exemple, des articles de WIKIPEDIA sur le même sujet, ou des recettes de cuisine, ou de la documentation technique de produits similaires).

Plusieurs méthodes ont été proposées pour extraire (ou construire) un corpus parallèle à partir de corpus comparables : elles utilisent les notions de *documents comparables* (Resnik and Smith, 2003) (Smith et al., 2010), de *phrases comparables* (Zhao and Vogel, 2002), (Munteanu and Marcu, 2005), et des ressources comme des bisegments ou des dictionnaires bilingues (Munteanu and Marcu, 2006), (Cettolo et al., 2010). Dans le cadre de sa thèse, Thi Ngoc Diep DO (Do, 2011) a aussi proposé 3 approches d'extraction de textes parallèles à partir de documents comparables.

1.1.3.2 *Travaux comparables sur la post-édition de TA*

1.1.3.2.1 *Introduction*

« Postediting/post-editing (PE) is by far most commonly referred to as a task related to Machine Translation (MT) and has been previously defined as the "term used for the correction of machine translation output by human linguists/editors." Another good summary statement indicates that "post-editing entails correction of a pre-translated text rather than translation 'from scratch'." In basic terms, the task of the post-editor is to edit, modify and/or correct pre-translated text that has been produced by a machine translation system from a source language into [a] target language[s]. »

– ALLEN, Jeffrey. 2003 (Allen, 2003)

Les premières études sur la post-édition sont apparues dans les années 80, liées aux outils associés aux systèmes de TA, puis aux outils associés aux MT (Senez, 1998). Les premières parurent au Japon, et seulement en japonais.

Dès les années 1985, il y avait une trentaine de systèmes de TA commerciaux au Japon, utilisables sur miniordinateur (AS-TRANSAC de Toshiba) ou poste de travail, ou PC (comme DUET-2 de SHARP), et cette activité était enseignée et pratiquée par des professionnels de la traduction. Ailleurs, les systèmes et leur usage étaient en retard d'un quinzaine d'années à cause de l'arrêt des financements provoqué par le rapport ALPAC (12/1966—1/1967). Le centre de recherche en TA de CMU ne fut ainsi créé qu'en 1985. C'est seulement en 1999 que, pour élaborer des directives appropriées et de la formation, les membres de l'*Association for Machine Translation in the Americas* (AMTA) et de l'*European Association for Machine Translation* (EAMT) créèrent un groupe d'intérêt pour la post-édition, le *Post-editing Interest Group*, ou *Post-editing SIG* (Allen, 1999).

Depuis les années 1990, avec le développement des systèmes de TA, la post-édition de traductions automatiques s'est développée dans le monde occidental, mais nettement plus lentement qu'au Japon. En fait, les gros bureaux de traducteurs professionnels ont longtemps préféré utiliser des systèmes à mémoires de traductions comme TRANSLATION MANAGER (TM2™) d'IBM, TRADOS, TRANSIT (de STAR), ou DEJA VU, même s'ils étaient couplés à des systèmes de TA, comme LMT pour TM2™. Selon (Boitet, 1996), la raison en est que les dictionnaires "main gauche" intégrés à ces systèmes et constamment mis à jour par les traducteurs suivaient l'évolution terminologique, alors que ceux des systèmes de TA évoluaient beaucoup plus lentement, de sorte que post-éditer des sorties de TA devenait plus long qu'utiliser les suggestions d'un système à mémoire de traductions (MT). Depuis que SYSTRAN a été mis sur minitel (en 1984), puis en service gratuit sur Internet (en 1994), ce sont les traducteurs indépendants qui se sont d'abord mis à utiliser des résultats de TA, en post-édition s'ils étaient assez bons, et comme "dictionnaire en contexte" sinon. Il y a quelques exceptions, dans des situations où les utilisateurs et les développeurs des systèmes de TA étaient en synergie. Ce fut le cas à la PAHO (*Pan American Health Organization*, Washington, et OMS, Genève, avec les systèmes SPANAM, ENGSPAN, puis PAHOMTS (Vasconcellos and León, 1985)), et aussi à Taipei (Hsu and Su, 1995).

L'autre cas où la post-édition de résultats de TA a depuis très longtemps été très appréciée, et utilisée par des traducteurs professionnels, est celui où le système de TA est spécialisé à un sous-langage précis. C'est le cas du système METEO au Canada (1981-2001)²⁰, dédié aux *bulletins*²¹ météorologiques, du système ALTFLASH de traduction des "brèves" du Nikkei (en mars 1998) (Uchino, Shirai et al., 2001), et surtout des "accélérateurs de traduction" (Pouliquen et al., 2013) construits par B. POULIQUEN à partir de mémoires de traductions contenant les bisegments produits par des traducteurs professionnels pendant plusieurs années (11 ans pour les systèmes construits avec Moses pour l'ONU à Genève).

Un nouvel usage de la post-édition est apparu depuis quelques années. Il s'agit non pas d'améliorer les traductions automatiques avant qu'elles ne soient "consommées", mais d'en post-éditer une partie pour augmenter et améliorer la MT à partir de laquelle un système de TA est construit. Ce sont donc des aides au développement et à l'amélioration de systèmes.

Nous présentons maintenant les systèmes destinés uniquement à aider les traducteurs, puis les systèmes destinés à améliorer des systèmes de TA.

²⁰ https://en.wikipedia.org/wiki/METEO_System .

²¹ et pas aux *situations* ni aux *alertes*.

1.1.3.2.2 Systèmes destinés uniquement à aider les traducteurs

Libellex. LIBELLEX²² est l'un des deux produits-phare de la société *Lingua et Machina* (L&M). Il s'agit d'une plate-forme de traitement multilingue d'aide à la traduction qui repose sur un système de gestion de mémoires de traductions. Il vise à constituer une mémoire collective d'entreprise pour tout ce qui concerne la communication multilingue. LIBELLEX intègre un module de traitement de la terminologie (import/export à partir de ressources du client, extraction de terminologie, gestion de terminologie, etc.), et un module pour créer et utiliser des systèmes de TA de type MOSES (alignement de corpus, entraînement, appel de systèmes de TA). Il peut être utilisé en mode SAAS²³ ou installé en mode licence sur des postes de travail informatiques.

SDL Trados Studio. La société TRADOS commercialise un outil de THAM appelé SDL TRADOS STUDIO, qui a environ 200 K utilisateurs (plus de 200 K licences du produit²⁴). Elle fournit des solutions innovantes en matière de logiciels de traduction à toute la chaîne logistique de traduction, y compris aux traducteurs indépendants, aux prestataires de services linguistiques, aux services internes de traduction et aux établissements universitaires. Cet outil se compose d'un gestionnaire de terminologie (SDL MULTITERM 2015), et d'un module d'alignement qui permet d'initialiser la MT utilisée à partir de documents déjà traduits.

1.1.3.2.3 Systèmes destinés à améliorer des systèmes de TA

MateCat. MATECAT (Schwenk et al., 2015) est un système visant à améliorer l'intégration d'un système de TA et de la traduction humaine (TH). L'objectif est d'améliorer la productivité des traducteurs professionnels et leur expérience du travail avec la TA. MATECAT propose une prétraduction par TA pour chaque phrase. Un système de TA basé sur MOSES est intégré dans MATECAT. Pour obtenir une bonne qualité de traduction et un gain de productivité, le système de TA est adapté au domaine (par exemple, la finance). MATECAT permet d'améliorer le système de TA en utilisant les résultats de post-édition, et il a intégré « Online model adaptation (Bertoldi et al., 2014) » de Moses pour améliorer dynamiquement le système de TA.

CASMACAT. CASMACAT (Alabau et al., 2013) (Cognitive Analysis and Statistical Methods for Advanced Computer Aided Translation) est un système pour la traduction assistée par ordinateur et l'étude scientifique de la traduction humaine. Pour l'amélioration des systèmes de TA (MOSES), CASMACAT propose trois méthodes efficaces en ligne pour mettre à jour les systèmes de TA : (1) Il infère de nouvelles règles à partir des segments post-édités, et les ajoute au modèle de traduction ; (2) il utilise les segments post-édité pour mettre à jour le modèle de langue ; et (3) il met à jour les paramètres discriminants de TA dans une étape MIRA (Chiang, 2012) et (Hasler et al., 2011).

1.1.3.3 Idées directrices pour faire évoluer SECTra_w/iMAG

1.1.3.3.1 Rendre SECTra configurable au niveau des corpus

Un SECTRA²⁵ configurable au niveau des corpus devrait permettre non seulement de définir le descripteur du corpus (le nom du corpus, les langues source/cible, le domaine, la taille, l'auteur, etc.), mais aussi de présenter la relation entre le corpus et les fichiers compagnons, de choisir le dictionnaire, de manipuler les corpus (fusion, séparation, comparaison, filtrage, etc.), le format d'import/export, le codage, etc.

²² <http://libellex.fr>

²³ https://fr.wikipedia.org/wiki/Logiciel_en_tant_que_service

²⁴ <http://www.translationzone.com/fr/about>

²⁵ Nom commun proposé dans la thèse de HUYNH Cong Phap (2010).

1.1.3.2 Trouver une modélisation

Nous souhaitons trouver une modélisation convenable pour les corpus multilingues, éventuellement complexes. Nous nous sommes inspiré de l'idée de structure de dictionnaire (les définitions sont présentées dans la thèse de M. MANGEOT (Mangeot, 2010)). Nous souhaitons construire une nouvelle version de SECTRA sur des corpus de traductions aussi variés et aussi gros que possible, comme JR-ACQUIS, EUROPARL, HANSARD, EOLSS augmenté par l'UNU-UNDL, ERIM, et surtout le corpus de brevets CLEF-IP (60 Go). Faire cela est nécessaire pour arriver à développer une méthode innovante et pertinente de modélisation des corpus multilingues alignés, et la partie du futur langage de programmation du système SECTRA qui permettra de décrire le type de chaque nouveau corpus (macrostructure, microstructure, comme pour les bases lexicales, et peut-être aussi "mésostructure" pour définir l'organisation des sous-documents, des fichiers satellites, et des documents ou fichiers compagnons).

1.1.3.3 Rendre SECTra programmable

Parmi les fonctions de SECTRA_W qu'on souhaite rendre programmables, il y a la synchronisation des données (en ligne ou hors ligne), l'import, l'export, la recherche, la sélection des données, et l'apprentissage incrémental de systèmes de TA à la MOSES.

Il faudrait que SECTra fournisse un support à la programmation, permettant de gérer et de traiter les corpus, et aussi des opérations motivées par l'utilisation dans des campagnes d'évaluation (par exemple, programmer certaines "mesures d'évaluation", comme BLEU, NIST ou TER), ou des opérations motivées par l'utilisation dans des projets de post-édition (par exemple, appel à un ou plusieurs systèmes de TA).

I.2 Améliorations de SECTra_w dans le cadre du projet Traouiero

I.2.1 Extension de fonctions existantes

Dans le cadre du projet TRAQUIERO, nous avons travaillé sur l'amélioration de SECTRA_W. Tout d'abord, nous avons amélioré la fonction d'appel de systèmes de TA. Elle permet de faire exécuter ou réexécuter la traduction automatique, par tel ou tel système, sur tout un corpus ou sur une sélection de documents.

Ensuite, nous avons refait la conception de la communication entre TRADOH et SECTRA_W, et réalisé la procédure d'appel des systèmes de TA via TRADOH. Enfin, nous avons participé à l'introduction de la possibilité pour un utilisateur d'ajouter à TRADOH des plugins permettant d'augmenter la collection des systèmes de TA appelables, et nous avons réalisé et intégré plusieurs plugins, dont un permettant d'appeler n'importe quel système MOSES.

1.2.1.1 Appel de systèmes de TA sur un corpus ou sur une sélection de documents

Pour ajouter de nouvelles traductions automatiques dans SECTRA_W, nous lui avons ajouté une nouvelle fonctionnalité. Dans la Figure 4, nous montrons comment choisir la langue cible et le système de TA dans l'interface de *Translate Corpus*. Une traduction via une langue pivot est aussi possible (cocher l'option *Pivot langage*).

seCTra_w

Home Import Evaluation Post-edition TM Admin Translation Contact us

Corpus name: demo2

All	MTs	Pivot language	Target language	Translation selection
<input type="checkbox"/>	<input type="checkbox"/> Google <input type="checkbox"/> Systran <input type="checkbox"/> MosesLIG	No pivot language	-select-	Start

	No	Document name	Source language	Pivot language	Target language	Operate
<input type="checkbox"/>	1	DOC1	french	No pivot language	-select-	Start
<input type="checkbox"/>	2	DOC2	english	No pivot language	-select-	Start
<input type="checkbox"/>	3	DOC3	english	No pivot language	-select-	Start
<input type="checkbox"/>	4	DOC4	english	No pivot language	-select-	Start
<input type="checkbox"/>	5	DOC5	english	No pivot language	-select-	Start
<input type="checkbox"/>	6	DOC6	english	No pivot language	-select-	Start
<input type="checkbox"/>	7	DOC7	french	No pivot language	-select-	Start
<input type="checkbox"/>	8	DOC8	english	No pivot language	-select-	Start
<input type="checkbox"/>	9	DOC9	english	No pivot language	-select-	Start
<input type="checkbox"/>	10	DOC10	chinese	No pivot language	-select-	Start
<input type="checkbox"/>	11	DOC11	french	No pivot language	-select-	Start
<input type="checkbox"/>	12	DOC12	french	No pivot language	-select-	Start
<input type="checkbox"/>	13	DOC13	french	No pivot language	-select-	Start
<input type="checkbox"/>	14	DOC14	french	No pivot language	-select-	Start
<input type="checkbox"/>	15	DOC16	french	No pivot language	-select-	Start

Figure 4 : Interface de « Translate Corpus »

I.2.1.2 Extension : utilisation de TRADOH et passage éventuel par une langue intermédiaire

Dans la conception de SECTRA_W, TRADOH est un important composant pour la gestion des appels à des systèmes de TA. Mais, en fait, il n’y avait pas en 2011 de vraie communication entre SECTRA_W et TRADOH. Dans le cadre du projet TRAQUIERO, nous avons reconstruit la fonctionnalité d’appel de systèmes de TA. Nous avons voulu pouvoir traiter les mêmes paires de langues avec GT et SYSTRAN, par défauts utilisés avec la même priorité pour produire les prétraductions.

Pour beaucoup de paires, GT passe en fait par un « pivot textuel » anglais. Mais d'autres systèmes, comme SYSTRAN, ne le font pas, car cela dégrade considérablement la qualité (linguistique et d'usage). Nous avons ajouté la possibilité pour un utilisateur de demander (via TRADOH) une traduction "double" passant par un texte dans une langue « pivot » qui peut être l'anglais ou toute autre langue (par exemple, le français pour faire de l'italien-espagnol). Cette possibilité est illustrée dans la Figure 4.

Dans la Figure 5, on voit une interface de TRADOH permettant de définir la langue pivot. Par exemple, le serveur SYSTRAN (V7) gracieusement prêté par SYSTRAN ne supporte pas la traduction du français vers le chinois. On prend l’anglais comme langue pivot pour obtenir la traduction en chinois (français→anglais, anglais→chinois).

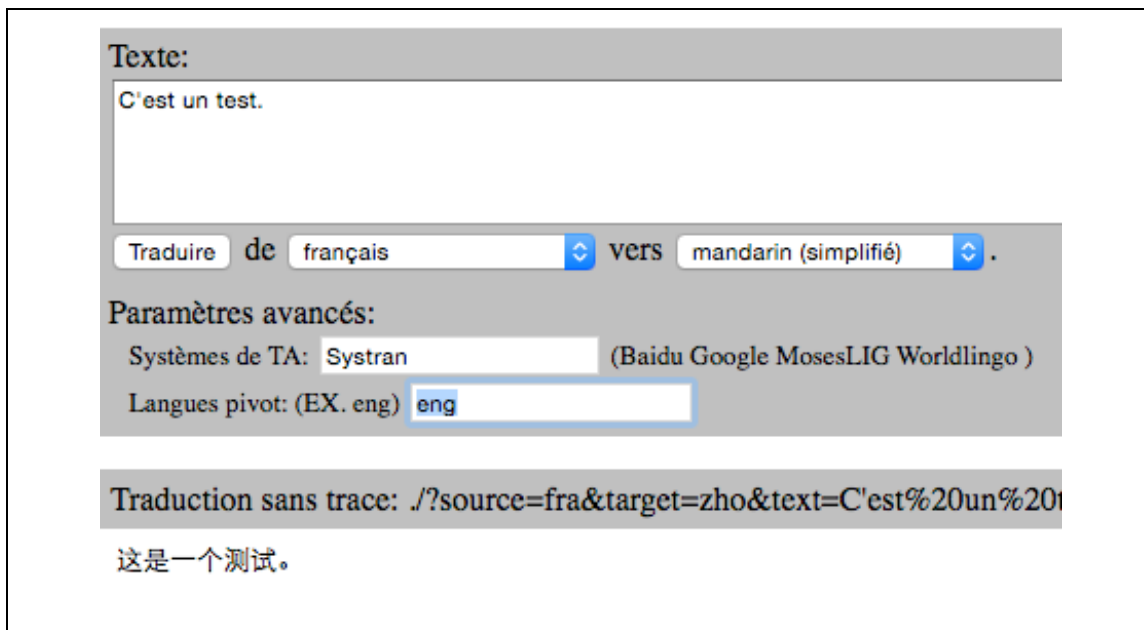


Figure 5 : Interface de TRADOH

I.2.1.3 Ajout par l'utilisateur de systèmes de TA utilisables pour les corpus associés à une MT donnée

Depuis 2004, de plus en plus d'environnements de création de systèmes de TA statistique en source ouvert (comme APERTIUM (Forcada et al., 2011), MOSES (Koehn, Hoang et al., 2007), JOSHUA (Ganitkevitch et al., 2012), etc.) ont été développés par des laboratoires ou des associations. Les utilisateurs utilisent ces boîtes à outils pour entraîner leurs systèmes de TA sur leurs corpus parallèles.

Pour mettre en œuvre un système de TA qu'il a créé, un utilisateur peut utiliser le plugin de TRADOH que nous avons créé pour appeler des systèmes de TA de type MOSES. Nous avons aussi créé un plugin permettant d'appeler les systèmes de TA de type Ariane-G5 (Guilbaud, 1999) disponibles sur la plate-forme HELOÏSE (Berment and Boitet, 2012).

Pour l'instant, nous avons utilisé notre plugin Tradoh-Moses pour activer dans SECTRA_W le système de TA anglais → hindi MOSES-CFILT-EN_HI de l'IITB (Bombay) et nos systèmes MOSESLIG français ↔ chinois.

Dans la Figure 4, nous montrons comment nous avons ajouté MOSESLIG pour appeler nos systèmes de TA français ↔ chinois. Quand l'utilisateur veut traduire un document ou un corpus du français vers le chinois (ou l'inverse), il coche la case MOSESLIG, et le système de TA français ↔ chinois sera appelé.

I.2.2 Aspects de génie logiciel

I.2.2.1 Sélections paramétrables

Dans le cadre du projet TRAQUIERO, nous avons travaillé sur la "programmabilité" de SECTRA_W. Nous avons ainsi introduit et implémenté les *sélections paramétrables*. Le but est, par exemple, de sélectionner une « bonne » partie d'une MT et de l'exporter pour construire un système MOSES ou pour faire de l'amélioration incrémentale (voir Chapitre VI).

Nous avons ajouté la fonction « *export* » dans SECTRA_W pour exporter les segments post-édités. Nous avons implémenté un type de prédicat de sélection assez fréquent, paramétrable dans l'interface présentée dans la Figure 6.

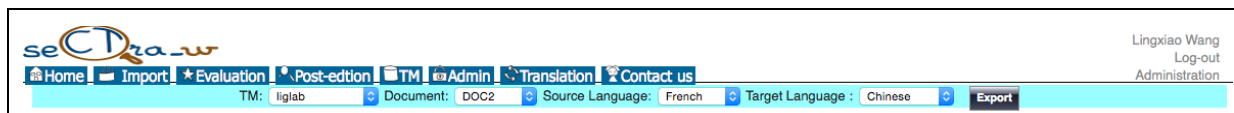


Figure 6 : Options de la fonction "export"

Par exemple, on peut utiliser l'expression de sélection suivante :

(\$niveau=3 et \$score>=14) ou (\$niveau=4 et \$score>=12) ou (\$niveau=5 et \$score>=11)

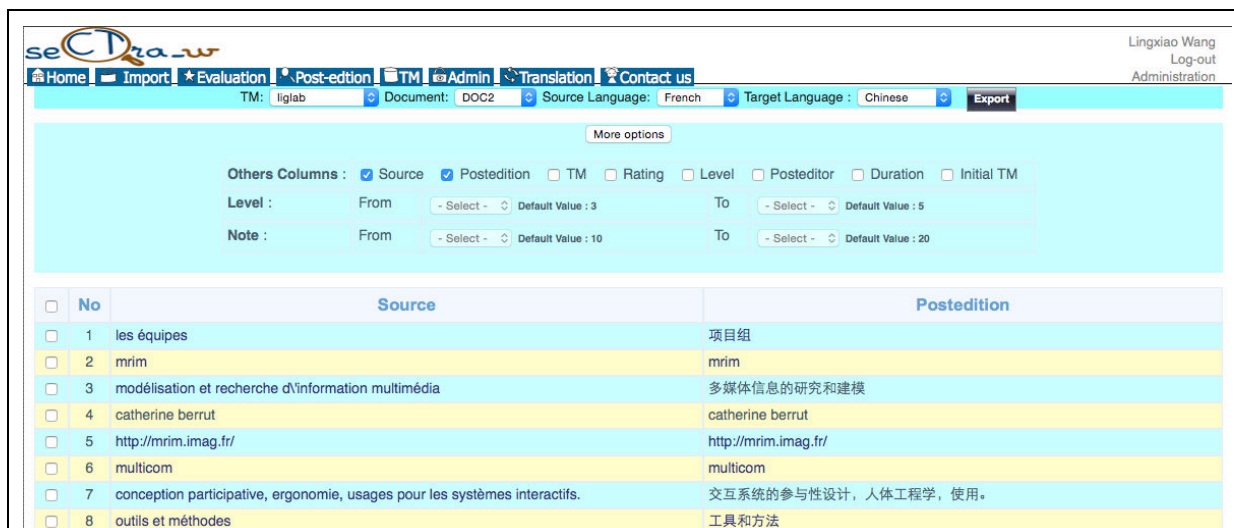


Figure 7 : Interface de sélection paramétrable dans SECTra_w

Pour aller plus loin, il faut définir ici un « petit morceau de langage » permettant de programmer n'importe quelle expression de sélection. Pour cela, nous avons enrichi l'API de SECTRA_W pour pouvoir maintenir et publier des « variables SECTRA » comme \$niveau, \$score, \$source, \$MT, \$posteditor, \$length, \$detected_language, \$date(s), etc. (Figure 7).

1.2.2.2 Clarification du traitement des "occurrences" des segments

Une page Web, qui évolue dans le temps, est représentée dans SECTRA_W par un « pseudo-document » qui contient des références à tous les « segments » déjà vus dans cette page, ainsi que le « squelette » de l'instance courante de la page.

Le problème vient surtout du fait qu'un segment peut avoir été trouvé auparavant dans une autre page Web, correspondant à un autre pseudo-document, et que, si on ne le trouve pas dans la MT du pseudo-document courant, il faut (on croit qu'il faut) aller le chercher dans les MT des autres pseudo-documents. Cela conduit à N recherches s'il y a N pseudo-documents et si le segment est nouveau... Or on devrait trouver la réponse par une seule recherche dans la MT globale.

À cause du problème algorithmique lié à ce mauvais choix de structure de données, il y a en fait une petite MT par pseudo-document (cette MT contient les segment source, les traductions obtenues par TA, et les post-éditions du document), et pas seulement un ensemble ou une liste des segments déjà vus.

Nous avons proposé une nouvelle structure de MT partagée pour résoudre ce problème, sans détruire la structure existant de SECTRA_W. Cette MT contient les segments validés par les modérateurs. Quand un nouveau segment arrive, SECTRA_W cherche d'abord la traduction dans la MT du pseudo-document en cours. S'il n'en trouve pas, il le cherche dans la MT

partagée. S'il le trouve, la traduction sera affichée sur la page Web traduite, sinon, il appellera le système de TA.

I.2.2.3 Analyse de fonctions dans XWiki, améliorant la modularité

SECTRA_W a été construit sur XWiki. Sa programmation utilise donc les langages associés (GROOVY²⁶, VELOCITY²⁷, JAVASCRIPT, HTML). La version de XWiki utilisée par SECTRA_W est la version 1.3, alors que la version actuelle est la 7.2. Elle nécessite l'installation de TOMCAT 5 dont on ne trouve plus d'installateur. Il a fallu procéder à l'installation complète « à la main » de TOMCAT 5. Des notes résument les manipulations effectuées, avec quelques explications, pour pouvoir les comprendre et les reproduire ou les adapter aux situations à venir.

L'installation de la version de XWiki sur laquelle fonctionne le service a donné aussi lieu à la rédaction d'une note détaillée, en particulier pour bien mémoriser ce qui doit être fait pour définir correctement les connecteurs entre le serveur TOMCAT et les services déployés par XWiki. Comme l'utilisation de XWiki nous paraît en fait peu adaptée à l'implémentation de logiciels comme SECTRA_W et IMAG, même si on passait à la version 7.2, et comme un tel passage serait très lourd, nous avons préféré (et spécifié lors du projet TRAQUIERO) nous orienter vers une réimplémentation plus modulaire, et indépendante de XWiki.

Pour faire évoluer SECTRA_W, nous avons donc commencé à travailler sur l'augmentation de sa modularité. Tout d'abord, nous avons fait une liste pour le code GROOVY qui est dans les pages XWiki (car ce code est sauvé dans des pages Web), puis nous avons analysé les fonctions contenues dans ce code. Enfin, nous avons classé les fonctions.

I.2.3 Travail de spécification

Dans le cadre du projet TRAQUIERO, nous avons fait une spécification pour améliorer la structure de la base de données de SECTRA_W. La nouvelle base de données de SECTRA_W est centrée sur le métier « gérer des corpus multilingues de textes multilingues alignés ». Les 5 entités principales identifiées et implémentées dans le modèle conceptuel sont listées ci-dessous et on leur a rajouté 2 entités (*pseudodocument* et *segment multilingualisé multilingue*), qui manquaient dans la version courante de SECTra_w, ainsi que des précisions sur la dernière entité (Traduction). Si l'on veut avoir une représentation en « réseau corporal », comme on le fait pour des « réseaux lexicaux », il faudra donc spécialiser la notion de « segment tout court », inutilisable telle quelle, en « segment source » et « segment cible » ou « traduction segmentale ».

On pourra alors considérer qu'un segment source est un nœud d'un graphe, relié par une relation (directe ou indirecte) de traduction à des segments cibles, et que le sous-graphe connexe obtenu forme exactement un *segment multilingualisé contextualisé*. Du point de vue des traducteurs et des représentations « à la TMX », il faut cependant considérer un segment multilingualisé contextualisé comme un objet identifié comme tel, à structure complexe, et muni d'un historique propre.

La structure logique de la nouvelle base de données de SECTra_w illustrée dans la Figure 8 ne prend en considération l'utilisateur que dans la table provenant de l'association plusieurs à plusieurs entre les segments et leurs contextes. La stratégie de programmation choisie pour développer SECTra_w prévoit des modifications, en particulier en ce qui concerne les colonnes de ses tables²⁸.

²⁶ <http://www.groovy-lang.org>

²⁷ <http://velocity.apache.org>

²⁸ Citation d'un rapport de TRAQUIERO, écrit par Valérie Bellynck.

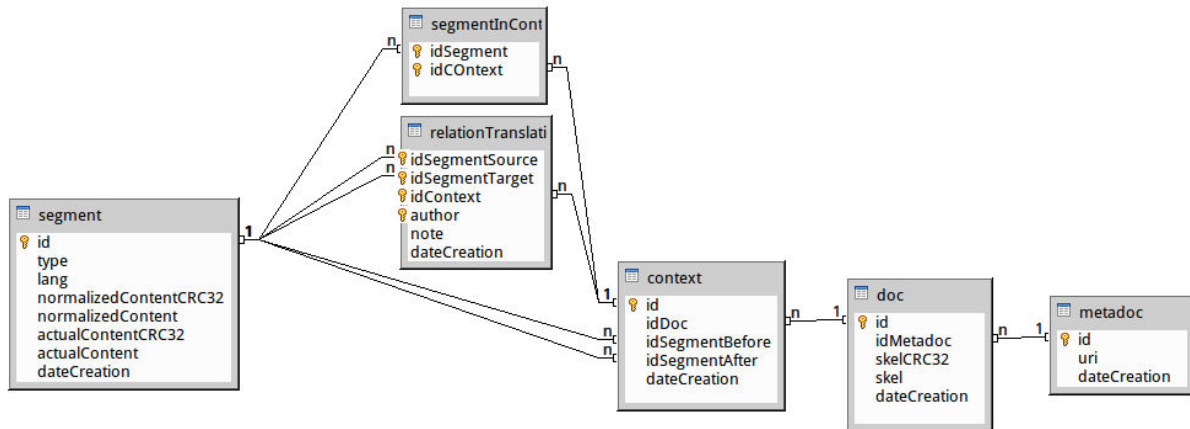


Figure 8 : Structure logique d'une base de données de corpus multilingues

I.3 Améliorations des iMAG dans le cadre du projet Traouiero

I.3.1 Paramétrisation

I.3.1.1 Réalisation d'un configurateur d'iMAG par C. P. HUYNH

Dans le cadre de la thèse de C.P. HUYNH, on a construit SECTRA_W pour permettre de mettre en œuvre des paramétrages « externes » au niveau du paramétrage et de la configuration par les non-informaticiens :

- Optimisation de l'interface d'évaluation subjective selon les critères d'évaluation d'une campagne, et le réglage de l'interface (nombre de segments, de colonnes, de langues, etc.).
- Filtrage de données post-éditées selon la date, l'auteur, la qualité (niveau traductionnel, note de qualité), etc.
- Recherche et remplacement interactif dans une sélection arbitraire (ensemble de segments, ensemble de documents).
- Récupération de données à partir de MT en fonction de divers paramètres.
- Lancement automatique des appels de systèmes extérieurs (TRADOH++, PIVAX, etc.).

Au niveau du paramétrage et la programmation par des contributeurs informaticiens, SECTRA_W permet d'écrire facilement des scripts pour le traitement de données, tels que des scripts pour l'import et l'export de corpus. Par exemple, A. FALAISE a écrit un script permettant l'appel depuis l'extérieur de la fonction d'import de SECTRA_W pour importer un corpus du projet SURVITRA et un autre du projet OMNIA. Cependant, nous n'avons pas encore développé les langages (langage de conditions booléennes, langage de flots de travail, langage narratif) dans la version actuelle de SECTRA_W. Nous prévoyons de les intégrer dans la prochaine version.

I.3.1.2 Améliorations prévues mais pas encore réalisées

Il faudrait améliorer le logiciel iMAG pour que, en fonction des informations sur le site S et des paramètres de l'accès multilingue à S, il appelle l'iMAG-S avec les paramètres adéquats. Ces paramètres concernent surtout les langues (lesquelles sont post-éditables, modérables, post-éditées, modérées, ou exclues...) et les systèmes de TA (lesquels sont souhaités, avec quels paramètres, avec quelles préférences entre systèmes pour chaque couple de langues, avec quels chemins traductionnels pour chaque couple, etc.).

Sur l'interface, le paramétrage d'une iMAG doit pouvoir permettre à l'administrateur d'un site Web « élu » de personnaliser la présentation (message de bienvenue, feuille de style), et de choisir les interacteurs qu'il veut proposer à ses visiteurs, selon leur rôle. Par exemple, la feuille de style pourrait permettre de générer une présentation « à la NEON », avec le texte source écrit au-dessus du texte cible, ligne par ligne, et la palette pourrait contenir les noms des contributeurs à chaque segment, et/ou son origine (niveau de ☆ à ☆☆☆☆☆), et/ou sa note (de 0 à 20), etc.

I.3.2 Travail de spécification

I.3.2.1 iMAG-Relais

Pour mieux gérer les rôles des utilisateurs entre composants différents, Valérie Bellynck a proposé un nouveau composant dit « RELAIS », qui aura pour rôle d'assurer l'identification vis-à-vis des différents services, en octroyant aux utilisateurs, de façon transparente, les droits suffisants pour assumer les tâches associées à leur rôle dans chacun d'eux, en fonction des paramétrages choisis par l'administrateur référent d'une passerelle iMAG, ou par ses administrateurs délégués. Nous reprenons ici une partie de son texte.

Lorsque le problème du transfert des droits a été identifié, le rôle du relais devait être de relayer les droits entre les différents services utilisés dans une architecture multiservices dans lesquels les utilisateurs peuvent avoir des rôles et des droits différents selon le service exploité. Il s'agit de transmettre ces droits sans introduire de faille de sécurité dans les services exploités.

Non seulement les types des rôles peuvent être différents [quelqu'un peut être client dans un site commercial de souscription d'offres de services, et modérateur dans un site contributif visant à construire collaborativement de la connaissance], mais les termes utilisés pour identifier les rôles peuvent être identiques tout en ne correspondant pas à des accès identiques aux mêmes fonctions.

I.3.2.2 Tableau blanc

Dans l'implémentation actuelle des iMAG, la plupart des composants sont des services Web, utilisés comme des serveurs. On constate deux inconvénients : (1) ces services ne sont pas autonomes (il faut les relancer à la main ou indirectement par des scripts), et (2) on n'a pas de vision sur ce qui se passe dans le système global, et qui peut être compliqué. En effet,

- il y a des appels mutuels. Par exemple, SECTRA_W appelle PIVAX en lui envoyant des segments, pour en recevoir des minidictionnaires en HTML, et PIVAX peut appeler SECTRA_W pour en obtenir la prétraduction et la post-édition des exemples d'usage contenus dans des articles de dictionnaires.
- il y a des exécutions de tâches *en boucle infinie*, qui sont pour l'instant limitées à des files d'attente de tâches considérées comme élémentaires²⁹. Par exemple, un corpus peut être en train d'être traduit de français en allemand, un document EOLSS peut être en train d'être *déconverti* depuis UNL³⁰ vers le français ou le russe (par appel à des serveurs de TA, via TRADOH++), et plusieurs documents et mémoires de traductions peuvent être en cours de post-édition ou d'évaluation humaine.

Donc le but de cette tâche est de :

²⁹ D'après C. BOITET, « Il s'agit de précalculer des traductions automatiques et des minidictionnaires de façon incrémentale, et une infinité *équitable* de fois si on avait l'éternité devant soi, de façon à mettre à jour ces résultats, sachant que les systèmes de TA et les ressources lexicales évoluent constamment. »

³⁰ Les documents EOLSS exploités sous SECTRA_W ont des fichiers compagnons UNL qui associent un graphe UNL à chaque segment, et ces graphes sont stockés comme des annotations d'un type particulier dans SECTRA_W.

- construire, comme cela a déjà été envisagé dans le projet (ANR) OMNIA, un dispositif de contrôle du système complet par *tableau blanc* et *chef d'orchestre*³¹ ;
- organiser les modules principaux comme des *agents à gros grain*, capables d'initiative, et en particulier capables de se relancer seuls, et de déterminer eux-mêmes ce qu'ils ont à faire en faisant appel au tableau blanc.

I.3.2.3 Architecture en agents à gros grains

Lorsque le système à modéliser est trop complexe, une analyse globale révèle rapidement ses limites. Il est souvent plus simple et rapide de chercher des solutions à des problèmes plus locaux, selon le principe classique en informatique du « *diviser pour régner* ». Par exemple, pour réguler le trafic aérien, de trop nombreux paramètres et contraintes influent sur le résultat de l'ensemble pour qu'une approche globale permette une résolution fiable.

En revanche, des solutions locales aident à résoudre élégamment et efficacement de telles difficultés. Les résultats obtenus sont souvent acceptables, bien qu'il soit au mieux difficile et au pire impossible de prouver cela théoriquement. L'approche multi-agents permet d'obtenir plus simplement des résultats sur des problèmes complexes, en contournant les difficultés qu'auraient des algorithmes globaux pour les gérer.

I.3.2.4 Files d'attente

Dans les cas où plusieurs clients (SECTRA_W/IMAG, TRADOH, PIVAX, etc.) veulent demander un traitement particulier, la communication par messages est souvent intéressante. En effet, elle permet de découpler les demandes des clients (comme l'appel à des systèmes de TA), et les traitements par les serveurs. Dans certains cas, les traitements à effectuer sont plus importants que le simple calcul d'une page. Il est alors beaucoup plus intéressant de passer par une liste de tâches partagées par les clients et les serveurs.

L'idée d'utiliser un système général de gestion de files d'attente sur le Web n'a finalement pas pu être réalisée lors du projet TRAQUIRO, faite de temps. Nous l'avons reprise, dans le cadre de cette thèse et de celle de Y. ZHANG. Nous en parlons plus en détail vers la fin de ce mémoire (Voir IX.1.3.4). Disons seulement ici que nous utilisons ACTIVEMQ³² de la fondation APACHE.

³¹ Cette idée a été proposée et expérimentée il y a déjà longtemps par Ch. Boitet et M. Seligman (Boitet, Blanchon et al., 1994) Boitet, C., Blanchon, H., Seligman, M. and Bellynck, V. (1994). Evolution of MT with the Web. Proc. *Conference "Machine Translation 25 Years On"*. pp. 1-13., pour la TA de parole hétérogène, et est devenue récemment un axe de recherche reconnu, avec par exemple un atelier à LREC-2010 sur les *flux de traitements linguistiques*.

³² <http://activemq.apache.org/>

Chapitre II Travail sur des aspects plus conceptuels et définition de nouvelles fonctionnalités

Résumé

Nous présentons dans ce chapitre quelques réflexions sur la modélisation des corpus de traductions, dont la structure est bien plus complexe que celle des "corpus parallèles" simples (comme le BTEC) utilisés en TA statique pour l'apprentissage, et qui sont en fait des MT. Le but poursuivi est, à terme, de séparer les concepts de MT et de corpus. À partir de cette analyse et des travaux plus pratiques décrits plus haut, nous avons aussi travaillé sur la conception de nouvelles fonctionnalités.

II.1 Modélisation de corpus de traductions variés

II.1.1 Variété des corpus de traductions et esquisse d'une méthode de formalisation

Les notions de corpus et de corpus de traductions ont été définies plus haut (Voir I.1.1.2.2). Dans un corpus de traductions, si l'on peut faire correspondre chaque unité du texte en langue source (l_1) avec chaque unité de texte en langue cible (l_2) au niveau des paragraphes, phrases ou mots, ce corpus de traductions est considéré comme un *corpus parallèle*.

On a aussi introduit plus haut la notion de *corpus comparable* (Déjean and Gaussier, 2002). Deux corpus comparables, dans deux langues l_1 et l_2 , "parlent de la même chose", mais ne sont pas des corpus de traductions. Au mieux, on peut y trouver des phrases en relation de traduction, et déterminer le sens de traduction. Comme l'a montré E. DELPECH dans sa thèse (Delpech, 2013), l'espoir selon lequel on pourrait les utiliser pour trouver efficacement de la terminologie bilingue a été largement déçu. C'est pourquoi nous ne nous intéressons qu'aux corpus de traductions.

II.1.1.1 Variété des corpus de traductions

La première chose à prendre en compte si l'on veut traiter les corpus de traductions est leur grande variété, tant au niveau des contenus que des formats, des structures et des usages.

II.1.1.1.1 Variété des contenus

Texte. Un *corpus écrit* est un corpus de textes, s'il contient essentiellement des textes. Il peut être constitué de documents différents (tableau, extraits de textes, etc.). Il peut être présenté comme un livre, une page Web, ou un document (en WORD, ODS, LATEX, INTERLEAF, etc.).

Pour chaque document source, les textes et les traductions des segments peuvent être sauvegardés dans des fichiers parallèles séparés, un par langue, ou bien à l'intérieur d'un même fichier, la représentation de chaque segment contenant alors le texte source et sa ou ses traductions. Les corpus BTEC, EUROPARL, JR-ACQUIS, MULTIUN sont formés de fichiers parallèles, alors que le corpus de brevets CLEF-IP est formé de fichiers multilingues au niveau de chaque segment.

Parole. Un *corpus de parole* est un corpus constitué d'enregistrements de données orales, et éventuellement de leurs transcriptions, qui peuvent être considérées comme des annotations. À la différence de l'écrit qui n'utilise qu'un seul support, l'oral associe le plus souvent la parole enregistrée à une représentation écrite et codée (transcriptions, traductions, annotations).

Un corpus de dialogues bilingues interprétés comme le corpus ERIM contient non seulement les tours de parole des participants et leurs traductions (orales) produites par un interprète,

mais aussi des tours de parole monolingues, qui forment des sous-dialogues de clarification entre l'interprète et un participant. Un système comme SECTra doit permettre de "rejouer" ce genre de corpus, et de créer les transcriptions écrites des segments. Dans SECTRA_W, cette possibilité existe, à un niveau très « rustique », et il est possible de "déléguer" cette activité à des environnements plus adaptés (par exemple, contenant des reconnaisseurs de parole et des aligneurs texte-signal).

Visuel. Un corpus visuel est un corpus contenant des vidéos ou images, et ce type de corpus peut contenir les textes pour compléter, décrire, ou présenter les objets. Un film en vidéo contient d'ailleurs la plupart du temps, sur une ou plusieurs pistes, des sous-titres en plusieurs langues.

II.1.1.1.2 Variété des formats

À cause de la variété des usages, il existe plusieurs formats de corpus. Le format le plus utilisé est le texte brut (.txt). Pour un corpus monolingue, les textes sont sauvegardés dans un ou plusieurs fichiers .txt monolingues.

Pour un corpus bilingue/multilingue, ce format manque de balises indiquant les relations entre les phrases, et leurs langues. Donc, d'habitude, les phrases sont sauvegardées dans des fichiers séparés parallèles, un par langue. Les corpus de traductions au format .txt sont très utilisés pour l'entraînement de systèmes de TA comme MOSES. Le texte source et le texte cible sont alignés (un paragraphe par segment), et sauvegardés dans $N+1$ fichiers s'il y a N langues cibles.

L'autre type de format le plus utilisé est « le format XML », ou plus exactement des instances de XML, comme TMX³³, TEI³⁴, DocBook³⁵, NITF³⁶, etc., correspondant à des DTD ou à des schémas XML définis par les professionnels du traitement des documents, de la traduction, et de la localisation.

Dans le cas d'un corpus de dialogues oraux, les enregistrements sonores sont sauvegardés dans des fichiers de format audio comme .wav, .mp3, .flac, etc., et chaque tour de parole est représenté par un fichier XML contenant les métadonnées associées, parmi lesquelles l'url du fichier audio.

De même, un corpus "visuel", dont les documents sont des vidéos ou des textes contenant des vidéos, contient les vidéos dans des fichiers .avi, .mkv, .mp4, etc, et des images dans des fichiers .svg, .jpeg, .png, etc. Dans un film en vidéo, un sous-titre constituant un segments et ses traductions peut être considéré comme d'un segment multilingualisé.

II.1.1.1.3 Variété des structures

Structure simple. La structure la plus simple est celle où un document est une suite de segments, en format texte, avec uniquement un identificateur par segments par segment. Les traductions dans N langues sont alors contenues dans N fichiers de même structure.

Corpus avec annotations dans des fichiers compagnons. Dans ce type de structure, chaque fichier (source ou traduction) peut avoir un ou plusieurs fichiers « compagnons » contenant des annotations des « externes ». Cette structure peut conserver le corpus original, et ajouter l'information pour la description du corpus. Par exemple, le corpus ERIM contient dialogues

³³ https://en.wikipedia.org/wiki/Translation_Memory_eXchange

³⁴ <http://www.tei-c.org/index.xml>

³⁵ <http://www.docbook.org>

³⁶ <http://www.gwg.nga.mil/ntb/index.html>

bilingues oraux interprétés. Les dialogues sont sauvegardés dans les fichiers *.wav*, et chaque fichier *.wav* est accompagné par un fichier *.xml* qui décrit le contenu du fichier *.wav*.

Corpus avec des annotations internes. Les mots et/ou les segments sont étiquetés par des balises XML, ou par des chaînes spéciales (ex : *Nous_Pr1p1*, *avions_VbIndImp1p1*, etc.). Les étiquettes peuvent fournir des informations de divers ordres, comme catégories syntaxiques, lemmes (forme canonique du mot fléchi), âge et sexe du locuteur, niveau d'études, etc.

Corpus arboré. Il s'agit de corpus parsés, contenant des informations sur la structure des phrases (corpus EOLSS, par exemple, voir II.1.3.3).

II.1.1.1.4 Différents usages

Les corpus de traductions différents aussi selon ce qui on veut en faire.

Corpus de traductions pour la TA

- Les corpus d'entraînement sont utilisés pour entraîner un système de TA MOSES (EUROPARL, MULTIUN, etc).
- Les corpus de développement sont utilisés pour modifier des poids dans l'étape de réglage (*Tuning*) en TA statique
- Les corpus de test sont utilisés pour évaluer les systèmes de TA et calculer des scores comme BLEU, NIST.
- D'autres corpus sont utilisés pour l'alignement, l'extraction de terminologie bilingue, etc.

Corpus de traductions pour apprenants. Un tel corpus contient des productions écrites et/ou orales faites par des apprenants d'une langue seconde servent à décrire l'interlangue et donc les difficultés des apprenants servent aussi à élaborer une typologie des erreurs pour l'utilisation dans un système de vérification grammaticale, ou production d'éditions bilingues annotées d'ouvrages littéraires.

II.1.1.2 Exemples

II.1.1.2.1 Corpus d'évaluation du projet TRANSAT

En 2005-2007, le GETA a participé au projet TRANSAT dans le cadre d'un contrat avec France Telecom R&D. À cette occasion, C. P. HUYNH a créé la première version de SECTRA_W, et nous l'avons utilisée comme support d'une campagne d'évaluation ayant pour but d'évaluer l'utilité des systèmes de traduction automatique commerciaux dans le domaine de la traduction de la parole.

Dans cette campagne, nous avons utilisé des données issues du corpus BTEC pour construire deux corpus sur deux types de tâches relatives au tourisme, ainsi que des dialogues collectés par l'équipe MULTICOM du CLIPS (Blanchon et al., 1999). Ce premier corpus source (anglais) a été extrait du corpus BTEC, et contenait un ensemble de 2224 tours de parole tous différents, en accordant une priorité aux tours de parole relatifs aux problèmes de santé. Le second corpus source (anglais), qui a aussi été extrait du corpus BTEC, était constitué de 2000 tours de parole dans le domaine de la restauration.

Les traductions candidates étaient les traductions de REVERSO TRANSLATOR 10 et de SYSTRAN (version 4). Les traductions de référence ont été produites par la post-édition de ces traductions candidates. La composition du corpus est détaillée dans l'Annexe 1.

Dans cette campagne d'évaluation, le GETA avait proposé un protocole d'évaluation pour le projet TRANSAT, donc ce corpus contient le résultat de mesures, comme la distance d'édition (Damerau, 1964, Levenshtein, 1966, Wagner and Fischer, 1974) en mots/caractères/pondérée,

BLEU, NIST, et aussi des mesures subjectives (de fluidité et d'adéquation). (Voir la définition des mesures dans l'Annexe 2).

II.1.1.2.2 Corpus B@BEL UNESCO

Le corpus UNESCO/B@BEL (Boitet, Boguslavskij et al., 2007) a été créé dans le cadre d'un contrat de recherche avec l'UNESCO. Nous avons traduit le site Web B@BEL de l'anglais vers 4 langues cible (français, russe, espagnol, et chinois). Les prétraductions ont été produites par SYSTRAN v5 PREMIUM (anglais→chinois/français/espagnol) et ETAP3³⁷ (anglais→russe). De plus des graphes UNL ont été construits semi-automatiquement.

Ce corpus contient environ 43200 mots (173 pages standard), et se présente comme un ensemble de multisegments, qui sont sauvegardés dans les fichiers texte, un par langue. Chaque segment source est aligné avec les 4 autres segments cible. Le corpus B@BEL contient 906 graphes UNL, et un échantillon contient 50 exemples (environ 50 pages standard). Un exemple tiré du corpus B@BEL est donné dans l'Annexe 3.

II.1.1.2.3 Corpus EOLSS

Le projet EOLSS/UNESCOL a été réalisé de février 2008 à octobre 2008 dans le cadre d'un contrat entre L'ASSOCIATION CHAMPOLLION et la fondation UNDL³⁸, visant à la traduction de l'anglais vers les 5 autres langues officielles de l'UNESCO (*français, espagnol, russe, chinois, arabe*) de 25 articles de l'EOLSS (*Encyclopedia of Life Support Systems*).

Le corpus EOLSS contient 25 documents sur l'eau et l'écologie de l'encyclopédie EOLSS, représentant environ 220K mots (13676 segments) ou 880 pages standard. Chaque document est constitué d'un fichier HTML (*.aspx*), d'un fichier compagnon *.unl*, et de fichiers satellites (images, icônes, et autres hors-texte) (Figure 9). Le fichier *.unl* contient des graphes UNL représentant des segments découpés à partir du fichier *.aspx* correspondant. Le lien de téléchargement du corpus EOLSS est dans l'Annexe 4.

³⁷ Institute For Information Transmission Problems (Kharkevich Institute) : <http://cl.iitp.ru/etap3>

³⁸ The UNDL FOUNDATION is a private Swiss law Foundation with head office in Geneva, Switzerland, legally representing the United Nations Organization in protecting its property rights pertaining to the UNL language and system, and legally representing the EOLSS Publishers and the UNESCO Joint Committee in translating the EOLSS with the use of the UNL technology.

Présentation de fichier HTML	Fichier compagnon .unl
<p>The wave propagates from the source with a velocity of long gravity water waves in accordance with the equation</p> $C_G = (g H)^{1/2}, \quad (1)$ <p>where g is the acceleration due to gravity, and H is the depth of the basin. Because the average depth of the world ocean is 4 km, the typical velocity of tsunami in the ocean is 200 m s-1 or 720 km h-1.</p> <p>Such a wave, propagating with the velocity of an airplane, may traverse the Pacific ocean in 10-12 hours and bring down a wall of water 10 m high with a velocity of more than 70 km h-1 upon a calm ocean beach.</p>	<p>[S:44] {org} The ... equation {/org} {unl} ... {/unl} ;graphe unl représentant la phrase [S]</p> <p>[S:45] {org} where ... basin. {/org} {unl} ... {/unl} [S]</p> <p>[S:46] {org} Because ... is 200 HTM1 or 720 HTM2. {/org} {unl} ... {/unl} [S]</p> <p>[S:47] {org} Such a wave, ... 70 HTM1 ... ocean beach. {/org} {unl} ... {/unl} [S]</p>

Figure 9 : Fichier HTML et fichier compagnon .unl

Pour adapter les documents à SECTRA_W, nous avons changé *.aspx* en *.html*. Nous avons utilisé le système de TA SYSTRAN pour produire les prétraductions, puis nous avons post-édité les sorties de TA. Dans la Figure 10, on présente un document (*GROUND AND SOIL WATER CHARACTERISTICS*) parallèle anglais-français. Après la post-édition du document, nous avons créé la page Web traduite en langue cible, et un fichier *.unl* en langue cible.

Summary	Résumé
<p>Water bedded under the earth's surface in the crust is called ground water. Geological structures holding groundwater are located in different physico-chemical conditions—this determines the different aggregative states of water: liquid, solid, and gaseous</p> <p>Liquid water is located in the upper part of the earth's crust under relatively low temperature and pressure. Water deeper than 20-70 km is in gaseous state under high pressure. Water in the solid state is located in the zone of permafrost.</p>	<p>L'eau sous la surface terrestre dans la croûte s'appelle la nappe phréatique. Les structures géologiques contenant des eaux souterraines sont situées dans des contextes physico-chimiques différents, ce qui détermine les différents états d'agrégation de l'eau : liquide, solide, et gazeux.</p> <p>L'eau liquide est située dans la partie supérieure de la croûte terrestre sous une température relativement basse et une pression. L'eau à plus de 20-70 kilomètres de profondeur se trouve à l'état gazeux sous haute pression. L'eau à l'état solide se trouve dans la zone de permafrost.</p>

Figure 10 : Document 2 traduit de l'anglais vers le français (*GROUND AND SOIL WATER CHARACTERISTICS*)

II.1.1.2.4 Corpus ERIM

Le corpus ERIM est un corpus de dialogues bilingues (français↔chinois, vietnamien, hindi, tamil) oraux interprétés. Un segment est constitué d'un tour de parole, et éventuellement d'un texte descriptif et d'un texte transcrit.

Voici un type de descripteur d'un dialogue français↔vietnamien du corpus ERIM.

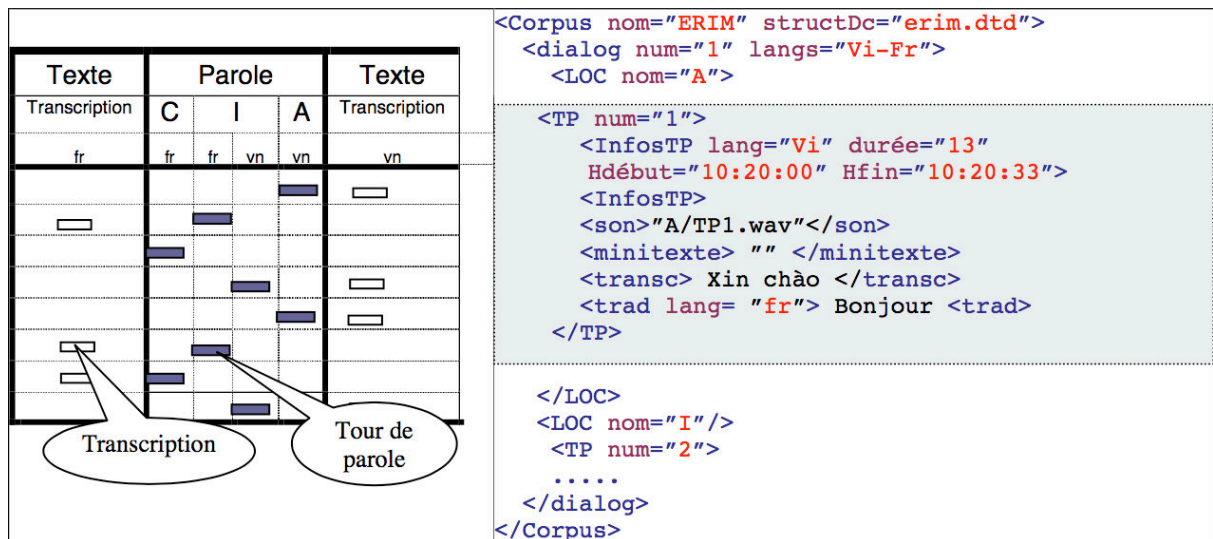


Figure 11 : Exemple de structure et de description d'un dialogue du corpus ERIM

La partie français↔vietnamien du corpus ERIM contient 4 séances de dialogue. Chaque séance contient 3 répertoires, 3 fichiers *.wpl* (par exemple, *french.wpl*, *vietamese.wpl*, et *dialogue.wpl*) et un fichier de transcription (comme le fichier *xml* dans la Figure 11). Chaque répertoire est nommé par le nom du locuteur (par exemple, *Cantona*, *Tung*, et *Interp*) ; il contient les tours de parole et les descripteurs. Les fichiers *.wpl* enregistrent les chemins des fichiers son. Par exemple, le fichier *french.wpl* (ou *vietamese.wpl*) enregistre les chemins des tours de parole en français (ou en vietnamien) émis/reçus par le locuteur français (ou vietnamien) et l'interprète (Figure 12). Le fichier *dialog.wpl* enregistre tous les tours de parole par ordre chronologique.



Figure 12 : Exemple du fichier *french.wpl* et *vietamese.wpl*

Le répertoire du locuteur contient les fichiers *.wav* et 3 descripteurs (3 fichiers *.xml* : *DSD.xml*, *DSL.xml*, *TIER.xml*).

- ***DSD.xml*** (*Description Session Dialogue*), enregistre l'information de la session de dialogue.
- ***DSL.xml*** (*Description Session Locuteur*), enregistre l'information du locuteur.
- ***TIER.xml*** (*Traces des Informations Emises et Reçues*), enregistre (ou trace) la session de dialogue, ce qui en fournit un historique.

II.1.1.2.5 Corpus de brevets CLEF-IP 2011

Le corpus CLEF-IP contient des brevets, physiquement stockés comme une collection de fichiers XML encodant des documents de brevets. Voir le détail au VII.2.1.

II.1.1.3 Structure interne des corpus est également hétérogène

Tableau 5 : Comparaison de l'organisations logiques, physiques, et interne de quelque corpus

Corpus	Vue logique	Organisation physique	Organisation interne dans un corpus
TRANSAT	ensemble de multisegments	Fichiers de texte par langue.	L'entrée contient des segments avec des séparateurs de segments et des scores.
B@bel	Ensemble de multisegments	Fichiers de texte par langue	Listes de paires d'énoncés <Seg_Source_L1, Seg_Cible_L2>.
EOLSS	Ensemble de documents	Chaque document correspond à un dossier, contenant 2 fichiers, <i>.html</i> et <i>.unl</i> .	Le fichier <i>.html</i> est le « fichier principal », le fichier <i>.unl</i> est le « fichier compagnon », et est utilisé pour guider la segmentation du fichier <i>.html</i> .
ERIM	Ensemble de dialogues contenant des tours de parole avec leur descripteurs.	Chaque dialogue correspond à un répertoire, contenant des fichiers son (<i>.wav</i>), des fichiers transcrits (<i>.xml</i>), et des fichiers texte (<i>.txt</i>)	Un fichier son (<i>.wav</i>) attaché à un fichier de l'annotation selon une convention de nommage.
CLEF-IP 2011	Ensemble de documents (<i>.xml</i>)	Chaque document contient les segments multilingues.	Un fichier XML

II.1.2 Esquisse d'une méthode de modélisation des corpus de traductions

Pour modéliser la structure des données d'un corpus de traductions, nous nous sommes inspirés des bases lexicales et du modèle de données LEMON (*LEXicon Model for ONtologies*). Nous souhaitons trouver une structure générique pouvant décrire tous les corpus de traductions, et une façon de décrire les relations entre de tels corpus. Nous distinguons les différents niveaux de structuration de ces corpus, et comparons les corpus existants pour identifier et lister les informations propres à chaque corpus de traductions (métadonnées).

II.1.2.1 Notions essentielles pour décrire un corpus

Pour décrire un corpus, nous utiliserons les notions de *segment*, *document*, *corpus*, etc. définies dans la thèse de C. P. HUYNH et rappelées plus haut (pp. 21-24).

Les *métadonnées* sont par essence des informations sur les données. Elles servent le plus souvent à caractériser un objet par des informations homogènes, relativement à une collection d'objets. Pour formuler et définir correctement les métadonnées, il faut donc bien connaître à la fois l'objet à décrire et les caractéristiques de la collection dans laquelle il doit être déposé.

Nous avons choisi de décrire les corpus de traductions avec les informations du DUBLIN CORE METADATA INITIATIVE³⁹. Le DUBLIN CORE (*DC*) est un schéma de métadonnées générique, et il permet de décrire des ressources numériques ou physiques et d'établir des relations avec d'autres ressources. Il comprend officiellement 15 éléments relatifs à la forme de description (titre, créateur, éditeur), à la thématique (sujet, description, langue, etc.) et relatifs à la propriété intellectuelle.

II.1.2.2 Niveau des métadonnées

Pour chaque entité (*corpus*, *document*, et *segment*), il contient les métadonnées adaptées. Chaque entité peut contenir des métadonnées génériques comme le norme de métadonnées dans le DC.

1. *Title* : nom donné à la ressource (celui par lequel elle est connue officiellement).
2. *Subject* : sujet du contenu de la ressource, décrit par un ensemble de mots-clés, par des phrases, ou par un code de classification.
3. *Description* : une description du contenu de la ressource. Peut contenir un résumé, une table des matières, une référence à une représentation graphique du contenu, ou un texte libre présentant le contenu.
4. *Publisher* : une entité responsable de la diffusion de la ressource, dans sa forme actuelle.
5. *Contributor* : la ou les entités qui ont contribué à la création du contenu de la ressource.
6. *Date* : une date associée à un événement dans le cycle de vie de la ressource.
7. *Type* : la nature ou le genre du contenu de la ressource.
8. *Format* : la matérialisation physique ou numérique de la ressource.
9. *Identifier* : référence non ambiguë à la ressource dans un contexte donné.
10. *Source* : référence à une ou plusieurs ressources à partir de laquelle la ressource actuelle a été dérivée.
11. *Language* : la langue du contenu intellectuel de la ressource.
12. *Relation* : référence à une autre ressource qui a un rapport avec cette ressource.
13. *Coverage* : portée ou couverture spatio-temporelle de la ressource.
14. *Rights* : information sur les droits associés à la ressource.

La norme de métadonnées du DC peut être utilisée dans différentes entités. Par exemple, quand on fait la description d'un corpus dans *corpus*, il faut préciser des métadonnées comme le nom du corpus, le sujet du contenu de la ressource, une description du contenu de la ressource, etc. Pour *un document*, il y a aussi les métadonnées de présentation. En même temps, il existe les autres métadonnées dans différentes entités. Par exemple, pour un segment post-édité dans une MT, il a un attribut « *temps de post-édition* ».

II.1.2.3 Niveau structurel

Pour la proposition d'une structure générique, nous étudions différents corpus, et nous proposons 3 niveaux de structure abstraite d'un corpus : macro-, micro-, mésostructure.

La *macrostructure* décrit l'organisation du corpus au niveau plus haut. Par exemple, un corpus contient les documents monolingues, ou un corpus contient deux documents parallèles. Il y a plusieurs types de *macrostructure*.

- *Macrostructure-monolingue*. le corpus contient un ou plusieurs documents monolingues.

³⁹ <http://dublincore.org/>

- *Macrostructure-multilingue*. Le corpus contient un ou plusieurs documents multilingues.
 - *Macrostructure-parallèle*. Le corpus contient les documents parallèles.
 - *Macrostructure-comparable*. Le corpus contient les documents comparables.
 - *Macrostructure-mixte*. Le corpus contient les documents contenant des segments multilingues.

La *microstructure* décrit l'organisation des éléments dans un corpus. Par exemple, un document contient 100 segments français, et chaque segment est sauvegardé par ligne, ou un document du CLEF-IP contient les 3 champs <title>.

La *mésosstructure* décrit les relations entre les éléments. Par exemple, "le segment B est la traduction du segment A", ou "le segment A a une description dans le document B".

II.1.3 Validation au niveau des métadonnées

II.1.3.1 Métadonnées pour le corpus d'évaluation à la TRANSAT

Les métadonnées du corpus d'évaluation à la TRANSAT contiennent 2 parties. La première concerne la description des segments monolingues extraits à partir du corpus BTEC, et l'autre concerne la description des résultats de TA avec les données d'évaluation.

Tableau 6 : Métadonnées du corpus BTEC (les segments extraits)

Métadonnées	Valeur	
<i>Nom de corpus</i>	BTEC	
<i>Projet</i>	TRANSAT	
<i>Langue source</i>	anglais	
<i>Nombre de tours de parole</i>	Partie I : 2223	4224 au total
	Partie II : 1998	
<i>Nombre de phrases</i>	Partie I : 2373	4486 au total
	Partie II : 2113	
<i>Nombre de mots</i>	Partie I : 14K (14001)	22419 au total
	Partie II : 10K (10418)	
<i>Nombre de mots par tour de parole</i>	moyenne = 6,2; min =1; max =30	
	moyenne = 5.2; min = 1; max = 22	
<i>Nombre de mots par phrase</i>	moyenne = 5,9; médiane = 5; min =1; max = 28	
	moyenne = 4,93; médiane = 5; min = 1; max = 19	
<i>Domaine</i>	Santé ; Accidents, perte, vol ; Itinéraire ; Transports, et Location de voiture ; Non spécifique	

Tableau 7 : Métadonnées des données d'évaluation à la TRANSAT

Métadonnées		Valeur
<i>Langue cible</i>		français
<i>Nom du système de TA</i>		Systran / Reverso
<i>Version du système de TA</i>		Version 4 (Systran) Version 10 Pack Monde (Reverso)
<i>Évaluation subjective à la NIST</i>	<i>Fluidité</i>	(F1) Formulation parfaitement compréhensible sans effort, que le style soit écrit ou oral. (F2) Formulation acceptable à l'oral, éventuellement compréhensible en faisant un effort. (F3) Formulation non acceptable.
	<i>Adéquation</i>	(A1) Toute l'information est transportée. (A2) Presque toute l'information est transportée. (A3) La moitié de l'information est transportée. (A4) Peu d'information est transportée. (A5) Aucune information n'est transportée, ou il y a un contresens.
<i>Évaluation objective fondée sur la distance d'édition</i>	<i>Distance d'édition en mots</i>	Par exemple : $D_c=20$, $D_w=7$, $D=9.6$ Dc : Character distance Dw : word distance D : sentence distance
	<i>Distance d'édition en caractères</i>	
	<i>Distance d'édition pondérée</i>	
<i>BLEU</i>		Par exemple : 34.36%
<i>NIST</i>		Par exemple : 4,08
<i>Évaluateur</i>		Par exemple : Hervé Blanchon, Jean-Philippe Guillaud

II.1.3.2 Métadonnées pour le corpus Unesco/B@bel

Dans le cadre du projet UNESCO/B@BEL, on a traduit le texte du site Web B@BEL (équivalant 173 pages standard, 43200 mots) de l'anglais vers le chinois, le français, le russe, et l'espagnol. Ce corpus parallèle a été produit par la post-édition, donc il contient les métadonnées relatives à la description de la post-édition.

Tableau 8 : Métadonnées du corpus UNESCO-B@bel

Métadonnées	Valeur	
<i>Nom de corpus</i>	B@bel	
<i>Nom du projet</i>	B@bel UNESCO	
<i>Type de corpus</i>	Corpus parallèle multilingue	
<i>Méthode de création</i>	Post-édition	
<i>Langue source</i>	Anglais	
<i>Langue(s) cible(s)</i>	Chinois, français, russe, espagnol	
<i>Nombre de mots</i>	43200	
<i>Système de TA</i>	Systran v 5.0 Premium (en→fr/es/ch) ETAP3 (Labo IPPI) (en→ru)	
<i>Nombre de l'UNL graphe</i>	906 graphes UNL	
<i>UNL graphe</i>	Par exemple : voir l'Annexe 3	
<i>Échantillon de graphes UNL</i>	50 graphes UNL (environ 5 pages standard)	
<i>Temps de post-édition sur les segments</i>	<i>Temps de post-édition (par segment)</i>	Par exemple : voir l'Annexe 3
	<i>Temps de post-édition (au total)</i>	25 minutes / page standard (Systran) 20 minutes / page standard (ETAP3)

II.1.3.3 Métadonnées pour le corpus EOLSS

Le corpus EOLSS contient 25 documents. Chaque document est constitué de 2 fichiers *.aspx* (*.html*), et d'un fichier compagnons *.unl*, et de fichiers satellites (images, icônes, et autres hors-texte). Le fichier *.unl* contient les graphes UNL représentant des segments découpés à partir du fichier *.html* correspondant. Pour la description du corpus EOLSS, nous essayons de présenter les métadonnées en 3 niveaux de structure abstraite au niveau plus haut (macro-, micro-, mésostructure).

Les métadonnées de *macrostructure* ont pour but de décrire l'information globale du corpus (Voir Tableau 9).

Tableau 9 : Métadonnées du corpus EOLSS au niveau de la macrostructure

Métadonnées	Valeur
<i>Nom du corpus</i>	EOLSS (Encyclopedia of Life Support Systems)
<i>Langue source</i>	Anglais
<i>Nombre des documents</i>	25
<i>Système de TA</i>	Systran V6
<i>Nombre des mots source</i>	220K
<i>Méthode de création</i>	Post-édition
<i>Lien original</i>	http://www.eolss.net
<i>Lien source</i>	http://www.undl.org/unl-eolss/unldoc.html
<i>Nom du document</i>	Par exemple : D2_E2_03_05_TXT
<i>Nom du fichier html</i>	Par exemple : D2_E2_03_05_TXT_English.html D2_E2_03_05_TXT_French.html
<i>Nom du fichier UNL</i>	Par exemple : D2_E2_03_05_TXT.unl
<i>Nom du fichier satellite</i>	Par exemple : D2_E2_03_05_TXT_skeleton.txt

Les métadonnées de *microstructure* décrivent la structure des données dans un document. Nous présentons un exemple du document 2 (*Ground and soil water characteristics*). Il y a les métadonnées pour les fichiers HTML, un fichier source avec son fichier cible a les mêmes métadonnées. Dans le Tableau 10, on prend le fichier source comme un exemple.

Tableau 10 : Métadonnées d'un fichier HTML au niveau de la microstructure

Métadonnées	Valeur
<i>Nom du fichier source</i>	D2_E2_03_05_TXT_English.html
<i>Type du fichier</i>	.html
<i>Nombre de segments</i>	502
<i>Nombre de mots</i>	6911
<i>Titre</i>	"Ground and soil water characteristics"
<i>Langue source</i>	anglais
<i>Auteur</i>	A.G. Kocharyan
<i>Adresse</i>	Laboratory of Water Quality, Water Problems Institute, Russian Academy of Sciences, Moscow, Russia
<i>Mots-clés</i>	"physically bound water, free water, water structure, soil and ground water"
<i>Résumé</i>	"Water bedded under the earth's surface in the crust is called ground water....."
<i>Chapitre relatifs</i>	"Chapter 2"
<i>Glossaire</i>	"Artesian water : all ground water, excluding subsoil water, bedded between"
<i>Bibliographie</i>	"Bloch A.M. (1965). Water structure and geological processes Leningrad, "Nedra(This book reveals the links between the physical structure and chemical peculiarities of water)....."
<i>Esquisse biographique</i>	"Andrey G. Kocharyan graduated from the Geological Faculty, Institute of Chemistry and Oil, Baku, in 1960. He received his Ph. D. in geochemical techniques of deposits' search in 1968 at Baku State University, Geological faculty."
<i>Date de livraison</i>	2007-12-03
<i>Date de la dernière mise à jour</i>	2008-06-26

Dans le Tableau 11, on donne les métadonnées du fichier .unl.

Tableau 11 : Métadonnées d'un fichier UNL au niveau de la microstructure

Métadonnées	Valeur
<i>Nom du fichier</i>	D2_E2_03_05_TXT_en.unl
<i>Type du fichier</i>	.unl
<i>ID</i>	{S:1}
<i>Phrase source</i>	{org} Agriculture water reuse and health {/org}
<i>UNL</i>	{unl} and(health(icl>state):0U.@entry, reuse(icl>act):0K) obj(reuse(icl>act):0K, water(icl>liquid):0E) mod(water(icl>liquid):0E, agriculture(icl>activity):02) {/unl}

Les métadonnées de la *mésosstructure* (Tableau 12) visent à montrer la relation entre les éléments du corpus, Par exemple, on montre la relation de traduction existant entre une phrase source et une phrase cible, avec éventuellement l'information complémentaire.

Tableau 12 : Métadonnées d'un corpus EOLSS au niveau de la mésostructure

Métadonnées	Valeur
<i>ID</i>	\$\$_sent_1
<i>Phrase source</i>	AGRICULTURE WATER REUSE AND HEALTH
<i>Langue source</i>	anglais
<i>Phrase post-éditée</i>	RÉUTILISATION DE L'EAU D'AGRICULTURE ET SANTÉ
<i>Langue cible</i>	français
<i>Système de TA</i>	Systran
<i>Sortie de TA</i>	RÉUTILISATION ET SANTÉ DE L'EAU D'AGRICULTURE
<i>Post-éditeur</i>	GETALP/XWikiGuest
<i>Niveau de postéditeur</i>	3
<i>Date</i>	2008-09-26 11:14:40.0
<i>Temps de post-édition</i>	17s
<i>Score de qualité</i>	17/20
<i>Version</i>	1

II.1.3.4 Métadonnées pour les corpus ERIM

Le Tableau 13 montre des métadonnées du corpus ERIM au niveau de *la macrostructure*.

Tableau 13 : Métadonnées d'un corpus ERIM au niveau de la macrostructure

Métadonnées	Valeur
<i>Nom du corpus</i>	ERIM français-vietnamien
<i>Type du corpus</i>	Parole + fichier accompagnons
<i>Support</i>	.wav + .xml
<i>Langues de parole</i>	français ↔ vietnamien
<i>Nombre de tours de parole du corpus</i>	Par exemple : 4 séances contenant 102 tours de parole dans "Réservation d'hôtel"
<i>Date de création</i>	"2006-09-08"
<i>Date de la dernière modification</i>	"2006-09-10"
<i>Nombre de tour de parole par séance</i>	Par exemple : 22 tours de parole dans la séance "Réservation d'hôtel 060908-1344"
<i>Locuteur</i>	Par exemple, <i>Tung / Eric Faffiote</i>
<i>Langue du locuteur</i>	français/vietnamien
<i>Interpréteur</i>	...

Le Tableau 14 montre des métadonnées du corpus ERIM au niveau de *la microstructure*.

Tableau 14 : Métadonnées de la séance dans le corpus ERIM a au niveau de la microstructure

Métadonnées	Valeur
Nom	“Réservation d’hôtel 060908-1344”
Sujet	“Réservation d’hôtel”
Langues de parole	A et B (Par exemple, français ↔ vietnamien)
Nombre de tour de parole	41
Date de création	“2006-09-08”
Date de la dernière modification	“2006-09-10”
Nombre de tours de parole en langue A	11
Nombre de tours de parole en langue B	11
Locuteur	Par exemple, <i>Tung et Cantona</i>
Langue du locuteur	Par exemple, français / vietnamien
Interprète	...
Taille de parole	06 : 0908 mn

Le Tableau 15 montre les métadonnées du corpus ERIM au niveau de la mésostructure. À ce niveau, les métadonnées ont pour but présenter l’information sur les fichiers .wav (parole).

Tableau 15 : Métadonnées d’un corpus ERIM au niveau de la mésostructure

Métadonnées	Valeur
Nom (de l’unité de parole)	“TP1” (Tour de parole 1)
Lien entre fichiers	TP1_fr.wav + TP1_fr.xml
Langues de parole	français
Rôle	Par exemple, <i>descripteur TP1_fr.xml</i>
Locuteur	Par exemple, <i>Cantona</i>
Sexe du locuteur	Par exemple, <i>male</i>
Nationalité	français
Langue du locuteur	français
Attributs	<i>Taux d’échantillonnage : 22000</i> <i>Bits par échantillonnage : 16</i> <i>Canaux audio : 1</i>
Heure de début	Par exemple, <i>10 :20 :00</i>
Heure de fin	Par exemple, <i>10 :20 :33</i>
Durée	13
Transcription associée	<i>TP1_fr.txt</i>
Date de création	“2006-09-08”
Date de la dernière modification	“2006-09-10”

II.1.3.5 Métadonnées pour les corpus de brevets CLEF-IP 2011

Du point de vue des métadonnées, le corpus CLEF-IP 2011 est un exemple spécifique. Il contient non seulement les métadonnées génériques, mais encore les métadonnées du document de brevet. On sépare des métadonnées en 2 parties. La première partie est constituée par les métadonnées du corpus CLEF-IP, et la seconde par celles concernant sur un document de brevet.

Tableau 16 : Métadonnées du corpus CLEF-IP 2011

<i>Métadonnées</i>	<i>Valeur</i>
<i>Nom du corpus</i>	<i>CLEF-IP 2011</i>
<i>Auteur</i>	<i>EPO</i>
<i>Type de fichier</i>	<i>Xml</i>
<i>Langues</i>	<i>anglais, français, allemand</i>
<i>Langue source</i>	<i>multilingues</i>
<i>Langue cible</i>	<i>multilingues</i>
<i>Nombre de documents</i>	<i>> 2,5M</i>
<i>Type de document de brevets</i>	<i>xml</i>
<i>Nombre de brevets</i>	<i>> 1,5M</i>
<i>Nombre de fichiers</i>	<i>3,5M fichiers xml</i>
<i>Version</i>	
<i>Date de création</i>	
<i>License</i>	<i>Creative Commons Attribution-NonCommercial-ShareAlike 3.0 Unported License.</i>

Pour présenter les métadonnées d'un document de brevet, nous prenons un document (*EP-0000007-B2.xml*) comme exemple. Voici les métadonnées principales ; une version complète est donnée dans l'Annexe 6.

Tableau 17 : Métadonnées d'un document de brevet

Métadonnées	Valeur										
Nom du fichier	EP-0000007-B2.xml										
ucid (ID unique)	EP-0000007-B2										
Pays	EP (EPO)										
Numéro du fichier	0000007										
Genre	B2										
Langue	fr										
Date de publication	19841121										
Date de production	20100220										
ID de famille	Identificateur de la famille DOCDB (utilisation pour le regroupement des familles simples)										
ID de référence du fichier	interne (utilisé par le fournisseur de données)										
Nom du document de brevet	EP-78200026-A										
Retrait	Retrait d'un document (oui est la seule valeur utilisée)										
Page du document	Utilisé dans les données EP/WO pour la section de recherche-rapport. Cet élément concerne les références à la description et les revendications (amendées) du document.										
Relation	"addition, division, continuation, continuation-in-part, continuing-reissue, reissue, us-divisional-reissue, reexamination, us-reexamination-reissue-merger, correction, utility-model-basis, us-provisional-application, related-publication"										
document-id	pays, numéro de document, genre, date, texte, langue										
abstract	langue et source										
description	<table border="0"> <tr> <td>invention-title</td> <td>best-mode</td> </tr> <tr> <td>technical-field</td> <td>mode-for-invention</td> </tr> <tr> <td>background-art</td> <td>industrial-applicability</td> </tr> <tr> <td>disclosure</td> <td>sequence-list-text</td> </tr> <tr> <td>description-of-drawings</td> <td>heading img imgref p</td> </tr> </table>	invention-title	best-mode	technical-field	mode-for-invention	background-art	industrial-applicability	disclosure	sequence-list-text	description-of-drawings	heading img imgref p
invention-title	best-mode										
technical-field	mode-for-invention										
background-art	industrial-applicability										
disclosure	sequence-list-text										
description-of-drawings	heading img imgref p										
claims	langue, type déclaration, revendication texte de revendication référence de revendication										
claim	id, numéro type de revendication										

II.2 Conception de nouvelles fonctionnalités et développements en cours

II.2.1 Nouvelles fonctionnalités

II.2.1.1 Intégration de « minidictionnaires » associés aux segments

Pour faciliter l'utilisation des ressources lexicales, nous avons intégré des « minidictionnaires » à l'interface d'IMAG et de SECTRA_W. Pour les ressources « finalisées », à savoir les traductions des mots et des termes qui sont présents dans des traductions (de segments) de qualité suffisante (pas les sorties brutes de TA, mais les post-éditions).

Parmi ces traductions, on peut encore distinguer celles qui sont « recommandées » pour le projet en cours. Pour chaque segment, on dispose aussi des informations lexicales potentiellement utiles, dans la liste des ressources indiquées, et on peut sauvegarder le résultat dans une structure de données (minidictionnaires) associée à ce segment. La structure de minidictionnaires est donnée dans l'Annexe 7.

II.2.1.2 Création d'une API pour appeler le système CREATDICO

CREATDICO est un intergiciel de consultation de dictionnaires. Avec cet intergiciel, on peut envoyer un texte ou un segment pour demander une consultation pour chaque lemme du texte ou du segment. Pour faire les lemmatisations, CREATDICO appelle LEXTOH (intergiciel de lemmatisation). L'API de CREATDICO a 10 paramètres :

Tableau 18 : 10 paramètres de l'API de CREATDICO

<code>typeEntry</code> (type d'entrée),	<code>input</code> (type d'entrée)
<code>output</code> (type de sortie)	<code>ls</code> (langue source en ISO 639-2)
<code>lc</code> (langue cible en ISO 639-2)	<code>serv</code> (serveur de service)
<code>dico</code> (nom de dictionnaire)	<code>lemmat</code> (outils d'analyse morphologique)
<code>trace</code> (trace pour débogage)	<code>formule</code> (affichage de l'interface de formulaire)

La sortie de CREATDICO est une sortie dédiée à la structure "MS (MINIDICO_SECTRA)" (Voir l'Annexe 7). Pour cette sortie, il y a un fichier de configuration qui permet de définir des valeurs par défaut des paramètres, par exemple :

Tableau 19 : Exemple de fichier de configuration de CREATDICO

```
typeEntry = txt
lemma = xip (condition : ls = {fra, eng})
lemma = jieba (condition : ls = {zho})
etc.
serv = pivax
dico = * // commentaire : tous les dictionnaires de pivax
```

Dans notre cas, on n'utilise jamais les paramètres : `trace` et `formule`, et ils ne sont pas obligatoires. Voici un exemple de l'utilisation de l'API de CREATDICO :

Tableau 20 : Exemple d'un lien pour l'utilisation de l'API de CREATDICO

```
http://atoum.imag.fr/getalp/Services/Web/CREATDICO/via-Ci-Hai/CREATDICO/index.php?output=MS&input>HelloWorld!&ls=eng&lc=zho
```

Pour faciliter l'utilisation des dictionnaires, nous avons ajouté un panneau de dictionnaires à l'interface de SECTRA_W. Quand l'utilisateur fait la post-édition, il peut double-cliquer sur un mot, puis choisir les dictionnaires existants dans ce panneau.

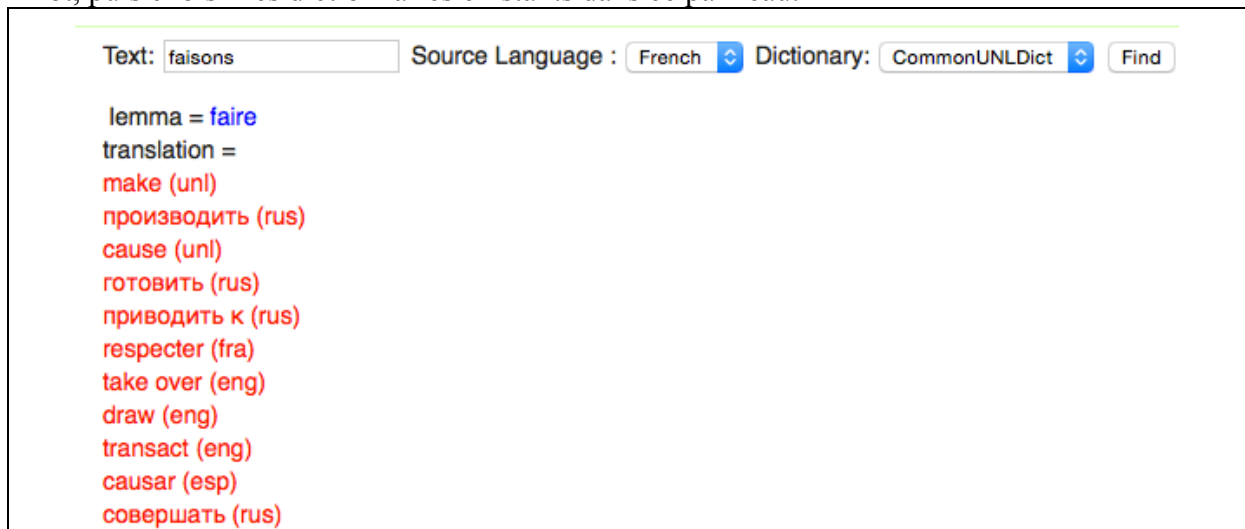


Figure 13 : Capture d'écran de panneau de dictionnaires ajouté à SECTra_w

Le résultat est affiché dans le champ en bas. Par exemple, si on double-clique sur le mot « *faisons* », il est ajouté automatiquement dans le champ de saisie. On choisit la langue source et le dictionnaire, et on clique sur le bouton « *Find* ». Le résultat est affiché en couleur rouge (Figure 13).

II.2.1.3 Manipulation de la MT

Sur un segment. Le résultat des systèmes de TA a évolué lors des mises à jour. Sur SECTRA_W, il faut permettre à l'utilisateur de supprimer les traductions de TA, et à de rappeler les systèmes de TA de la dernière version. Si la post-édition n'est pas bonne, s'il existe les erreurs sur le segment post-édité, il faut aussi permettre à l'utilisateur de le supprimer. Pour réaliser cette fonction, nous avons ajouté 2 boutons avant la traduction de TA. Le bouton « *Clean* » permet de supprimer la traduction de TA sur l'interface de SECTRA_W, et aussi sur la MT. Le deuxième bouton « *Get* » permet de recalculer (retraduire) le segment source par TA. Si la post-édition n'est pas bonne, on peut le supprimer par bouton « *Delete* ». (Voir Figure 14).

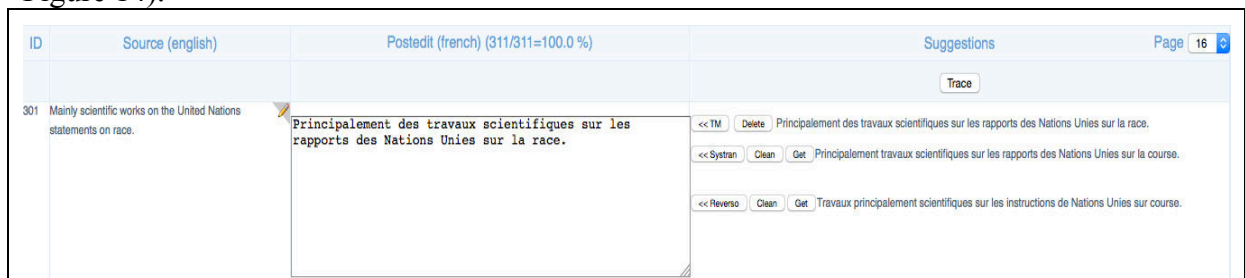


Figure 14 : Interface de SECTra_w intégrant les boutons « *Delete* », « *Clean* » et « *Get* »

Sur un corpus. On a besoin de (re)traduire une sélection de segments, éventuellement tous, depuis l'interface de SECTra_w. Pour réaliser cette fonction, nous avons créé une interface permettant de choisir les segments source et un ou plusieurs systèmes de TA. L'utilisateur peut cocher les segments sources, et les traduire par le système de TA choisi. Sur cette interface, il peut aussi importer des résultats de TA dans la MT associée. Il demande de préparer les segments source et les traductions, et les phrases sont alignées par ligne.

Voici une image d'écran (Figure 15).

The screenshot shows the SECTra_w web interface. At the top, there is a header with 'Corpus name: mt_incmoses'. Below this is a navigation bar with tabs for 'All', 'MTs', 'Pivot language', 'Target language', and 'Translation selection'. The 'All' tab is active, showing a table of 13 documents. Each document row includes a checkbox, a number, a document name, source language, pivot language, target language, and an 'Operate' button. To the right of the table is a 'Translation selection' panel with various input fields and buttons.

No	Document name	Source language	Pivot language	Target language	Operate
<input type="checkbox"/>	1 DOC1	english	No pivot language	--select--	Start
<input type="checkbox"/>	2 DOC2	french	No pivot language	--select--	Start
<input type="checkbox"/>	3 DOC3	french	No pivot language	--select--	Start
<input type="checkbox"/>	4 DOC4	english	No pivot language	--select--	Start
<input type="checkbox"/>	5 DOC5	chinese	No pivot language	--select--	Start
<input type="checkbox"/>	6 DOC6	chinese	No pivot language	--select--	Start
<input type="checkbox"/>	7 DOC7	chinese	No pivot language	--select--	Start
<input type="checkbox"/>	8 DOC8	chinese	No pivot language	--select--	Start
<input type="checkbox"/>	9 DOC9	french	No pivot language	--select--	Start
<input type="checkbox"/>	10 DOC10	french	No pivot language	--select--	Start
<input type="checkbox"/>	11 DOC11	french	No pivot language	--select--	Start
<input type="checkbox"/>	12 DOC12	french	No pivot language	--select--	Start
<input type="checkbox"/>	13 DOC13	french	No pivot language	--select--	Start

Figure 15 : Traduction des segments sélectionnés et ajout à la MT

II.2.2 Programmabilité

II.2.2.1 API avec CREATDICO, TRADOH, SEGDOC

SECTra_w a montré une possibilité de l'utilisation de l'API, et il permet l'accès à différents types de services. Pour appeler CREATDICO, Tradoh et SEGDOC, on a intégré l'API sur SECTRA_W, et les fonctions sont réalisées en Groovy. On peut suivre une instruction pour remplir les paramètres dans un lien. Quand on accède ce lien avec nos valeurs pour les paramètres, on peut recevoir le résultat immédiatement sur une nouvelle page Web. Par exemple, on peut appeler les systèmes de TA via TRADOH sur l'API de SECTRA_W (Figure 16). On remplit les paramètres (nom de TA, langue source/cible, et texte source), et le système génère un lien pour obtenir la traduction.

The screenshot shows the SECTra_w web interface. At the top, there is a logo for 'seCTra_w' and a navigation bar with links for 'Home', 'Import', 'Evaluation', 'Post-edition', 'TM', 'Admin', 'Translation', and 'Contact us'. Below the navigation bar is a section titled 'Call Tradoh'. A warning message is displayed in a yellow box, indicating that parameters are required for the API call.

Call Tradoh

⚠ Please give parameters:
 MT (this is MT name),
 source language(source),
 target language(tran),
 source text(sourcetext),
 context (context- this parameter is option).
 Please contact with [technical support](#) for more details

Figure 16 : Exemple de l'API « Call Tradoh »

II.2.2.2 Vers un langage de commandes pour SECTra_w

La définition d'un tel langage est particulièrement intéressante parce qu'il permettra de réaliser toutes les opérations sur les corpus, à la fois en ligne de commande et en une formulaire, et aussi de générer des codes pour les interfaces, ce qui accélérera leur développement tout en assurant leur homogénéité.

L'exemple le plus utilisé de langage déclaratif dédié à la sélection d'éléments et à la commande d'actions à exécuter sur toutes les données sélectionnées est SQL. Le langage à définir est bien de cette nature.

Par exemple, en MySQL, la commande qui sélectionne tous les enregistrements de la relation mémorisée dans la table `demo2`.`segments` et dont le champ `creation_date` est compris entre '2013-01-16' et '2013-01-18' est :

```
SELECT * FROM `demo2`.`segments` WHERE `creation_date` > '2013-01-16' AND `creation_date` < '2013-01-18'.
```

Le langage d'interrogation des données permettant les sélections, la définition de sélections dans les corpus doit s'en inspirer.

Par exemple :

```
Selection 3 = $seg \in demo2 \and $creation_data \in 2013-01-06 .. 2013-01-18 \and {{ $level = 3 \and $score > 13}}
```

L'analyseur du langage en question sera faite en ANTLR (Parr and Quong, 1995), générateur de compilateurs dont le GETALP a une grande expérience, à cause du projet Ariane-Y. Les règles de grammaire décrivant le langage et transmises à ANTLR peuvent définir des actions.

II.2.3 Modularité

II.2.3.1 Extraction de codes depuis XWiki vers Java dans des fichiers externes

Dans I.2.2.3, nous avons un peu parlé du problème de la programmation en XWiki, et nous avons listé et classifié les fonctions de SECTRA_W/IMAG dans le cadre du projet Traouiero. Pour réaliser la modularité de SECTRA_W/IMAG, nous avons aussi extrait les codes de XWiki (en Groovy). Comme Groovy est proche de Java, nous les avons reprogrammés en Java, et les avons publiés sous forme de packages de Java (*.jar*).

Par exemple, l'extraction de la fonction « Segmentation de la page Web » est détaillée à http://tools.aximag.fr/prepaTradSeg/_README.htm.

La version initiale de ce service avait seulement pour objectif de fournir aux IMAG le prétraitement d'une page Web. D'abord, on a l'URL de la page ou un contenu HTML (sérialisation du DOM fournie par un navigateur). Ensuite on appelle un outil de segmentation d'un contenu HTML en segments, avec le service <http://tools.aximag.fr/enrichir/seghtml> et un paramétrage de l'appel à ce service, tel que le résultat soit fourni dans un format attendu par l'uniformisateur d'appels à des systèmes de TA. Cela produit la même forme de sortie que celle qu'aurait fournie le service GT pour le contenu source, avec l'habillage de chaque segment source identifié par le traitement du segmenteur par l'intercalage de 2 nœuds imbriqués dans le DOM.

II.2.3.2 Vers un "serveur corporal pour la TA" avec lien en cours SECTra-Ariane

ARIANE-G5 est un générateur de systèmes de TA reposant sur 5 langages spécialisés pour la programmation linguistique (LSPL)⁴⁰. Chacun de ces langages est compilé, et les tables internes produites sont données en paramètres au "moteur" du langage.

⁴⁰ Les LSPL sont conçus pour implémenter les outils linguistiques. Au CETA, l'idée de construire un système de TALN basé sur plusieurs LSPL est venue des travaux de B. Vauquois, G. Veillon, J. Veyrunes et leurs collègues des années 1962-1967. Cette idée s'est également imposée dans (Hutchins, 1986) Hutchins, W. J. (1986). Machine translation: past, present, future, Ellis Horwood: Chichester.ou encore (Arnold, Balkan et al., 1994) Arnold, D., Balkan, L., Meijer, S., Humphreys, R. and Sadler, L. (1994). Machine translation: An introductory guide. NCC Blackwell. 200 p..

Bien que le système ARIANE-G5 ait déjà une fonction de post-édition et de gestion de corpus, on voudrait déléguer à un « serveur corporal » le soin de gérer les corpus. La motivation de cette délégation est identique à celle des connaissances lexicales où on pourrait profiter des avantages d'un système sur le Web comme SECTRA_W, et surtout avancer vers une unification des ressources corporales.

Dans le système ARIANE-G5, à l'entrée et à la sortie du processus de traduction, l'unité de traduction est une simple chaîne de caractères. Les chaînes doivent être codées en EBCDIC. Dans l'implémentation ARIANE-H (HELOÏSE de V. Berment (BERMENT AND BOITET, 2012)), elle sont en ASCII.

En ARIANE-G5, l'organisation est la suivante :

- fichier texte source.
- fichier de description arborescente (utilisé pour calculer les unités de traduction).
- fichier de segmentation en unités de traduction, selon une fenêtre paramétrable par le linguiste (200-300 caractères).
- fichiers d'arbres de résultats intermédiaires possibles.
 - chaque fichier selon la chaîne d'exécution/production
 - fichiers selon les phrases. Chaque fichier contient la liste des arbres fournis par cette phase, un pour chaque unité de traduction.
- fichiers résultats de traduction automatique d'une ou plusieurs chaînes d'exécution. Le numéro de la chaîne d'exécution est inclus dans le nom de fichier.
- fichiers de révision.

Avec un corpus, on a 2 fichiers « spéciaux », \$DESCR contenant la liste des textes de ce corpus avec le format <nom externe, nom interne>, et un autre contenant une liste hiérarchique de séparateurs (au maximum 8).

ARIANE-G5 a 3 types de fichier de résultat.

Résultat intermédiaire (liste d'arbres décorés) est un format d'échange entre les étapes du système. Ce format représente des arbres décorés, les formes et UL utilisées.

Résultat brut de la traduction est la concaténation des traductions des unités de traduction, et on n'a plus accès aux unités individuelles.

Révision. Pour un texte donné, il y a un seul fichier de révision par langue et par chaîne de traduction.

La version actuelle de SECTRA_W a été construite pour être utilisable aussi bien par des programmes que par des utilisateurs humains. Pour l'intégrer à ARIANE-G5, il faut :

- définir et implémenter certaines fonctions sur les corpus, existant en ARIANE-G5, mais pas encore intégrées à SECTRA_W ;
- définir une syntaxe appropriée pour échanger avec ARIANE-G5 des segments, des documents, des corpus et des mémoires de traductions (formats d'import-export)
- définir une structure interne dans ARIANE-Y pour les données "corporales".

En particulier, la révision (en construction) ARIANE-Y doit pouvoir envoyer des résultats de traitement,

- sur des segments (par exemple des arbres d'analyse ou des graphes UNL, à ranger comme des « annotations » au même titre que les TA et les PE),
- sur des documents (graphe de chaînes d'arbres, par exemple).

Chapitre III Variété des iMAG et de leurs usages : de l'accès multilingue à la création de bons corpus bilingues et à la traduction littéraire contributive de qualité

Introduction

Nous faisons ici le point sur l'ensemble des passerelles iMAG créées depuis 2009, et décrivons leurs différents usages. En ce qui nous concerne, nous en avons créé une instance, essentiellement destinées à la création de bons corpus bilingues pour la TA français-chinois, à des expérimentations sur d'apprentissage incrémentale, et à un projet personnel de traduction littéraire.

III.1 Liste (avec les MT associées)

Une passerelle iMAG peut être utilisée pour supporter non seulement la post-édition des sites Web et des documents, mais aussi l'expérimentation, la traduction, l'évaluation, et le support d'accès multilingue aux sites Web commerciaux.

Depuis 2009, nous avons créé 213 iMAG sur 20 MT, et notre plate-forme supporte plus de 8 langues source et 30 langues cible. Il y a plusieurs types de sites Web "élus" (accédés en multilingue), qui couvrent de nombreux domaines. Dans cette section, on montre des exemples (par liste) d'iMAG des différents types.

Laboratoire/Université. Nous avons créé des iMAG pour les sites Web de laboratoires et d'universités, semblable à celui de notre laboratoire (LIG), ISCC, LICIA⁴¹, NICT⁴², etc. et pour les sites Web d'universités sont comme l'Université Joseph Fourier (UJF), IITB, etc. (Tableau 21).

Tableau 21 : iMAG pour les sites Web de laboratoires et d'universités

<i>ID</i>	<i>Nom de site</i>	<i>Lien</i>	<i>Langues source</i>	<i>Langues Cible</i>	<i>TA</i>	<i>MT</i>
<i>LIG-Lab</i>	<i>Laboratoire d'Informatique de Grenoble</i>	<i>http://www.liglab.fr</i>	<i>French</i>	<i>Chinese , English , French ,</i>	<i>Moses+ Google</i>	<i>liglab</i>
<i>IPPI-IITP</i>	<i>Институт проблем передачи информации (Institute for Information Transmission Problems)</i>	<i>http://iitp.ru</i>	<i>Russian</i>	<i>Chinese , English , French ,</i>	<i>Google</i>	<i>demo</i>
<i>ISCC</i>	<i>Institut des sciences de la communication CNRS / Paris-Sorbonne / UPMC</i>	<i>http://www.iscc.cnrs.fr</i>	<i>French</i>	<i>Arabic , Bulgarian , Chinese , Croatian , Czech , Danish , Dutch , English ,</i>	<i>Systran</i>	<i>demo</i>
<i>UJF</i>	<i>Université Joseph Fourier</i>	<i>http://www.ujf-grenoble.fr</i>	<i>French</i>	<i>Arabic , Belarusian , Bulgarian , Catalan , Chinese , Czech , ...</i>	<i>Google</i>	<i>demo2</i>

⁴¹ <http://licia-lab.imag.fr/index.php/fr-FR>

⁴² <http://www.nict.go.jp/>

<i>IITB</i>	<i>Indian Institute of Technology BomBay</i>	<i>http://www.iitb.ac.in</i>	<i>English</i>	<i>..., Hungarian , Icelandic, Japanese , Korean , ...</i>	<i>Google</i>	<i>demo2</i>
-------------	--	------------------------------	----------------	---	---------------	--------------

Sociétés. Le site AXIMAG.FR⁴³ supporte des iMAG dédiées aux sites Web d'organismes et de sociétés (commerciales ou non), comme AcXys, AMIES, IDB, etc. Ce type d'iMAG dédié aux sites contient souvent une MT spécifique pour un site Web.

Tableau 22 : iMAG pour les sites Web d'organismes et de sociétés

<i>ID</i>	<i>Nom de site</i>	<i>Lien</i>	<i>Langues source</i>	<i>Langues Cible</i>	<i>TA</i>	<i>MT</i>
<i>AcXys</i>	<i>AcXys</i>	<i>http://www.acxys.com/</i>	<i>English</i>	<i>Arabic , Chinese , Dutch , French , German , Greek , Hungarian , Italian , Japanese , Korean , Malay , Polish , Portuguese , Russian , Spanish , Thai , Turkish , Vietnamese</i>	<i>Google</i>	<i>mt_acxys</i>
<i>iMAG-AMIES-privee</i>	<i>AMIES</i>	<i>http://www.agence-maths-entreprises.fr/</i>	<i>French</i>	<i>Chinese , English , German , Japanese , Portuguese , Russian , Spanish</i>	<i>Google</i>	<i>mt_amies</i>
<i>IDB</i>	<i>Islamic Development Bank</i>	<i>http://www.isdb.org</i>	<i>English</i>	<i>Arabic , Chinese , French , German , Portuguese , Russian , Spanish , Swahili</i>	<i>Google+Systran</i>	<i>demo</i>
<i>Bio-Clean</i>	<i>Bio-Clean</i>	<i>http://projet-bioclean.3beesonline.com</i>	<i>French</i>	<i>Arabic , Chinese , German , Portuguese,</i>	<i>Google</i>	<i>demo2</i>

Projets et expérimentations. Le site aximag .fr peut aussi être utilisé comme une plate-forme pour la partie de post-édition ou d'évaluation de support multilingues d'un projet, ou pour des expériences de recherche. Les iMAG de projet sont EOLSS, AKENOU, HOMERICA, etc., et les iMAG d'expérimentation sont POWERS, BEMBOOK, etc.

Tableau 23 : iMAG pour des projets et des expérimentations

<i>ID</i>	<i>Nom de site</i>	<i>Lien</i>	<i>Langue source</i>	<i>Langues cible</i>	<i>TA</i>	<i>MT</i>
<i>EOLSS</i>	<i>EOLSS</i>	<i>http://www-clips.imag.fr/geta/User/christian.boitet/iMAGs-tests/EOLSS</i>	<i>English</i>	<i>Chinese , French , Japanese , Russian , Spanish</i>	<i>Google</i>	<i>demo2+ eolss</i>
<i>Powers</i>	<i>The book of me</i>	<i>http://www-clips.imag.fr/geod/User/laurent.besacier/TRANSLATION-EXP/The_Book_of_Me.html</i>	<i>English</i>	<i>French , Romanian</i>	<i>Moses</i>	<i>mt_powers</i>
<i>BEMBO OK</i>	<i>BEMBO OK</i>	<i>http://www.kenrico.com/media/bembook/21/21.htm</i>	<i>English</i>	<i>Arabic , Bengali , Bulgarian , Chinese , French , German , Hindi , Hungarian , Italian , Japanese , Malay , Marathi , Polish , Portuguese , Russian , Spanish</i>	<i>Google+Systran</i>	<i>demo1</i>

⁴³ <http://service.aximag.fr/xwiki/bin/view/home/imag>

III.2 Commentaires sur les utilisations actuelles

III.2.1 Accès à des sites Web d'organismes ou de sociétés

Plusieurs sites Web de sociétés, de laboratoires, ou d'universités n'offrent pas d'accès multilingue, ou supportent un ou deux langues étrangères. Mais au laboratoire ou à l'université, il y a beaucoup de chercheurs ou étudiants étrangers, et le site Web devrait distribuer l'information à tous. Il y a une vraie demande pour des sites Web multilingues. Le serveur aximag.fr supporte plusieurs sites Web de sociétés, de laboratoires, ou d'universités. Par exemple, en 2011, nous avons créé une iMAG (*iMAG-LIG-LAB*) pour le site Web de notre laboratoire, en lui attachant une MT dédiée (*liglab*). Au départ de *iMAG-LIG-LAB*, nous avons environ 10K segments source en français. Nous avons fait la post-édition du français vers le chinois. Après première passe de post-édition, nous avons obtenu environ 2K segments post-édités en chinois, et les pages Web principales du laboratoire ont été bien présentées en chinois (en juin 2013)⁴⁴.

Pour aider les étudiants étrangers à accéder aux sites Web d'universités, nous avons expérimenté la technique iMAG pour traduire les sites Web de l'UJF. Les étudiants étrangers parlant (ou comprenant) le français ont aidé à post-éditer les pages Web du français vers leurs langues maternelles.

Des iMAG ont été créées non seulement pour des organismes d'enseignement (comme LIG, ISCC, ANRT, etc) et les universités, mais aussi pour les sites Web commerciaux. Pour les sites Web comme ACXYS, nous traduisons les segments source de l'anglais vers 18 langues étrangères. Le site Web ACXYS contient les segments traduits en français, allemand, russe, et polonais. Nous les ajoutons dans la MT dédiée (*mt_acxys*) pour la cohérence entre le site officiel et l'iMAG.

L'iMAG a aussi un cas particulier. Depuis mai 2010, la Métro propose l'intégralité de son site web en version chinoise et anglaise, les deux langues les plus parlées au monde. Nous avons créé une iMAG dédiée, et une mémoire dédiée (*lametro*) pour le site Web de La Métro, permettant d'obtenir des textes traduits en constante amélioration. Son atout majeur réside dans une mémoire spécifique au site : conservant les segments post-édités, celle-ci les réutilise ensuite en tenant compte du type de vocabulaire et des termes utilisés sur le site. La qualité de traduction est ainsi en constante amélioration.

III.2.2 Aide à la traduction de documents : rapports, parties de thèse, manuels...

Le but essentiel d'une iMAG est d'offrir un bon support pour l'accès multilingue à un site Web. Mais cette technique peut aussi nous aider à la traduction de documents, si on fournit le document au format HTML. Nous avons travaillé sur la post-édition des rapports, de parties de thèse, manuels, etc.

Beaucoup d'universités publient leurs matériels éducatifs gratuitement, mais ils ne sont généralement pas disponibles dans plus d'une ou deux langues, ou disponibles moyennant un supplément. L'approche iMAG offre une alternative rapide, pratique et de très faible coût pour obtenir des versions multilingues de matériel éducatif, comme les manuels scolaires qui sont convertis dans un format compatible iMAG. Un bon exemple est l'iMAG de démonstration pour le livre « Bioelectromagnetism » de Jaakko Malmivuo et Robert Plonsey, visitable à <http://service.aximag.fr/xwiki/bin/view/imag/BEMBOOK>.

⁴⁴ Tout le travail s'effectue sur une MT dédiée (*liglab*). Après juillet 2014, on a créé une autre iMAG pour le site du LIG, utilisant la MT générique (partagée) *demo2*.

Un autre exemple, est l'IMAG *xan-fr*⁴⁵, utilisée pour traduire des documents du français vers l'anglais. Le dernier travail est la post-édition du résumé de ma thèse (voir Figure 17). On a utilisé GT pour prétraduire le texte, puis fait la post-édition sur le texte anglais.

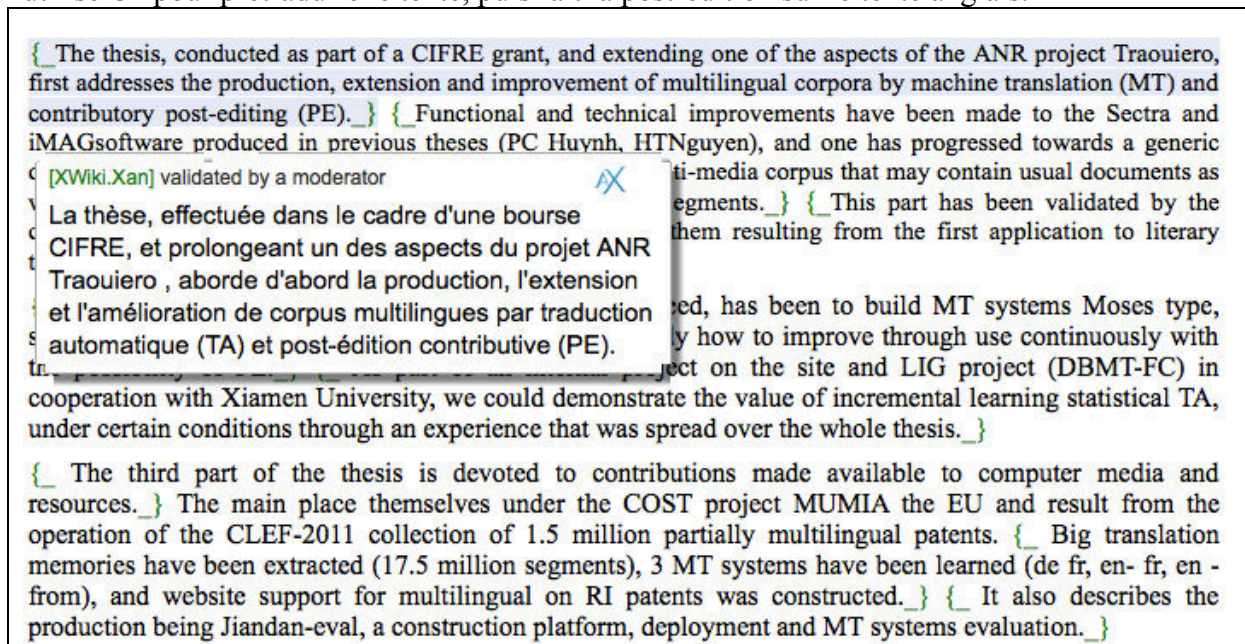


Figure 17 : Post-édition d'un document français accédé en anglais (résumé de la thèse de Lingxiao WANG)

III.2.3 Accès multilingue à des documents pédagogiques : MACAU

R. Kalitvianski et C. Boitet (Kalitvianski et al., 2012) ont lancé le projet MACAU (Multilingual Access & Contributive Appropriation for Universities) pour permettre aux étudiants étrangers de l'UJF d'accéder dans leur langue aux supports pédagogiques, par amélioration contributive des "prétraductions" par TA. Pour l'accès multilingue dans MACAU-OFI, nous avons créé une IMAG « dédiée » à ce projet.

L'idée de base est de ne pas chercher à construire différentes versions, une en chaque langue, mais de garder l'unicité de la version « originale », et de la « reconstruire » dans différentes langues en choisissant une « meilleure traduction » (ou suggestion de traduction) pour chaque segment dans une MT (mémoire de traductions) associée. Pour l'instant, nous supposons que la version originale de chaque document est dans une seule langue, mais, dans le futur, on utilisera un détecteur de langue, et il pourra y avoir des segments écrits en différentes langues dans le document original.

Le projet MACAU-OFI consiste à essayer d'approfondir cette idée en allant dans deux directions : d'abord, permettre aux étudiants de participer en contribuant par des fragments de notes de cours ou d'exercices, et ensuite en enrichissant l'outil par l'intégration de ce que nous appellerons une « ontologie du domaine », permettant en particulier, dans le cas d'OFI, de viser à l'autoformation par utilisation d'outils de simulation d'automates, grammaires, et graphes.

Au cours des deux dernières années universitaires, nous avons encouragé les étudiants de Master 1 d'informatique de l'UJF à produire des documents de cours sur la complexité calculatoire. Chacun a eu un compte sur aximag.fr. Certains ont beaucoup contribué, d'autre moins. Voici les nombres d'étudiants concernés, par langue maternelle.

⁴⁵ <http://service.aximag.fr/xwiki/bin/view/imag/xan-fr>

Tableau 24 : Nombre de langue du projet MACAU (06/2013)

<i>Langue</i>	<i>Nombre</i>
<i>Chinois</i>	<i>7</i>
<i>Arabe</i>	<i>2</i>
<i>Russe</i>	<i>1</i>
<i>Anglais</i>	<i>2</i>

Les fichiers reçus étaient dans les formats doc, odt, pdf, tex, html. Certains étaient des cours complets, d'autres ne contenaient que quelques chapitres.

Tableau 25 : Statistiques de documents dans MACAU (06/2013)

	<i>tex</i>	<i>doc, docx</i>	<i>odt</i>	<i>pdf</i>	<i>html</i>
<i>Nombre de fichiers</i>	<i>16</i>	<i>5</i>	<i>5</i>	<i>19</i>	<i>7</i>
<i>Cours complets</i>	<i>2</i>	<i>1</i>	<i>0</i>	<i>0</i>	<i>0</i>
<i>Chapitres disponibles (hors cours complets)</i>	<i>1-8</i>	<i>1-8</i>	<i>4, 7</i>	<i>1-7</i>	<i>-</i>
<i>Autres</i>	<i>-</i>	<i>1 corrigé d'un quick, 1 fichier de notes explicatives</i>	<i>-</i>	<i>-</i>	<i>Site avec notes sur Pseudo-Pascal, RAM, 2SAT et exercices</i>

III.2.4 Évaluation

Depuis 2007, les fonctionnalités de SECTRA ont évolué pour l'adapter à des usages différents.

Le premier usage de SECTRA_w a été la « campagne d'évaluation » du projet TRANSAT. Il a été utilisé avec succès, fin 2007, dans le cadre du projet TRANSAT de FT R&D. Avec SECTRA_w, après avoir importé un corpus source, et éventuellement les traductions de référence, on peut appeler plusieurs systèmes de TA, stocker leurs résultats, et demander à des juges d'effectuer l'évaluation subjective (fluidité, adéquation). SECTra_w fournit plusieurs méthodes d'évaluation objective (NIST, BLEU, etc.), et permet aussi d'effectuer l'évaluation objective liée à la tâche, en permettant à des participants de post-éditer les résultats de systèmes de TA, et en mesurant une distance d'édition (et/ou le temps de post-édition). Les résultats post-édités peuvent être ajoutés à l'ensemble des traductions de référence, ou le constituer s'il n'y a pas de références.

Le deuxième usage est l'évaluation pour améliorer des systèmes (AI). Nous avons utilisé SECTRA_w/MAG pour évaluer les résultats des systèmes de TA, puis nous faisons la post-édition pour créer une MT de bonne qualité. Nous avons utilisé cette méthode pour la construction des ressources français-chinois, comme la MT sur l'énergie et la MT du LIG-LAB. Notre système de TA a été amélioré par AI (Voir Chapitre IV et Chapitre V).

Enfin, SECTRA_w/MAG a été utilisé pour supporter l'expérimentation et la recherche d'une formule de prédiction de choix optimal par Haozhou WANG (Wang, 2015).

III.3 Utilisations plus novatrices : production de bons corpus parallèles, et post-édition de textes littéraires pour l'auto-apprentissage ou pour la traduction contributive

III.3.1 Production de « corpus parallèles » de qualité

Thanks to SECTra_w in-built system of annotation of each translation or post-edition of a segment by a reliability level (from * to *****) and a quality score (0..20), one can extract from the TM associated to a website S a subset verifying any predicate based on levels and scores.

To implement that, we have introduced and implemented into SECTra_w the notion of selection. A selection is defined intentionally (by a predicate) or extensionally (by an explicit list), and can be named, for later recall.

Take for example the TM of the website of Greater Grenoble (La Métro) that contains 2500 web pages, or about 30000 segments. More than half have been pre-translated and post-edited into Chinese for the Shanghai Expo in 2010. We may select a "quite good part" of this TM by creating the selection:

```
TM-lametro-extract-good = TM_select (lametro, [level=3 & score >=13 | level=4 & score >=12 | level=5 & score >=11]).
```

The following example shows an even simpler extraction, from the French-Chinese part of the Demo2 TM associated with iMAG-Doc_Par_jour shown on. The predicate is simply [level=3 & score >=13], and its parameters can be directly chosen through the GUI.

<input type="checkbox"/>	No	Pseudo Doc	Source	Cible	Stars	Notes
<input type="checkbox"/>	1	DOC16	la salle du haut conseil située au 9ème étage de l'Institut du monde arabe, dans le 5ème.	位于巴黎5区的阿拉伯世界博物馆10楼高级理事会议厅。	3	20
<input type="checkbox"/>	2	DOC16	le 24ème étage de la tour zamansky, l'université pierre et marie curie, sur le campus de jussieu, dans le 5ème.	位于巴黎5区的皮埃尔和玛丽·居里大学加希耶校区的扎曼斯基大楼25楼。	3	20
<input type="checkbox"/>	3	DOC16	le 18ème étage de la bibliothèque françois mitterrand, dans le 13ème.	位于巴黎13区的法国国家图书馆密特朗官19楼。	3	20
<input type="checkbox"/>	4	DOC16	le 6ème étage de l'hôtel industriel de dominique perrault, dans le 13ème.	位于巴黎13区的多米尼克·佩罗工业馆7楼。	3	20
<input type="checkbox"/>	5	DOC16	la nuit blanche 2012 permettra aux visiteurs de découvrir la ville lumière d'en haut grâce à 15 belvédères normalement fermés au public.	2012巴黎不眠夜将使参观者可以从15个平时不对外开放的平台发现欣赏巴黎这座光影之城。	3	20
<input type="checkbox"/>	6	DOC16	jk rowling	jk罗琳	3	20
<input type="checkbox"/>	7	DOC16	en effet, jk rowling, qui a créé notre sorcier à lunettes, pourrait se replonger dans l'univers d'harry potter.	事实上, "创造了"我们那位眼镜魔法师的jk罗琳, 有可能会写哈利波特的魔法世界里的续集。	3	20
<input type="checkbox"/>	8	DOC16	c'est en 2011 que la saga de nos célèbres sorciers s'est achevée, à l'issue du septième livre intitulé « harry potter et les reliques de la mort » qui a été divisé en deux parties au cinéma.	2011年, 第七本《哈利波特和死亡圣器》分为上下两部电影, 上映结束之后, 我们这著名的哈利波特系列魔法小说完结了。	3	20
<input type="checkbox"/>	9	DOC16	cinq ans après le dernier opus de cette série à succès, jk rowling revient avec une nouvelle œuvre, pour adultes cette fois.	在这成功的系列小说最后一章完结的5年之后, jk罗琳带着她的新作品回归了。这次是给大人看的小说。	3	20
<input type="checkbox"/>	10	DOC16	et la star des librairies a su entretenir le mystère.	这位图书史上的明星会将这光芒延续下去。	3	20
<input type="checkbox"/>	11	DOC14	la joie	欢乐	3	20
<input type="checkbox"/>	12	DOC16	« une bourgade apparemment idyllique mais qui va faire face "aux tourments les plus violents ».	"一个表面看起来非常美好的田园小镇, 但是它将会面临最猛烈的动荡。"	3	20

Figure 18 : Extraction of a "good" TM from a TM produced by "natural" post-edition

The selection obtained can then be exported, as 2 parallel files (source and post-edition) in a simple XML format (Figure 19). SECTra_w also provides additional information (TM, Last updated, Duration of post-editing, post-editor, etc.), and other available download formats (TMX, TXT, and CSV). These data can be used later to "feed" an empirical Moses-based MT system that will become specialized to that website⁴⁶.

⁴⁶ We are running such an experiment but cannot describe it here for lack of space.

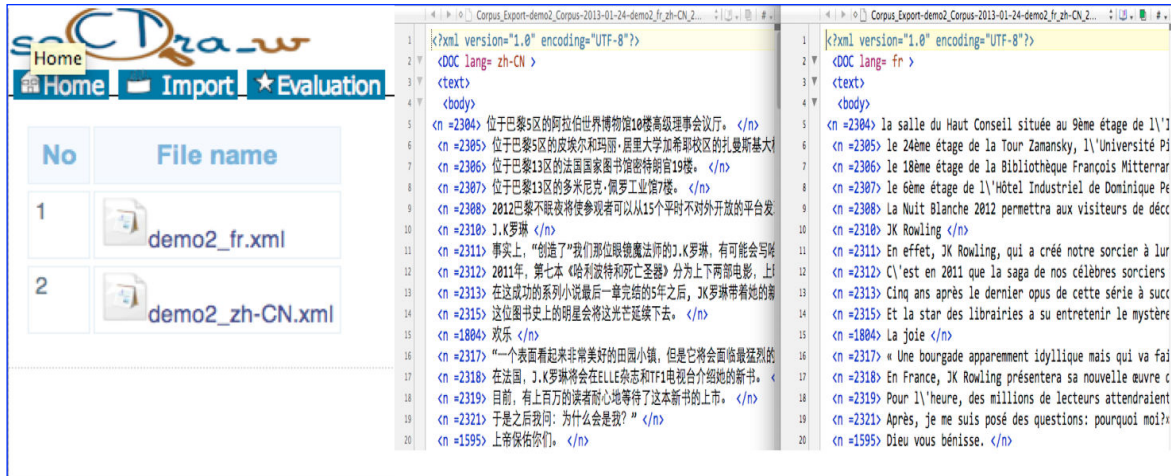


Figure 19 : Export of a « good » part of a TM

That possibility is very interesting in the current context. It has been proven that MT systems can be specialized to sublanguages and produce outputs of very high usage value (Chandioux, 1988) (Isabelle, 1987). That means that the outputs are quite readable, and very cheap to post-edit to produce professional quality output.

In recent experiments with a Paris-based multilingual content processing firm, a Moses instance built from a high proportion of a 300K bi-segment TM mixed with a standard parallel corpus extracted from EuroParl (Koehn, 2005) got a BLEU (Papineni et al., 2002) score of about 70%. At this high level, BLEU correlates with usage value: it takes typically 10-15 minutes only to post-edit the equivalent of 1 standard page (250 words, or 400 kanjis), instead of 1 hour to produce a draft translation. But that method works only if a parallel corpus specialized to the sublanguage at hand is available, and that is quite rare in practice.⁴⁷

The situation is similar if the considered MT system is built by an “expert” method (as TAUM-METEO and then METEO).

For example, there is no available parallel Chinese↔French corpus for e-mails, chats, and short technical notes. Building a parallel corpus from scratch is not an option because of the cost of the operation and the scarcity of translators knowing both languages and the technical terms.

Using an iMAG offers a graceful way to solve that difficulty. Whatever MT systems are available, one can begin without any delay to start the bilingual service needed (a web-based chat, for example), routing messages and documents through web pages, and using iMAGs to make them accessible (and improvable) in the desired languages. After a while, the TM-S dedicated to the (empirically defined) sublanguage of S will contain enough “good” bi-segments to extract them and use them to build a specialized instance of an MT system (for example, a specialized Moses-S system⁴⁸).

An important point here is that, in order to encourage end users to post-edit, post-editing should be made very simple and user-friendly. One should refrain from transforming it into a debugging environment for some MT systems. That would also go against the principle to be open to as many MT systems as possible.

⁴⁷ Remember: in 2001, Language Weaver (LW) claimed « to be able to produce an MT system overnight » from a large enough parallel corpus. While that was undoubtedly true, LW produced actually only 4 MT systems in 4 years... because parallel corpora corresponding to the translation needs of solvable clients were and are hard to find.

⁴⁸ We have built a French-Chinese Moses system for iMAG-LIG, based on 12000 already post-edited segments.

III.3.2 « Voyage au centre de la terre » de Jules Verne

Nous avons post-édité 21 chapitres du roman « VOYAGE AU CENTRE DE LA TERRE » de Jules Verne du français vers le chinois. Dans la Figure 20, nous présentons une capture d'écran de post-édition de ce roman.



Figure 20 : Capture d'écran de iMAG français→chinois pour « Voyage au centre de la terre »

Tableau 26 montre la statistique de la post-édition des 21 chapitres.

Tableau 26 : Statistique sur 21 chapitres de « Voyage au centre de la terre »

Chapitre	Segments	Mots	Mots/Seg (moyenne)	Page_std (250 mots)	Temps SECTra (secondes)	Temps/page_std en minutes.
CH1	76	1377	18,12	5,51	1009	3,10
CH2	104	1391	13,38	5,56	1565	4,69
CH3	44	592	13,45	2,37	742	5,22
CH4	106	1362	12,85	5,45	1477	4,52
CH 5-6	279	4160	14,91	16,64	3536	3,54
CH 7-9	399	6213	15,57	24,85	3534	2,37
CH 10-12	319	5154	16,16	20,62	4941	4,01
CH 13-15	254	5407	21,29	21,63	3190	2,46
CH16-18	333	4890	14,68	19,56	3252	2,77
CH19-21	202	2932	14,51	11,73	1942	2,76
Au total	2116	33478	15,49	133,91	27188	3,38

III.3.3 « The Book of Me » de Powers

L'expérimentation de post-édition de « THE BOOKS OF ME » de Richard Powers sur iMAG est pilotée par Laurent Besacier en 2014. Dans cette section, Nous montons la statistique, le résultat et l'évaluation. Citons ici (Besacier, 2014).

Corpus et statistiques de post-édition

L'œuvre, composée de 545 segments et 10731 mots est divisée en trois blocs identiques. Le Tableau 27 résume le nombre de mots des données source et cible [TA ou PE⁴⁹]. Sans surprises, un ratio supérieur à 1,2 est observé entre cible française [TA] et source anglaise. On constate cependant que ce ratio tend à diminuer après post-édition de la sortie française.

Tableau 27 : Corpus source, cible traduite et cible corrigée

Itération (nb. seg)	Anglais (nb. mots)	TA Français (nb. mots)	PE Français (nb. mots)
It.1 (184)	3593	4295	4013
It.2 (185)	3729	4593	4202
It.3 (176)	3409	4429	3912
Total (545)	10731	13317	12127

Le point de vue des lecteurs sur la traduction post-éditée

Neuf lecteurs ont accepté de lire l'œuvre traduite et ont répondu à un questionnaire, toujours ouvert sur fluidsurveys.com⁵⁰. La version pdf de l'essai traduit ainsi que le fichier tableur rassemblant les résultats du sondage sont également rendus disponibles dans github. Après trois questions permettant de mieux cerner le profil du lecteur, une première partie (5 questions) interroge les lecteurs sur la lisibilité et la qualité du texte littéraire traduit. Une seconde partie (7 questions) vérifie que certaines subtilités du texte ont été bien comprises.

Le point de vue du traducteur de R. Powers

Pour finir cette étude pilote, un dixième lecteur a été sollicité : le traducteur français de l'auteur, J-Y Pellegrin, enseignant chercheur à Paris-Sorbonne. Son avis est résumé ici sous la forme de questions-réponses. Le manque de place ne nous permet pas de commenter ces remarques mais nous pensons qu'elles sont assez explicites pour être délivrées en l'état.

Lisibilité ?

"Le texte auquel vous êtes parvenu restitue une image fidèle du contenu de l'article de Powers. Le pari de la lisibilité est gagné et certains passages (notamment ceux qui portent sur les aspects scientifiques de l'expérience décrite) sont très convaincants."

Imperfections ?

"Il reste bien sûr des imperfections, des lourdeurs, voire des erreurs ponctuelles, qui appellent une correction"

Principales erreurs ?

"Le défaut le plus répétitif, celui dont souffre d'ailleurs le travail de tout traducteur débutant, est le calque syntaxique, là où le français structure différemment la phrase .../... On comprend, mais ça ne sonne pas vraiment français"

"Autre défaut assez fréquent, la perte des idiomatismes du français au profit d'anglicismes. Parfois ces anglicismes peuvent être plus dérangeants lorsqu'ils

⁴⁹ La post-édition utilisée ici est obtenue après chaque itération du processus; la dernière étape de révision n'est donc pas prise en compte à ce stade.

⁵⁰ https://fluidsurveys.com/surveys/manuela-cristina/un-livre-sur-moi-qualite-de-la-traduction/?TEST_DATA=

flirtent avec le français comme dans « connaissances actionnables » (p. 18) au lieu de « connaissances pratiques / utilisables ». "

"Un troisième défaut tient à la non prise en compte de certains repères culturels.../...Par exemple, Powers fait plusieurs références à la topographie de Boston qui donnent lieu à des inexactitudes dans la traduction : « la rivière Charles » par exemple (p. 12) qui n'est pas une rivière mais plutôt un fleuve ; c'est pourquoi on traduira par « la Charles River » ou simplement « la Charles »"

Ce texte pourrait-il servir de base de départ à un traducteur littéraire professionnel ?"

"Instinctivement, je serais tenté de répondre non pour l'instant, parce que, dès son premier jet, le traducteur possède des réflexes qui lui permettent de produire un texte plus « propre » que celui auquel vous êtes parvenu .../... Cependant, ce traducteur passera plus de 25 heures à produire les 42 feuillets de 1500 signes correspondant au texte de Powers. À raison de 7 feuillets par jour en moyenne, il faut 6 journées de 8h pour venir à bout du texte .../... Si, en revanche, je pouvais ne travailler que sur votre texte, en oubliant complètement celui de Powers parce que j'aurais la garantie que votre traduction ne comporte aucune erreur, ni oubli, ni aplatissage par rapport à l'original, mais qu'elle demande simplement à être améliorée, rendue plus fluide, dans un français plus authentique, les choses seraient différentes et le gain de temps sans doute considérable."

III.3.4 « IITB : Monastery, Sanctuary, Laboratory » de Rohit Manchanda

Le livre retraçant cinq premières décennies de l'IIT-Bombay, « MONASTERY, SANCTUARY, LABORATORY », a été terminé et annoncé le 5 septembre 2008. Rédigé par Rohit Manchanda, et publié par Macmillan Inde, le livre retrace le parcours de l'IIT-Bombay depuis ses débuts jusqu'à présent sous la forme d'un récit historique.

EN 2012, à l'occasion d'un séjour à l'IIT Bombai professeur invité, C. Boitet a créé une IMAG pour traduire ce livre vers le hindi à l'aide de volontaires. Dans la Figure 21, on montre un chapitre étant traduit de l'anglais en hindi.



[WikiXan] in validation process

It's an arena an eagle wheeling up in the sky would see as a large bowl of landscape, amphitheatre-like, its northern rim as if chipped away to form the saddle the contingent has crested.

[Please suggest a better translation](#)

एम

कट्टर 10, 1959. बंबई के एक चौकी पूर्वोत्तर में पुरुषों के एक छोटे से टीम झड़व यह एक तेज छोड़ दिया ले, वे ऊपर इलाके की ओर में एक खड़ी वृद्धि झड़व इसके काठी, आर अचानक उन्हें पहल, अपने सही करने का। एक विशाल, शांत विस्ता, शहर के आद्यांगिक क्षेत्र में वे के माध्यम से चला गया हैं के विपरीत में एक अध्ययन materializes. यह एक क्षेत्र के आकाश में Wheeling ईगल परिदृश्य, अखाड़ा की तरह, इसके उत्तरी रिम के रूप में यदि दूर chipped हैं काठी फार्म आकस्मिक Crested हैं की एक बड़ी कटोरी के रूप में देखना होगा है. इसकी पृष्ठभूमि हावी एक फैला हुआ झील, अमीबा के आकार glistens, मार्च आकाश के शानदार उपर नीला दर्शाती, झील garlanding हैं पहाड़ियों का एक रोलिंग वर्धमान है. पूरे क्षेत्र के पेड़ के साथ बिंदीदार है. यह पानी, पहाड़ी, पत्तेदार लोकेल हमारी टीम पर देख रहे हैं अपने स्वयं के ठोक के बाद यह नाम जलाशय, पर्वत झील के साथ पर्वत, inland की मेंड है.

सहूलियत से अपने विशेषाधिकार प्राप्त ईगल तबे समय तक फैला मानव अँख से आगे देखने के लिए पर्वत झील से परे पहाड़ियों के कारण उत्तर यह एक और नोट कर सकते हैं, बड़ी झील, इलाके के एक दूसरे कटोरा में cupped. फिर भी आगे, इस 2 झील के उत्तर, विहार झील, घने झूठ, उष्णकटिबंधीय जंगलों, बंबई के शहरी हलचल से हटा दिया जन्मी सोच के बावजूद किया जा रहा है उसके उपनगरों के द्वारा सभी पक्षों पर घेरे में है.

दो झीलों के बीच, पर्वत और विहार, आकार भूमि का एक पथ के निहित है, इस हवाई दृश्य में लगभग बढ़या दो टूठदार पंजे के साथ एक केकड़े के शरीर की तरह. यह एन्क्लेव, एक विशाल 550 एकड़ जमीन है कि हाल ही में जब तक था अर्ध वन भूमि, हमारे यात्रियों के बारे में हासिल करने के लिए गतव्य है. अब कुछ समय के लिए, एन्क्लेव का एक हिस्सा गतिविधि के साथ रहती है. हर सौ मीटर या तो, नए भवनों उपर जा रहा है डाल रहे हैं, या किसी भी कीमत पर अपने आरंभ दिखाई दे रहे हैं मूलधार में खोदा जा रहा है, plinths जगह में smoothed कर. भरा, unmetalled सड़कों, विशेष रूप से दिन के लिए बाहर नोकदार, आड़ी - तिरछी रेखाएँ आधार. सड़क से एक किलोमीटर के बारे में, एक बड़े उत्सव *shamiana* एक समाशोधन में खड़ा किया है.

स्थल पर पहुंचने समूह पहले से ही यहाँ इकट्ठे भीड़ के लिए काफी उत्साह के लिए कारण हैं: घेरा इसकी संख्या में स्वतंत्र भारत के प्रथम प्रधानमंत्री जवाहर लाल

Figure 21 : Exemple de post-édition d'un chapitre de « Monastery, Sanctuary, Laboratory: 50 Years of IIT-Bombay » de Rohit Manchanda

Partie B Construction de systèmes de TA spécialisés à des sous-langages en français ↔ chinois

Introduction

Pour construire des systèmes de TA fr-zh adaptés aux besoins potentiels des clients de L&M (surtout EDF et RENAULT), nous avons d'abord cherché une mémoire de traductions pour au moins un des sous-langages envisagés, mais il n'y en avait aucune. Nous avons alors produit (par post-édition de TA produites par GT) une MT de bonne qualité de 9000 segments, à partir de laquelle nous avons construit un système initial de TA fr-zh, meilleur que GT sur ce sous-langage.

En contexte de recherche, nous avons d'autre part construit trois systèmes de TA dans le cadre du projet TABE-FC. Les données d'entraînement contiennent deux parties. Une partie a été extraite de segments parallèles trouvés à partir de sites boursiers, l'autre partie a été obtenue via la post-édition, par deux étudiantes chinoises, des prétraductions (produites par GT) de pages de sites boursiers publiés en français ou en anglais. Nous avons alors construit des systèmes de TA pour les couples où nous avons assez de données (français-chinois, fr-en, en-zh), et nous les utilisons pour accéder en chinois aux pages Web de ces sites via la plateforme SECTRA_w/IMAG. Nous avons également créé des systèmes chinois-français à partir de MT utilisées « à l'envers », car nous n'avions personne capable de post-éditer du chinois → français. Les résultats, comme on peut s'y attendre dans ce cas de figure, sont très mauvais pour la tâche de compréhension, et tout juste utilisables pour la post-édition par des francophones comprenant très bien le chinois.

Enfin, nous avons étudié les conditions dans lesquelles on peut utiliser avec profit la technique d'apprentissage incrémental proposée dans MOSES (Koehn et al., 2007). Pour des sites Web comme celui du LIG, il semble qu'il faille partir d'environ 30000 segments déjà post-édités (600K mots, ou 2400 pages), avec des incréments des 200 à 300 bisegments, et des « étapes » de 20 incréments.

Dans le Chapitre IV, nous passons revue système de TA français ↔ chinois en contexte industriel, et décrivons nos efforts pour en construction un pour L&M. Dans le Chapitre V, nous présentons la suite de cette action, mimée cette fois-ici en contexte de recherche en coopération avec l'université de Xiamen. Enfin, dans le Chapitre VI, nous mettrons que la technique d'apprentissage incrémentale (AI) de MOSES utilisée avec profit, au moins dans certains conditions.

Chapitre IV Revue des systèmes TA français ↔ chinois en contexte industriel

Résumé

Ce chapitre présente une revue des systèmes de TA, directs ou passant par l'anglais, permettant de traduire entre chinois et français, ainsi que certains besoins dont nous avons eu connaissance à l'occasion de notre activité dans L&M.

Introduction

Beaucoup d'entreprises ont besoin de systèmes de TA français↔chinois pour aider à traduire des documents français ou des documents chinois. Elles demandent que le système de TA soit privé, et que la traduction soit vraiment fiable (pour leur sous-langages), même si la fluidité n'est pas parfaite. Pour pouvoir arriver à cela, nous avons commencé par passer en revue l'histoire du développement des systèmes de TA concernant le chinois, depuis les tout premiers travaux en Chine (dès 1957). Nous avons ensuite comparé 4 systèmes de TA opérationnels et disponibles sur le Web ou sur des serveurs privés sur la paire de langues français↔chinois. Aucun de ces systèmes n'est satisfaisant dans ces contextes, car (1) ou bien l'information « *sort* » de l'entreprise, (2) ou bien les licences en intranet sont trop chères, et (3) de toutes façons, la qualité des TA brutes est jugée insuffisante, et il n'y a pas moyen de l'améliorer en spécialisant le système au sous-langage concerné.

Après avoir présenté la demande de grosses sociétés et un état de l'art des systèmes actuels dont on peut penser qu'ils pourraient être utilisés par ce type de société (ou d'organisme) de façon opérationnelle, nous décrivons les travaux réalisés à L&M pour construire des systèmes de TA pour le chinois basés sur Moses, utilisables à terme par des clients potentiels.

IV.1 Demande de grosses sociétés

Lingua et Machina édite une application Web de gestion des contenus multilingues en entreprise appelée LIBELLEX⁵¹. Cette plate-forme intègre divers outils d'aide à la traduction (concordances bilingues, outils d'extraction et de gestion de terminologies, mémoires de traductions, systèmes de traduction automatique et outils de gestion de projets de traduction).

Les entreprises clientes de L&M, comme EDF et RENAULT, ont des filiales en Chine. Il y a beaucoup d'échanges internes aux entreprises entre la France et la Chine, comme les courriers électroniques, les comptes rendus, et les rapports. Ces textes doivent être traduits (en français ou en chinois), et pas en anglais, langue dans laquelle ni les uns ni les autres ne sont à l'aise.

Ces textes, et en particulière les rapports, contiennent généralement beaucoup de termes spécialisés. Dans de tels contextes, recourir à des traducteurs humains n'est pas une option, même si on avait de très bonnes mémoires de traductions. Il faudrait en effet qu'ils soient excellents et compétents dans la terminologie des domaines concernés, et donc très chers. Mais, de toutes façons, on n'arriverait pas à satisfaire la demande de « *temps réel* », et on n'a en réalité pas besoin de traductions de qualité professionnelle. Il suffit qu'elles soient assez compréhensibles et fiables, et surtout qu'on puisse améliorer « *en ligne* » celles qui contiennent des contresens et de faux termes.

⁵¹ <http://www.libellex.fr/>

Pour certains textes quotidiens comme les courriers électroniques et les réunions à distance, on a simplement besoin d'une traduction « compréhensible » et en temps réel, mais pas d'une traduction parfaite.

Les entreprises clientes ont besoin d'une solution moins chère, plus rapide, et donc automatisée, mais les systèmes existants ne sont pas satisfaisants. Tout d'abord, l'information privilégiée de l'entreprise est considérée comme secrète, confidentielle. On ne peut donc pas utiliser les systèmes publics comme GT, parce que l'entreprise ne veut pas divulguer des informations propriétaires. Ensuite, un système de TA français-chinois commercial, comme Systran Enterprise Server, est très coûteux (15000€⁵² pour déployer une instance de SYSTRAN ENTERPRISE SERVER 7 chez un client). Enfin, la qualité de traduction pour la paire de langues français↔chinois n'est pas suffisante. Le résultat de la TA est souvent incompréhensible.

Pour obtenir un système de TA dont le résultat brut (non post-édité) est jugé comme « *suffisant* » par une entreprise cliente potentielle de L&M, la seule solution est de « *personnaliser* » son système de TA français-chinois, et cela quel que soit le paradigme de TA utilisé. En ce qui concerne L&M, le choix s'est porté vers le paradigme de la TA statistique, et sur le développement à l'aide de l'outil MOSES.

IV.2 État de l'art de la TA du chinois

La recherche sur la TA du chinois a commencé dans la deuxième moitié des années 1950. Les chercheurs se sont concentrés principalement sur les 4 paires de langues chinois↔anglais et chinois↔russe, mais il n'y pas eu alors de recherches sur les 2 paires français↔chinois. La première expérience concernant à la fois le français et le chinois a été faite par le professeur Feng Zhiwei en 1981-82 sur une maquette chinois→français/anglais/allemand/japonais en ARIANE-78 (Feng, 1981).

Les éditeurs de systèmes de TA qui proposent le couple français↔chinois utilisent l'approche dite du « pivot textuel », c'est-à-dire qu'ils appliquent successivement les « *paires* » français↔anglais et anglais↔chinois, en utilisant l'anglais comme "pivot textuel". Mais les résultats sont bien pires que ceux déjà jugés très peu satisfaisants obtenus avec l'anglais.

La section IV.2.1 propose un rapide historique de la recherche sur la TA du chinois, et dans la section IV.2.2, nous testons quatre systèmes de TA qui proposent la paire de langues français↔chinois.

IV.2.1 Historique

La Chine est le quatrième pays qui s'est lancé dans la traduction automatique (TA), à la suite des Etats-Unis (en 1951), du Royaume-Uni (1955), du Canada (Booth à Saskatoon, venant d'Angleterre) et de l'Union soviétique (1954, Lyapunov et Bagrinovskaya, Novosibirsk). En 1959, des chercheurs de l'Institut de technologie de l'informatique et de l'Institut de linguistique de l'Académie Chinoise des Sciences (ACS, 中国科学院, Chinese Academy of Sciences) menèrent la première expérience de traduction automatique en russe→chinois.

Le plus important des projets de TA chinois→anglais a commencé à l'Université d'État de l'Ohio en juillet 1961. C'était un projet de recherche mis en place avec le soutien de la « National Science Foundation » (NSF) sous la direction de William S. Y. Wang. Au milieu des années 1960, Wang rejoignit le groupe de Berkeley et continua sa recherche dans le cadre du projet POLA, toujours sur la TA chinois→anglais.

⁵² http://en.wikipedia.org/wiki/Comparison_of_machine_translation_applications

En même temps, le centre de recherche d'IBM à Yorktown travaillait sur un système chinois →anglais, fonctionnant sur les mêmes principes que le système russe →anglais, et utilisant aussi le tout nouveau et très fameux disque photocopique (King and Chang, 1963).

En 1972, l'Université chinoise de Hong-Kong (香港中文大学, Chinese University of Hong Kong) proposa le système CULT (Loh and Kong, 1979), développé pour traduire des textes mathématiques du chinois vers l'anglais. Les résultats étaient très bons, grâce à une importante pré-édition manuelle.

Au début des années 1980, de nombreux instituts et universités ont été impliqués dans des recherches sur la TA. L'Institut de Technologie de Harbin (ITH, 哈尔滨工业大学, Harbin Institute of Technology) et l'Université du Nord-Est (UNE, 东北大学, Northeastern University) ont commencé leur recherches sur des systèmes de TA chinois ↔ anglais au milieu des années 1980. De plus, l'Université de Nankin (南京大学, Nanjing University) a commencé une recherche sur la TA japonais ↔ chinois durant cette période.

Pendant les années 1980, il y a eu une très grande activité commerciale au Japon, et au moins 30 sociétés informatiques (comme FUJITSU, HITACHI, NEC, SHARP, TOSHIBA, etc.) ont développé des logiciels de TA pour japonais ↔ anglais. Certains, notamment FUJITSU (avec ATLAS-II dû à H. Uchida), ont développé de gros prototypes pour d'autres langues, dont le chinois, mais ne les ont pas commercialisés à l'époque, car seul le couple anglais-japonais semblait pouvoir apporter un ROI réel (« Return On Investment » ou « retour sur investissement »).

En 1987-93, beaucoup d'instituts et d'universités chinois ont aussi participé au projet « *Joint study for Multilingual Machine Translation* », qui a été financé par le gouvernement japonais⁵³, et visait à un objectif ambitieux : produire un système de TA à « pivot » sémantique de haute qualité sur le modèle d'ATLAS-II pour cinq langues asiatiques (japonais, chinois, thaï, malais et indonésien) et l'anglais.

Durant la période 1980-2005, tous les systèmes de TA adoptèrent l'approche "experte" à base de règles, tandis que certains d'entre eux utilisaient aussi une approche à base d'exemples comme complément.

Il y a trois systèmes très connus : (1) MT-IR-EC, un système de TA anglais → chinois pour traduire les titres et les catalogues des journaux, développé par le Research Institute of Post and Telecommunication Science, (2) KY-1, un système de TA anglais → chinois développé par l'académie des sciences militaires, qui est aussi le cœur du premier système commercial de TA, TRANSTAR (Dong, 1990), et (3) HUAJIAN, un système de TA chinois → anglais développé par HUAJIAN CO. LTD.

De 1995 à 2005, beaucoup d'équipes de recherche et d'entreprises ont publié des systèmes de TA, comme GAOLI (GAOLI CO. LTD), CCID (CCID GROUP), et KINGSOFT QUICK TRANSLATION (KINGSOFT CO. LTD).

La recherche sur la TA statistique a commencé en Chine à partir de 2004-2005. En 2006, cinq équipes de recherche (l'Institut de technologie de l'informatique de l'ACS, l'institut d'Automatique de l'ACS, l'Institut du logiciel de l'ACS, l'Université de Xiamen et ITH) ont publié un système de TA statistique en source ouvert, SILK ROAD (Silkroad, 2006). En 2011, UNE a publié un nouveau système de traduction automatique en source ouvert, NIUTRANS (Xiao et al., 2012).

⁵³ Essentiellement, par l'ODA (Overseas Development Agency) du METI.

Aujourd'hui, beaucoup de systèmes de TA proposent le couple français↔chinois, comme GT, SYSTRAN, REVERSO, BING, etc., mais il est difficile de trouver un système de TA français↔chinois de qualité d'usage correcte, sans doute car ils passent presque tous par l'anglais (« *pivot textuel* »).

IV.2.2 Expérimentations

IV.2.2.1 Systèmes étudiés

Nous avons fait des expériences sur quatre systèmes de TA français→chinois pour évaluer la qualité de TA. Nous avons d'abord testé GT, puis trois systèmes de TA statistique entraînés par les boîtes à outils disponibles en source ouverte (MOSES v2.0, JOSHUA v5.0, NIUTRANS v1.3).

Pour la construction des 3 systèmes de TA statistique, nous avons utilisé les mêmes données d'entraînement et les mêmes données de test. Ces données ont été extraites du corpus MULTIUN⁵⁴ (Eisele and Chen, 2010). Nous avons pris 1M de bisegments fr-zh comme données d'entraînement, 1K comme données de développement (*Tuning*), et aussi 1K comme données de test. Tous les textes chinois ont été segmentés par le segmenteur *Stanford*⁵⁵ (Chang et al., 2008), et nous avons utilisé le segmenteur de MOSES (un script perl *tokenizer.perl*) pour segmenter le texte français.

Tableau 28 : Statistique sur les données

	Nb de Segments	Nb de mots français	Nb de caractères français	Nb de caractère chinois
Données d'apprentissage	1M	28 951 255	199 486 574	35 109 257
Données de développement	1K	21 758	151 369	32 165
Données d'évaluation	1K	25 251	176 492	31 752

Notre expérimentation contient 2 étapes. Dans première, nous évaluons les systèmes de TA, par rapport avec le score BLEU. Dans la deuxième, nous utilisons SECTRA_W pour nous aider à évaluer les TA à la qualité d'usage du résultat et au degré d'automatisme.

La qualité d'usage et le degré d'automatisme sont obtenus à partir des temps de post-édition en utilisant les formules proposées par C. Boitet dans ses cours au NII (Boitet, 2009).

⁵⁴ <http://opus.lingfil.uu.se/MultiUN.php>

⁵⁵ <http://nlp.stanford.edu/software/segmenter.shtml>

Tableau 29 : Formule d'évaluation de l'automatisme et de la qualité d'un système de TA

1. Automaticity (MT module only): (taken from NII lecture notes by Boitet, 2009)

$$A = 1 - \frac{T(\text{post_édition_hum})}{T(\text{interaction_hum})}$$

Ex : A = 83,3% if first draft takes 1h per standard page (of 250 words or 1550 characters) and human interaction takes 10 mn/page (minutes per page).

2. Quality (wrt HT) in %

$$Q = 1 - (2/100 \times \frac{Tpe_{total}(for\ the\ task)}{Thum_{estim}(for\ the\ task)} \times Thum_{std_{page}}(mn))$$

Ex :

Q = 40% if $Tpe_{total} = 30mn/p$ (8/20)

Q = 60% if $Tpe_{total} = 20mn/p$ (12/20)

Q = 90% if $Tpe_{total} = 5mn/p$ (18/20)

$Thum_{estim}$ est le temps moyen de traduction humain par page, estimé dans la tâche en cours. Sans plus d'information, on l'évalue à 60 mn.

IV.2.2.2 Expérimentation avec GT

Depuis 2007, GT fournit un service de traduction français↔chinois. Dans cette expérience, nous avons traduit un texte du français vers le chinois en utilisant GT. Le texte en entrée source contient 1K phrases françaises du corpus MULTUN. Dans le Tableau 30, nous montrons 3 phrases extraites de ces 1K. Les références ont été produites par nous (par PE).

Tableau 30 : Exemple de traduction de GT

ID	Segment	Traduction de GT	Référence	Trace
1	Le Secrétaire général souhaite vivement que le plan-cadre d'équipement soit achevé d'ici à la mi-2014 et le nécessaire sera fait pour atteindre cet objectif, en contrôlant bien la portée du projet, en effectuant rapidement les réinstallations et en suivant de très près chaque activité de sorte qu'elle soit réalisée dans les délais prescrits.	秘书长真诚希望，基本建设总计划通过2014年中期完成，必要的工作将实现这一目标，控制项目的良好范围，迅速进行搬迁及以下非常接近每个活动，使其上进行的时间。	秘书长的坚定目标是，到2014年年中完成基本建设总计划，并将通过控制规模、加快搬迁和加紧监测每项活动的时间表，全力实现这一目标。	秘书长的秘书长坚定真诚目标希望是，基本建设到2014计划年年通过中2014完成年基本建设中期总完成计划，必要并的工作将实现通过这一目标一控制项目规模的、良好加快范围搬迁一和迅速加紧进行监测搬迁每及项以下活动非常的接近时间表每个活动，使全力其实现上这进行二的目标时间。
2	Au cours des six dernières années, la Commission s'est engagée dans des négociations complexes concernant ce projet, qui ne sera réalisé au mieux et dans les limites du budget impartie que si le financement est rapide.	在过去的六年里，该委员会已从事有关项目复杂的谈判，这将是最好的实现，并在规定的预算范围内，如果资金是快速的。	过去六年来，本委员会就该项目进行了复杂的谈判。这个项目必须及时获得资金，才能有效运作不超出预算。	过去的六年里来，该本委员会已就从事该有关项目进行了复杂的谈判一。这个将项目是必须最及时好获得的资金实现，并才在能规定有效的运作预算不范围超出内预算一如果资金是快速的。
3	Afin de garantir que la	为确保满足委员	为了确保委员会	为了为确保满足委员会的不

ID	Segment	Traduction de GT	Référence	Trace
	Commission réponde aux attentes et marque une réelle différence au Burundi, tous ceux qui sont impliqués doivent examiner la meilleure façon de soutenir la mise en œuvre des engagements identifiés en tenant compte des différences dans les capacités et l'expertise.	会的期望，并在布隆迪的一个真正的区别，所有参与应考虑如何最好地支持确定的承诺的执行情况，同时考虑到在能力和专业知识的差异。	不负众望，给布隆迪带来切实的变化，所有相关方都必须考虑如何才能最好地支持助承诺的执行，并且铭记各方在能力和专门知识方面的差异。	负众望期望，并给在布隆迪的带来一个切实真正的区别变化，所有参与相关应方都必须考虑如何才能最好地地支持助确定的承诺的执行情况，同时并且考虑铭记到各方在能力和专业专门知识方面的差异。

Le score BLEU de GT est 38.25%.

IV.2.2.3 Expérimentation avec Moses

MOSES (Koehn, Hoang et al., 2007) propose l'ensemble des outils nécessaires à la construction d'un modèle de traduction. Un décodeur permet aussi d'utiliser ces outils afin de produire la traduction d'un texte source. C'est un outil sous licence libre.

Tout d'abord, nous calculons des alignements de mots en utilisant GIZA++ (Och, 2003), qui implémente les algorithmes des modèles IBM 1-5 (Brown et al., 1993) et HMM (Vogel et al., 1996). On utilise les alignements pour construire la table de traductions. Enfin, un modèle de réordonnancement est construit, contenant les informations sur les positions dans les phrases des mots traduits par rapport aux mots traduits précédemment. Le modèle de langue est construit à l'aide de l'outil IRSTLM (Federico et al., 2008).

La construction de le système de TA MOSES a pris 15 heures (15h 32mn) pour finir la procédure d'entraînement (du prétraitement jusqu'à l'évaluation BLEU), Le score BLEU est 36,72%.

IV.2.2.4 Expérimentation avec Joshua

JOSHUA (Li et al., 2009) est un décodeur développé d'abord pour utiliser le modèle hiérarchique. Il est accompagné de l'ensemble des outils nécessaires à son fonctionnement : alignement (avec GIZA++), construction de la table de traductions, décodage, optimisation des poids, minimisation d'erreur, et calcul du modèle de langue cible. Depuis la version 6, il supporte le modèle à fragments (*chunks*). Il intègre un segmenteur du chinois, et on utilise donc son script *pipeline.pl*⁵⁶ pour entraîner le système de TA. JOSHUA demande d'écrire les paramètres dans un fichier de configuration. Les paramètres sont dans Tableau 31.

Tableau 31 : Paramètres de configuration de Joshua

<code>\$JOSHUA/scripts/training/pipeline.pl // le script de pipeline</code>
<code>--rundir 1M //le répertoire de travail</code>
<code>--source fr // langue source</code>
<code>--target zh // langue cible</code>
<code>--corpus 1M/train/train // données d'entraînement</code>
<code>--tune 1M/tune/tune // données de développement</code>
<code>--test 1M/test/test // données de test</code>
<code>--lm-order 5</code>
<code>--aligner giza</code>

⁵⁶ <http://joshua-decoder.org/6.0/pipeline.html>

La construction du système JOSHUA a pris 18 heures (18 h 12mn) pour la procédure d'entraînement (du prétraitement jusqu'à l'évaluation de BLEU). Le score BLEU est 32,18%.

IV.2.2.5 Expérimentation avec NiuTrans

NIUTRANS (Xiao, Zhu et al., 2012) est une boîte à outils en source ouvert permettant d'entraîner un système de TA statistique. Il est développé en C++ par l'UNE (东北大学). Actuellement, NiuTrans supporte déjà le modèle syntagmatique (PBMT) et le modèle hiérarchique.

Pour adapter des données à NIUTRANS, tout d'abord, on doit prétraiter les données. Nos données sont d'abord segmentées en mots (*tokenisation*) par le segmenteur STANFORD. Ensuite, on utilise les scripts perl, fournis par NIUTRANS, pour normaliser les segments⁵⁷. Enfin on produit le fichier *alignement.txt*⁵⁸.

Après la préparation des données, nous avons pris 10 heures (10 h 37 mn, sans compte le temps d'alignement) pour entraîner le système. Nous avons obtenu un score BLEU de 33,19%.

IV.2.2.6 Description du résultat

Pour la comparaison de la qualité de traduction des systèmes de TA, nous avons entraîné les systèmes de TA dans ces mêmes conditions (les outils, les données et le matériel, sauf pour GT). Nous prenons en compte le temps d'entraînement, et le score BLEU. Le système de TA construit avec MOSES a la meilleure qualité de traduction parmi les systèmes de TA. C'est sans doute, grâce au corpus d'entraînement adapté au même domaine qu'il est un peu mieux que GT.

Le score BLEU n'est en pratique pas bien corrélé à la qualité d'usage de la TA. Pour l'évaluer, nous ajoutons les segments source et les résultats de TA dans SECTRA_W, puis nous post-éditons et calculons la distance de post-édition. Dans le Tableau 32, nous montrons un exemple de résultat de TA, et nous pouvons voir la « Trace » du calcul de distance entre la référence et le résultat de TA.

Tableau 32 : Comparaison d'exemples de traductions obtenues par TA et d'une référence

Source	le Secrétaire général souhaite vivement que le plan - cadre d ' équipement soit achevé d ' ici à la mi - 2014 et le nécessaire sera fait pour atteindre cet objectif , en contrôlant bien la portée du projet , en effectuant rapidement les réinstallations et en suivant de très près chaque activité de sorte qu ' elle soit réalisée dans les délais prescrits .
Reference	秘书长的 坚定目标是， 到 2014 年 年 中 完成 基本建设 总 计划， 并 将 通过 控制 规模、 加快 搬迁 和 加紧 监测 每 项 活动 的 时间 表， 全 力 实 现 这 一 目 标。
Google	秘书长 真诚 希望， 基本建设 总 计划 通过 2014 年 中 期 完成， 必要 的 工作 将 实 现 这 一 目 标， 控制 项目 的 良好 范围， 迅速 进行 搬迁 及 以下 非常 接近 每 个 活动， 使 其 上 进 行 的 时 间。
Trace (PE vers la référence)	秘书长的 秘书长 坚定 真诚 目标 希望 是， 基本建设 到 总 计划 通过 2014 年 中 期 年 中 完成 一 基本建设 必要 总 的 计划 工作， 并 将 实 现 通过 这 一 目 标 一 控制 项目 规模 的、 良好 加快 范围 搬迁 一 和 迅速 加紧 进行 监测 搬迁 每 及 项 以下 活动 非常 的 接近 时间 表 每 个 活 动， 使 全 力 其 实 现 上 这 进 行 一 的 目 标 时 间
Moses	秘书长 真诚 希望 这 项 框 架 计 划 的 将 于 Mi- 2014 年 完成， 并 将 采 取 步 骤 实 现 这 一 目 标， 通过 控制 的 项目 的 范围， 迅速 重新 安 置 或 采 取 行 动， 不 影 响 的 每 项 具 体 活 动 进 行 规 定 时 限 内 提 出。
Trace (PE vers la référence)	秘书长 真诚 的 希望 坚定 这 项 目 标 框 架 是 计 划， 的 将 到 手 2014 年 中 期 年 中 完成， 中 并 完成 将 基本建设 采 取 总 步 骤 计 划 实 现， 这 并 一 将 目 标， 通过 控制 的 规 模 项、 目 的 加快 范围 搬迁， 和 迅速 加紧 重新 监 测 安 置 每 或 项 采 取 活 动 行 动， 不 影 响 的 每 时 间 表 项， 具 体 全 力 活 动 实 现 进 行 这 规 定 二 时 限 目 标 内 提 出。

⁵⁷ <http://www.nlplab.com/NiuPlan/NiuTrans.YourData.html>

⁵⁸ <http://www.nlplab.cn/NiuTrans.Phrase.html>

Joshua	秘书长真诚希望这项框架计划设备完成在这里,mi 2014年之后、和必要的工作将为实现这一目标,在控制两个项目的范围、执行快速的重新安置,通过密切关注每个活动,在规定时间内提出。
Trace (PE vers la référence)	秘书长真诚的希望这项框架坚定计划目标设备是完成,在到这里,mi 2014年之后年一中和完成必要基本建设的总工作计划将,为并实现将这通过一目标,在控制两规模个项目的范围、执行加快快速搬迁的和重新加紧安置监测,每项通过活动密切的关注时间表每个,活动全力,实现在这规定二时限目标内提出。
NiuTrans	秘书长强烈希望计划框架的设备或已完成的在2014年Mi必须将实现这一目标,通过检查影响控制两个项目的范围、执行的快速的重新安装又密切注视活动,实现的规定时限。
Trace (PE vers la référence)	秘书长强烈希望计划框架的设备坚定或目标已是完成,的到在2014年Mi年必须中将完成实现基本建设这总一计划目标,并将通过检查控制影响控制规模两、个加快项目搬迁的和范围一执行加紧的监测快速每项重新活动安装的又时间表密切,注视全力活动,实现的这规定二时限目标。

Dans l'Annexe 8, on donne 50 segments en « vue SECTra/Post-édition », montrant pour chaque segment le texte source, la PE, et les TA en mode « Trace ». Voici un exemple de résultat d'évaluation sur le segment présenté dans le Tableau 33.

Tableau 33 : Exemple de résultat d'évaluation

Mots	TA	TPE	TPE/p.std	DistPE ($\alpha=0,2$; $\beta=0,8$)	Q
53	Google	127s	16,0 mn	Dc: 112 ; Dw: 60 ; D= 70,4.	68%
	Moses	123s	15,5 mn	Dc: 126 ; Dw: 66 ; D= 78	69%
	Joshua	153s	19,2 mn	Dc: 137 Dw: 68 ; D= 81,8	61,6%
	NiuTrans	142s	17,8 mn	Dc: 117 Dw: 64 ; D= 70	64,4%

Conclusion

Nous avons comparé les 3 systèmes de TA et GT sur le BLEU et la qualité d'usage. Nous avons choisi le "meilleur" système, et il peut traduire les phrases du français vers le chinois. Mais la qualité de traduction n'est pas satisfaisante. La qualité est limitée par la taille de corpus, le domaine de traduction, le lexique, etc. Comment construire un système de TA français-chinois en haute qualité ? Pour nous, c'est un vrai défi.

IV.3 Construction de systèmes de TA pour le chinois basés sur Moses en contexte industriel

L'étude précédente nous a montré qu'il n'y avait pas pour l'instant de système français↔chinois pouvant être utilisé tel quel, ou adapté rapidement, pour satisfaire les besoins des grandes sociétés en général, et des clients potentiels de L&M en particulier.

Nous avons donc essayé de construire nous-même un système français-chinois à partir d'une MT correspondant aux besoins d'au moins un client potentiel de L&M. Malheureusement, aucun n'avait de telle MT. Nous en avons donc construit une, mais, faute de ressources (en post-éditeurs), nous n'avons pas pu dépasser 9000 bisegments (dans ce cas, 112500 mots ou 450 pages standard). Les résultats ont été encore pires qu'avec les 4 systèmes étudiés plus haut. Notre hypothèse est qu'il aurait fallu disposer d'une MT d'au moins 30000 à 50000 segments.

IV.3.1 Choix du sous-langage et des couples à traiter

Comme le client potentiel le plus prometteur pour L&M était EDF, nous avons cherché à construire un système pour EDF. C'est une très grosse entreprise spécialisée dans le domaine

de l'énergie électrique, qu'elle soit produite dans des centrales nucléaires, hydrauliques, à charbon, à gaz, éoliennes, ou photovoltaïque.

Depuis une quarantaine d'années, EDF est implantée en Chine, où elle a construit des centrales nucléaires, et travaillé avec d'innombrables cadres, ouvriers et ingénieurs chinois. Beaucoup de documents ont été traduits, dans les deux sens, et nous espérons avoir accès à des documents parallèles, ou au moins à de grosses mémoires de traductions, en supposant que des outils comme SDL TRADOS ou DEJA VU avaient été utilisés pour produire ces traductions.

Nous comptons bien sûr choisir comme sous-langage objet du système de TA à construire celui correspondant à la MT la meilleure en qualité et la plus grande en volume. Malheureusement, nous n'avons rien pu obtenir du tout. Peut-être ces textes parallèles ou ces MT existent-ils et sont-ils cachés, peut-être n'ont-ils jamais été créés, nous n'en savons rien. Sachant que la traduction n'est presque jamais consolidée dans les comptes des entreprises, et est le plus souvent sous-traitée de manière opportuniste, il est possible que la seconde hypothèse soit la bonne.

Quoi qu'il en soit, en 2013, L&M n'avait pas pu avoir accès à un corpus parallèle ou à une MT d'EDF français→chinois permettant de développer un système de TA, qu'il s'agisse d'un système MOSES (il aurait fallu entre 20K et 30K "bons" bisegments) ou d'un système ARIANE (à règles et dictionnaires), pour lequel il aurait fallu des corpus parallèles ou comparables de 2K à 3K segments (pour l'étude typologique) et un dictionnaire bilingue de 10K à 20K entrées.

Nous avons alors décidé d'essayer de construire un système MOSES à partir d'un corpus parallèle que nous construirions nous-même à partir d'un corpus bilingue le plus "adapté" ou "vraisemblable" possible, puis, s'il était trop petit, à partir d'un corpus monolingue complémentaire que nous traduirions.

Nous savions que les besoins d'EDF étaient dans les deux sens (français↔chinois). Nous nous sommes concentré sur le sens français→chinois, car nous savions que nous pourrions nous-même évaluer et post-éditer les résultats, alors que nous n'avions personne dans notre environnement qui comprenne bien le chinois technique et soit de langue maternelle française.

Cependant, nous avons aussi fait quelques essais en chinois→français, "pour voir", en nous disant que, si les résultats étaient encourageants, nous pourrions peut-être écrire des parties de cette thèse en chinois, les faire traduire par GT ou par notre système, et les faire ensuite réviser par des chercheurs du laboratoire, compensant leur ignorance du chinois par leur connaissance du domaine. Comme on pouvait s'y attendre, cet espoir a été totalement déçu, et nous ne nous étendrons pas sur cet essai.

IV.3.2 Recherche infructueuse de corpus parallèles adaptés

Notre première idée a été d'essayer d'extraire un corpus bilingue français→chinois concernant un des domaines d'EDF à partir des corpus parallèles librement disponibles sur le Web.

La performance d'un système de traduction automatique statistique (TAS) (Koehn, 2009) dépend fortement de la taille et de la qualité du corpus parallèle utilisé pour l'entraînement. Les ressources actuelles en corpus parallèles bilingues ou multilingues libres de droits proviennent généralement d'institutions internationales. C'est le cas du corpus « EUROPARL » (Koehn, 2005) extrait des délibérations du Parlement européen, du corpus « CANADIAN HANSARDS », contenant les transcriptions en français et en anglais des débats du Parlement canadien, et du corpus JRC-ACQUIS qui fournit une quantité comparable de textes législatifs européens en 22 langues (Steinberger et al., 2006). En ce qui concerne le chinois, il y a

beaucoup de corpus parallèles anglais→chinois et chinois→anglais, notamment celui du journal XINHUA NEWS (Graff et al., 2003), mais très peu de corpus français→chinois.

Or, pour construire un système de TAS français→chinois, il faut disposer d'un corpus parallèle français→chinois (dans le bon sens) pour entraîner les modèles. Nous nous sommes tourné vers le corpus parallèle MULTIUN français-chinois, qui a été construit par extraction du site Web des Nations-Unies, puis nettoyé et converti au format XML par Andreas Eisele et Yu Chen (Eisele and Chen, 2010) en 2010. En février 2013, la version alignée de ce corpus a été publiée sur le site Web OPUS⁵⁹ (Tiedemann, 2012). Voici les corpus que nous avons pu collecter à ce point (Tableau 34).

Tableau 34 : Corpus collectés en cherchant des corpus pour le français→chinois

Nom	Direction	Nb segments	Nb mots source	Mots fr / segment	caractères zh / segment	Mots fr / caractères zh
MultiUN	en-zh	8,8M	220,4M	24,97	71,31	285,54%
MultiUN	fr-zh	8,7M	243,8M	27,94	71,94	257,33%

Que pouvait-on en espérer ? *A priori*, peu, car, quand on entraîne un système de TAS avec un tel corpus "généraliste" pour la traduction dans un domaine précis, par exemple l'énergie, on obtient d'habitude de mauvais résultats. Nous avons fait l'expérience, qui a confirmé cette crainte.

Nous en avons conclu qu'un système de TA français-chinois entraîné seulement avec le corpus MULTIUN ne pourrait pas répondre aux besoins des entreprises clientes de L&M. En nous inspirant de publications mentionnant la possibilité de mélanger un petit corpus spécialisé à un grand corpus généraliste, nous avons alors décidé d'essayer cela, et de construire un corpus parallèle français-chinois spécialisé au sous-langage des notes techniques et des courriels concernant le domaine de l'énergie.

Pour cela, nous sommes parti des sites Web d'EDF en France et en Chine, car nous avons remarqué qu'ils contiennent beaucoup de segments français et chinois presque parallèles.

Nous avons extrait des textes français et chinois à partir de ces sites Web à l'aide de l'outil BOILERPIPE (Kohlschütter et al., 2010). Mais ces textes ne peuvent pas être utilisés tels quels pour entraîner un système de TA, il faut d'abord les segmenter, les nettoyer, les aligner, et enfin extraire des bisegments réellement parallèles (en relation de traduction).

Tout d'abord, nous avons débruité les textes, en supprimant les segments inutiles comme les liens (par exemple, <http://...>), les chiffres, les dates, etc.

Ensuite, nous avons normalisé l'encodage des caractères en transformant tout en UTF-8 (certains textes français étaient en codage ASCII MAC ou WINDOWS, et les textes chinois étaient le plus souvent en GB-2312-80).

Après avoir nettoyé le "bruit", nous avons procédé à l'alignement au niveau des segments (phrases ou titres) en utilisant l'outil LF ALIGNER⁶⁰. Nous avons finalement obtenu un corpus parallèle d'environ 3K bisegments. Un extrait en est donné à l'Annexe 9.

Cette petite quantité de données n'est pas suffisante pour entraîner un système Moses, mais, mélangée à ce que nous avons extrait de MULTIUN, elle a suffi pour améliorer un peu le

⁵⁹ <http://opus.lingfil.uu.se/MultiUN.php>

⁶⁰ <http://sourceforge.net/projects/aligner/>

système de TA, en lui faisant apprendre des termes comme « *le noyau des atomes* », « *l'hydraulique* », « *le charbon propre* », etc.

Remarque : Dans son rapport, le professeur Xiaodong SHI nous a proposé une méthode d'analogie (Lü et al., 2007) pour augmenter la quantité des données d'entraînement. L'idée était de voir si 9K bisegments déjà produits pouvaient être exploités pour extraire une quantité significative de segments similaires à partir du corpus MultiUN, bien que nos segments post-édités et le corpus MultiUN soient dans deux domaines très différents.

Nous avons essayé de le faire, mais le résultat a été mauvais. L'Annexe 10 montre un exemple des segments extraits qui contiennent certains mots-clés.

IV.3.3 Production de corpus par PE de résultats de Google

Pour augmenter la quantité des données d'entraînement, nous avons construit une MT par post-édition des sorties de la 1^{ère} version de notre système. C'est un bon moyen de produire des MT de bonne qualité (Wang and Boitet, 2013).

Pour post-éditer plus vite et mieux, nous avons utilisé la plate-forme SECTRA_w/IMAG. Nous avons divisé notre MT en domaines, comme les nouvelles, les reportages et les pages Web de Wikipédia. Nous avons transformé tous les textes monolingues en des fichiers html, nous les avons placés dans une hiérarchie de fichiers mise sur le serveur du laboratoire, et nous avons créé une iMAG dédiée à ces fichiers.

Les prétraductions ont été fournies par GT. Après la post-édition, nous avons sélectionné, pour construire notre MT, les bisegments que nous estimions adaptés à notre besoin. Cette sélection était basée sur le niveau de fiabilité (d'une étoile "☆" à cinq étoiles "☆☆☆☆☆") et sur la note de qualité (de 0 à 20) associés à chaque segment et à chaque langue cible dans la MT. Le prédicat de sélection était:

(fiab = 3 && score ≥ 12) || (fiab = 4 && score ≥ 11) || (fiab = 5 && score ≥ 10)

Nous avons finalement obtenu 6000 segments parallèles. En les ajoutant aux 3000 segments parallèles extraits à partir des sites Web d'EDF, nous avons au total collecté 9000 bisegments (environ 450 pages standard) de qualité suffisante pour construire un système. Cependant, nous n'avons pas de garantie que ce "noyau" représente bien le sous-langage, qui pour nous reste inconnu et inconnaisable, des notes techniques d'EDF sur l'énergie (Tableau 35).

Tableau 35 : Exemples de bisegments français→chinois parmi les 9000 collectés ou produits

Français	Chinois
Charbon propre	清洁煤发电
EDF Asie	EDF 亚洲
Activités	业务概览
Charbon propre	洁净煤
En Chine, le charbon représente près de 80% de la production d'électricité et devrait continuer d'occuper une place majoritaire dans l'avenir (plus de 60 % à l'horizon 2020).	中国煤电约占全国总发电量的 80%，今后还会继续占有主导地位（预计 2020 年占 60%以上）。
Pour limiter les impacts sur l'environnement, la Chine développe des centrales à charbon à haut rendement moins polluantes. En s'appuyant sur ses compétences d'ingénierie, EDF prend part à ces projets. Ils permettent au Groupe de consolider et de développer son expérience pour faire face aux besoins qui pourraient émerger en Europe dans l'avenir.	为了减轻煤电对环境的影响，中国致力于发展高效、低污染的燃煤电厂。法国电力集团以专业技能为依托，参与中国洁净煤火电项目。通过参与项目建设，法国电力集团将巩固和发展其火电技术，应对欧洲未来可能出现的需求。
EDF a signé plusieurs accords de coopération avec des	法国电力集团已与国电、三峡集团、大唐等多家国

Français	Chinois
producteurs nationaux d'électricité, portant sur le développement conjoint de projets électriques, par exemple les Groupes de Trois Gorges, Guodian, Datang, etc.	有大型电力公司签署了多项电力合作协议。
French Investment Guangxi Laibin Electric power Co (FIGLEC) - Chine est une filiale à 100 % du groupe EDF. La société est propriétaire de la centrale de Laibin B (d'une puissance de 720 MW), exploitée par SYNERGIE, aussi filiale d'EDF.	广西来宾法资发电有限公司是法国电力集团的全资子公司，拥有两台单机容量为 360 兆瓦的机组，总装机容量 720 兆瓦。

IV.3.4 Construction de systèmes français → chinois

IV.3.4.1 Composants

Nos systèmes de TA sont construits avec MOSES, qui fournit des outils optimisés pour réaliser l'entraînement, mais qui ne contient pas d'outil pour traiter le chinois. À L&M, nous disposions aussi de la boîte à outils MYRIAM, qui intègre le segmenteur du chinois de XELDA⁶¹ (Xerox Linguistic Development Architecture), et un programme java pour normaliser les phrases. Pour l'alignement des mots, nous avons utilisé l'outil MGIZA (Gao and Vogel, 2008), qui propose une implémentation efficace et parallèle de GIZA++.

IV.3.4.2 Paramétrisation et construction de 2 modèles de TAS

Les paramètres de ces systèmes ont été optimisés de manière usuelle avec l'outil MERT (Minimum Error Rate Training) (Och, 2003). Les traductions produites sont évaluées avec la mesure BLEU (Bilingual Evaluation Understudy) (Papineni et al., 2002), ainsi qu'avec notre « distance mixte de PE » (MTER), qui est bien corrélée à la qualité d'usage.

Nous avons ainsi défini deux modèles. L'un est un modèle générique qui est entraîné avec le corpus parallèle MULTIUN, et l'autre est un modèle spécifique entraîné avec notre mémoire de bisegments post-édités. Quand on a de nouveaux segments post-édités, on n'a pas besoin de réentraîner le modèle avec toutes les données. On met seulement à jour le modèle spécifique. C'est très utile dans le cas où l'on a un très gros modèle générique qu'on souhaite utiliser dans un nouveau système, sans avoir à le réentraîner.

Tableau 36 : Comparaison des temps d'entraînement de Moses

Système de TA	Quantité de données	Temps d'entraînement	Temps du réentraînement avec 10K nouveaux bisegments
Un modèle	2M+9K bisegments	17 h 35mn	19 h 21mn
Modèle générique	2M+9K bisegments	13 h 10mn (10h30+2h40)	2 h 40mn
Modèle spécifique			

Nous avons ajouté 10K nouveaux bisegments, puis comparé les temps du réentraînement. Notre stratégie est de gagner 86,2% du temps sur le réentraînement avec 10K nouveaux segments. Le Tableau 37 présente les caractéristiques du serveur utilisé pour l'expérimentation.

⁶¹ <http://www.xrce.xerox.com/About-XRCE/History/Historical-projects/XeLDA>

Tableau 37 : Configuration de la machine

Nombre de processeurs	4
Nombre de cœurs par processeur	2
Thread	4
Mémoire	8G

IV.3.4.3 Choix des données pour le modèle général

Le corpus MultiUN français-chinois contient 9,7 M de phrases parallèles, composé d'environ 300 millions de mots français et environ 600 millions de caractères chinois, soit environ 315M mots du dictionnaire). Ce corpus contient des segments composés de chiffres, de numéros, ou encore de dates ; ces segments ne sont pas utiles pour notre système. Nous les avons donc éliminés. À la fin, il nous restait environs 8,3 M bisegments.

Nous n'avons pas utilisé tout le corpus pour entraîner le modèle générique, parce que nous avons d'abord fait une expérience pour trouver la quantité appropriée des segments. Nous avons testé l'entraînement du système Moses français-chinois avec des quantités différentes (1M, 2M, 3M, 4M, et 5M) de bisegments, et calculé le score BLEU pour chacune de ces quantités. Les résultats sont présentés dans le tableau ci-dessous (Tableau 38). Les scores de 2M vers 5M sont égale les mêmes. On prend en compte le temps d'entraînement, et on choisit 2M bisegments pour entraîner le modèle générique.

Tableau 38 : Scores BLEU pour différentes tailles du corpus d'entraînement

Nombre de bisegments	BLEU	Temps d'entraînement
1M	49,78%	16 h
2M	52,48%	18 h
3M	52,25%	23 h
4M	54,31%	32 h
5M	55,52%	42 h

Nous avons prétraité les bisegments (le corpus MULTIUN et la MT) en utilisant MYRIAM. Plus précisément,

- nous avons d'abord normalisé les segments (en convertissant les entités HTML, en séparant les ponctuations, et en ôtant la casse).
- nous avons ensuite utilisé XELDA pour marquer les segments français et chinois.

IV.3.4.4 Construction de systèmes de TA

Nous avons entraîné le modèle général construit uniquement à partir du corpus MultiUN (2M). Le tableau ci-dessous montre quelques chiffres concernant ce corpus.

Tableau 39 : Statistiques sur le corpus MultiUN

	Français	Chinois
Segment	2M	
Mots	58M	50M
Caractères	398M	90M
Pages_std	38K	226K

Le modèle de traduction est entraîné sur 2M segments tirés de MULTIUN. La partie cible du corpus a été utilisée pour produire un modèle de langue. Dans cette partie de l'entraînement, nous avons sauté l'étape usuelle d'optimisation des poids (tuning). Elle a été utilisée plus tard, pour construire un autre système, en association avec des modèles plus spécialisés.

Nous avons entraîné le modèle spécifique avec notre MT (9000 bisegments français-chinois), et collecté 1,7M segments monolingues chinois pour produire le modèle de langue cible.

Les paramètres utilisés ont été les mêmes que ceux utilisés pour l'entraînement du modèle général, sauf que nous avons utilisé la phase d'optimisation des poids (tuning).

Pour construire un système d'essai en chinois→français, nous avons utilisé les mêmes données à l'envers, en faisant comme si, étant donné un couple (\$fr, \$zh) de segments tels que \$zh=trad(\$fr), on avait aussi \$fr=trad(\$zh)). Nous savons pertinemment que c'est faux (la relation de traduction entre phrases n'est pas symétrique), mais... c'était "mieux que rien", puisque nous n'avions personne qui puisse post-éditer des résultats de TA zh→fr.

IV.3.5 Évaluations et perspectives

Pour évaluer la progression de notre système de TA "vers une cible" (pour chaque segment post-édité, la cible est sa post-édition), nous avons calculé le score BLEU pour chaque système de TA. Certes, BLEU ne mesure pas la "qualité" (ni au sens de qualité linguistique, ni au sens de qualité d'usage), et *ne peut pas* la mesurer, comme cela a été montré dans le fameux article d'Osborne, Callison-Burch et Koehn (Callison-Burch et al., 2006) "*Re-evaluating the Role of BLEU in Machine Translation Research*". Mais BLEU exprime bien une similarité textuelle, et peut être "raisonnablement" utilisé pour évaluer la progression d'un système vers une certaine "cible". Pour cela, il vaut d'ailleurs mieux n'avoir qu'une seule cible ("*traduction de référence*") par segment source, plutôt que 5, 10 ou 15 dans certaines campagnes d'évaluation. Notons que BLEU donne une mesure globale et ne donne vraiment pas d'indication fiable au niveau des segments individuels.

Nous évaluons aussi la qualité d'usage à partir du temps de post-édition (en minute par page, mn/p). Comme SECTRA associe un chronomètre à chaque segment, nous disposons du *temps primaire de post-édition*, T_{pe_1} , pour chaque segment édité à travers cette interface. Nous disposons aussi du *temps total de post-édition*, T_{pe_tot} , pour des sessions de post-édition⁶². Nous pouvons en déduire, pour chaque segment, le temps total T_{pe_tot} ainsi que le *temps secondaire de post-édition*, T_{pe_2} , qui correspond au temps passé à chercher des équivalents dans les lexiques et bases terminologiques, ou à communiquer avec d'autres personnes pour trouver une bonne traduction d'une expression "hors dictionnaire". Typiquement, ce temps représente les 2/3 ou les 3/4 du temps total en contexte de traduction professionnelle.

Au début, nous avons post-édité 300 segments (Voir un exemple des données de test dans le Tableau 40), du français vers le chinois (à partir des prétraductions fournies par GT), et les avons utilisées comme données de test pour notre évaluation.

Tableau 40 : Exemple de données de test

ID	Segment	Traduction de Google	PE par humaine	Trace
26	Nucléaire : la Chine adopte l'EPR	核电：中国采用EPR	核电：中国采用欧洲压水堆技术	核电：中国采用 EPR 欧洲压水堆技术
27	Au Laos, le projet de centrale hydraulique Nam Theun 2 (1070 MW) est porté par la société de projet Nam Theun 2 Power Company (NTPC), dont le	在老挝，液压动力项目南屯2（1070兆瓦）支持的项目公司南屯2电力公司（NTPC），法国电力集团与	在老挝，中央液压草案南屯2号项目（1070兆瓦）的支持，该项目公司南屯2电力公司（	在老挝， 中央 液压 动力 草案 项目 南屯2号项目（1070兆瓦） 的 支持 的 该 项目 公司南屯2电力公司（NTPC），法国电力集团与

⁶² Quand nous et d'autres Chinois qui nous aident post-éditons, nous notons l'heure au début et à la fin d'une session, et aussi l'ensemble des segments post-édités durant la session. Nous en tirons un temps total moyen. Nous faisons l'hypothèse (qui semble vérifiée) que les temps sont proportionnels aux nombres de mots, et une simple règle de trois nous donne alors T_{pe_tot} pour chaque segment.

ID	Segment	Traduction de Google	PE par humaine	Trace
	groupe EDF est le premier actionnaire avec 40 % des parts.	40% 的股权的第一大股东。	NTPC)，法国电力集团是拥有40%股权的最大股东。	是拥有40%的股权的第一大股东。
28	D'une capacité de 715 MW, la centrale a été mise en service en février 2005. Elle bénéficie des technologies éprouvées des turbines les plus récentes, ainsi que des derniers retours d'expérience des centrales à gaz construites par EDF qui en a assuré la construction et la livraison « clé en main » et qui participe maintenant à son exploitation.	容量为715兆瓦，该厂已于2005年2月它已被证明的经验，最后返回建造的EDF最新的涡轮机和燃气电厂保证了技术建设和交付“交钥匙”，现在参与其运作。	该厂容量为715兆瓦，2005年2月初投入运行。这一项目的汽轮机采用了最新经过验证的技术，吸取了EDF燃气机组最新的反馈经验，以“交钥匙”模式承担工程的建设，和参与电厂的运行。	该厂容量为715兆瓦，该厂2005年于2005年月初2月投入它运行已。被这证明二项目的经验汽轮机一采用最后了返回最新建造经过验证的EDF技术最，新吸取的了涡轮机EDF和燃气电厂机组保证最新子的技术反馈建设经验和，交付以“交钥匙”模式承担工程的建设，现在和参与其电厂运作的运行。
29	Le développement de l'énergie nucléaire est un enjeu majeur pour la Chine et le reste du monde dans le cadre de la préservation de l'environnement et de la réduction de l'effet de serre.	核电的发展是中国和重大问题世界保护环境，减少温室效应的范围内。	发展核能对中国和世界都具有非常重要的意义，是保护环境和减缓温室效应重要的途径。	发展核能核电对的中国发展和是世界中国都和具有重大非常问题重要世界的意义，是保护环境一和减少减缓温室效应重要的范围途径内。

Le Tableau 41 présente les statistiques des données de test. Ensuite, nous avons continué à post-éditer des résultats de TA, mais plus ceux de GT : nous avons continué à mettre dans notre MT les résultats de GT, mais nous avons post-édité les résultats de nos systèmes, puis réinjecté ces nouvelles cibles comme des « références » dans le processus d'apprentissage, etc.

Tableau 41 : Statistiques des données de test

Nb de segments	Nb de mots Par segment (source)	Nb de p.std (source)	Nb de caractères par segment (cible)	Nb de p.std (cible)	Tpe_p.std
300	26,1	31,32	28,3	21,2	6,2 mn

Nous traduisons les segments source avec 3 systèmes, GT, système de TA entraîné avec le corpus MULTIUN (2M bisegments), et notre système combiné. Les segments post-édités sont utilisés comme les références. Les scores BLEU des systèmes de TA sont montrés dans le Tableau 42, et avec un exemple de traduction.

Tableau 42 : Score BLEU et exemples de sorties de systèmes de TA

Systeme	Temps de PE (mn/pstd)	Source : Être un leader du renouveau du nucléaire dans le monde Référence : 成为全世界核能复兴的领导者	
GT	22	Traduction	作为世界核复兴的领导者
		Trace	成为全世界作为核能在世界上的核复兴的领导者
Système de TA (MultiUN)	28	Traduction	受教育复兴和体面工作问题的一个世界核
		Trace	成为全世界受核能教育复兴和体面工作问题的领导者一个世界核
Système de TA (MultiUN+MT)	24	Traduction	成为一个全球核电复兴的领导者
		Trace	成为全世界一个核能全球核电复兴的领导者

Conclusion

Nous avons essayé de créer un système de TA français→chinois en utilisant le corpus MULTIUN et la MT, traduisant le contenu dans le domaine de l'énergie.

Au début, nous nous sommes limité à ces ressources (très peu de corpus français-chinois, en particulier, le corpus français-chinois adapté au domaine) pour construire un système de TA statistique avec MOSES. La traduction obtenue est totalement « incompréhensible ».

Ensuite, nous avons testé les systèmes existants comme GT, mais le résultat n'était pas non plus satisfaisant. Pour construire un système ayant une qualité de traduction acceptable, nous avons commencé à construire un corpus parallèle spécialisé à notre domaine. Nous avons collecté et extrait des segments parallèles à partir de sites Web, mais la quantité de segments parallèles n'était toujours pas suffisante. Nous avons alors créé des IMAG pour des sites Web correspondant à notre sujet, et nous les avons post-édités. Les résultats de TA ont été fournis d'abord par GT, puis par notre système de TA basé sur MOSES, en construction. Nous avons obtenu plusieurs milliers de segments parallèles.

Enfin, nous avons construit un système de TA français→chinois avec un corpus "mixé" (corpus extrait et MT), et la qualité de traduction s'est enfin acceptable. Notre hypothèse est que nous pourrions arriver à une qualité vraiment bonne si nous avions non pas 9000, mais entre 30000 et 50000 segments. Nous n'avons pas eu les ressources suffisantes pour la tester, et espérons trouver une situation le permettant dans le futur.

Chapitre V Construction de systèmes de TA pour le chinois avec Moses en contexte de recherche : le projet TAFE-FC

Introduction

Au laboratoire, notre recherche sur la TA fr-zh a été surtout menée dans le cadre du projet TAFE-FC, défini et commencé grâce à un séjour sabbatique d'un an du Pr Yidong Chen de XMU. Un des buts de ce projet est de pouvoir comparer différents « paradigmes » de TA, sans se limiter aux systèmes purement « MOSES ». Il s'agit aussi de systèmes « à règles » comme NEON (de XMU), ou de systèmes empiriques utilisant des statistiques ainsi que des informations sémantiques et pragmatiques (architecture du Pr. Yidong Chen (Chen et al., 2014)).

Pour construire des systèmes de TA français-chinois de bonne qualité en contexte de recherche, nous étions libre de sélectionner des sous-langages très restreints (comme des « brèves » de sites boursiers), ou pour lesquels nous aurions déjà un corpus parallèle de bonne qualité (comme pour le site du LIG). Nous avons travaillé sur les deux sous-langages mentionnés ci-dessus, et présentons ce qui a été fait pour celui des sites boursiers et économiques.

Après avoir extrait environ 1000 pages Web de ces sites, nous avons extrait environ 3000 bisegments, puis les avons filtrés et avons obtenu une MT de bonne qualité. La quantité de segments parallèles n'étant pas suffisante pour entraîner un système MOSES, nous avons créé 3 IMAG pour 3 sites Web boursiers, et nous avons post-édité les segments prétraduits par GT, grâce à la participation de 4 stagiaires chinois étudiants à l'UJF. Cela fait, nous avons pu construire des systèmes Moses, et les avons évalués (distance de post-édition, temps, ainsi que BLEU et NIST). Nous sommes arrivés à des résultats compréhensibles, mais pas encore très fiables.

V.1 Buts du projet TAFE-FC

Le projet TAFE-FC est un projet collaboratif mené par l'Université de Grenoble, LIG-GETALP et l'Université de Xiamen (XMU, 厦门大学, Xiamen University), et visant à créer des instances d'un nouveau type de système de traduction automatique statistique utilisant des ressources lexico-sémantiques et discursives.

V.1.1 Buts théoriques

Le but concret est de développer des systèmes de TAS français↔chinois pour des sites boursiers et économiques. Comme très peu de corpus ou de dictionnaires bilingues français↔chinois sont disponibles sur Internet, l'anglais est utilisé comme « pivot » pour construire les équivalents français↔chinois par transitivité. Outre la description générale de ce projet, nous décrivons les progrès sur deux axes de recherche liés à ce projet. Pour cela, nous utilisons une méthode, proposée par XMU, d'induction de probabilité fondée sur la similarité thématique, qui produit des tables de traduction français-chinois à partir de tables de traduction français-anglais et anglais-chinois. Pour disposer de bons corpus parallèles français-chinois, nous utilisons un système Web de post-édition collaborative qui peut déclencher l'amélioration incrémentale du système de TA en utilisant des mesures d'évaluation de TA et en extrayant la "meilleure partie" de la mémoire de traductions courante.

V.1.2 Buts pratiques

Ce projet se concentre sur deux objectifs pratiques. Nous prenons qu'une utilisation opérationnelle permettrait d'obtenir une excellente qualité. Le premier est de construire un service de systèmes de TA français↔chinois. Tout d'abord, nous nous restreignons à des sous-langages observés, puis nous récupérons et construisons la ressource français↔chinois, et créons finalement un système de TA.

À cause de la limitation des ressources français-chinois, nous essayons de construire le système de TA français-chinois via une langue pivot (anglais) pour faire l'extraction de la traduction équivalent.

Le deuxième objectif est intégré notre système de TA dans SECTRA_w/IMAG (via TRADOH). Pour cela, nous créons des IMAG dédiées pour les sites Web boursiers et économiques. Notre plateforme permet d'accéder aux sites Web en multilingue, et fournit déjà une traduction français→chinois de bonne qualité.

V.1.3 Définition du projet

Nous avons organisé ce projet avec une architecture à trois niveaux, contenant sept « missions » (Figure 22). Tout d'abord, nous travaillons pour la construction de ressources français-chinois orientées vers l'économie (corpus parallèles, ressources sémantiques bilingues, banques d'annotations discursives, etc.) en vue du développement. Ensuite, nous construisons un nouveau modèle de TA basé sur la sémantique et l'information de discours (niveau de la recherche fondamentale). Enfin, nous prévoyons de construire un système Web de TA, et de l'appeler par SECTRA_w/IMAG ou par JIANDAN-EVAL pour améliorer la qualité de traduction des sites Web économiques (niveau de l'application).

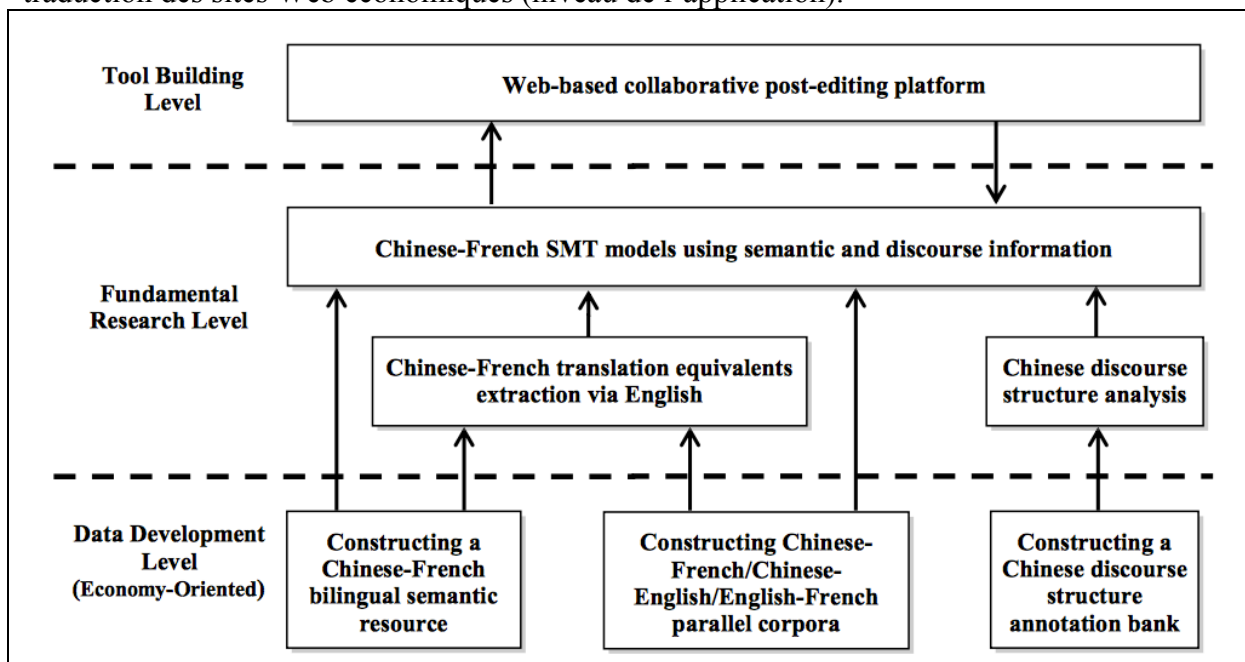


Figure 22 : Architecture à 3 niveaux et 7 « missions » du projet TABE-FC (Chen, Wang et al., 2014)

V.2 Constitution des corpus d'apprentissage

V.2.1 Recherche de sites et collecte de pages Web monolingues et bilingues

Comme déjà dit, il est difficile de construire un grand corpus parallèle chinois-français, en particulier pour le domaine économique.

Par conséquent, il est raisonnable de construire les ressources destinées à la construction de systèmes de TA chinois-français en utilisant l'anglais comme un pivot. Nous construisons un corpus anglais-chinois et le corpus anglais-français liés au domaine de l'économie. Ensuite, pour soutenir l'apprentissage des connaissances sur les transformations structurelles entre le chinois et le français, nous créons un corpus parallèle français-chinois par transitivité et révision (en PE).

Nous avons essayé de chercher une ressource existante pour construire des données d'apprentissage. Certains sites Web économiques contiennent des segments multilingues, par exemple, le site « BOURSE DE PARIS »⁶³ (français-anglais), le site « TMX »⁶⁴ de Toronto (anglais-français), « SZSE » de Shenzhen (chinois-anglais) (深圳证券交易所), et le site « abnnewswire.net »⁶⁵ (multilingue).

Nous avons téléchargé les pages Web parallèles (.html) (Figure 23), et puis fait l'alignement et l'extraction des segments parallèles sur les pages Web en utilisant LF ALIGNER.

Release time	Code	Stock Name	Document
02/10/2013 16:21	63701	JP#HSI RC1406G	Debt and Structured Products - [Expiry Announcement regarding Callable Bull/Bear Contract] Notice of Valuation of Residual Value of Callable Bull/Bear Contracts (63701) issued by J.P. Morgan Structured Products B.V. (35.20KB, PDF)
02/10/2013 16:30	63666	BP#HSI RC14080	Debt and Structured Products - [Expiry Announcement regarding Callable Bull/Bear Contract] Notice of Valuation of Residual Value of 300,000,000 European Style (Cash Settled) Category R Callable Bull Contracts 2013-2014 relating to the Hang Seng Index (the "CBBs") (Stock Code : 63666) (8.45KB, PDF)
02/10/2013 16:45	63674	UB#HSI RC1401H	Debt and Structured Products - [Expiry Announcement regarding Callable Bull/Bear Contract] Notice of Valuation of Residual Value of Callable Bull/Bear Contracts issued by UBS AG (21.44KB, PDF)
02/10/2013 16:45	63675	UB#HSI RC1407R	Debt and Structured Products - [Expiry Announcement regarding Callable Bull/Bear Contract] Notice of Valuation of Residual Value of Callable Bull/Bear Contracts issued by UBS AG (21.44KB, PDF)
02/10/2013 16:51	63667	CS#HSI RC1409N	Debt and Structured Products - [Expiry Announcement regarding Callable Bull/Bear Contract] Announcement on Valuation of Residual Value (80.78KB, PDF)
02/10/2013 16:52	63669	CS#HSI RC14090	Debt and Structured Products - [Expiry Announcement regarding Callable Bull/Bear Contract] Announcement on Valuation of Residual Value (80.06KB, PDF)
02/10/2013 17:20	63726	SC#HSI RC1406D	Debt and Structured Products - [Supplemental Listing Document of Callable Bull/Bear Contract] Supplemental Listing Document for Callable Bull/Bear Contracts (CBBs) to be issued by Standard Chartered Bank (428.36KB, PDF)
02/10/2013	63746	SC#HSI RP1403C	Debt and Structured Products - [Supplemental Listing Document of

時間排序	代號	股份名稱	文件類型
02/10/2013 16:21	63701	恒指準通四六牛G	債券及結構性產品 - [牛熊證到期公告] J.P. Morgan Structured Products B.V. 發行之可贖回牛熊證 (63701) 之剩餘價值估值通知 (163.89KB, PDF)
02/10/2013 16:30	63666	恒指法巴四八牛O	債券及結構性產品 - [牛熊證到期公告] 有關恒指指數300,000,000份2013-2014年歐式(現金結算)B類可贖回牛熊證(「牛熊證」)之估值剩餘價值通知(證券代號: 63666) (74.45KB, PDF)
02/10/2013 16:45	63674	恒指瑞銀四一牛H	債券及結構性產品 - [牛熊證到期公告] 瑞士銀行(UBS AG)於瑞士註冊成立之有限公司透過其倫敦分行發行之可贖回牛熊證(「牛熊證」)之估值剩餘價值通知 (118.16KB, PDF)
02/10/2013 16:45	63675	恒指瑞銀四七牛R	債券及結構性產品 - [牛熊證到期公告] 瑞士銀行(UBS AG)於瑞士註冊成立之有限公司透過其倫敦分行發行之可贖回牛熊證(「牛熊證」)之估值剩餘價值通知 (118.16KB, PDF)
02/10/2013 16:51	63667	恒指瑞銀四九牛N	債券及結構性產品 - [牛熊證到期公告] 剩餘價值估值通知 (142.39KB, PDF)
02/10/2013 16:52	63669	恒指瑞銀四九牛O	債券及結構性產品 - [牛熊證到期公告] 剩餘價值估值通知 (142.00KB, PDF)
02/10/2013 17:20	63726	恒指渣打四六牛D	債券及結構性產品 - [牛熊證之補充上市文件] 由渣打銀行發行之牛熊證之補充上市文件 (326.36KB, PDF)
02/10/2013 17:20	63746	恒指渣打四三熊C	債券及結構性產品 - [牛熊證之補充上市文件] 由渣打銀行發行之牛熊證之補充上市文件 (326.36KB, PDF)
02/10/2013 17:20	63748	恒指渣打四三熊D	債券及結構性產品 - [牛熊證之補充上市文件] 由渣打銀行發行之牛熊證之補充上市文件 (326.36KB, PDF)
02/10/2013 17:20	63749	恒指渣打四三熊E	債券及結構性產品 - [牛熊證之補充上市文件] 由渣打銀行發行之牛熊證之補充上市文件 (326.36KB, PDF)
02/10/2013 17:35	63725	恒指法巴四八牛P	債券及結構性產品 - [牛熊證之補充上市文件] 關於BNP PARIBAS ARBITRAGE ISSUANCE B.V.所發行之股份代號 63725 之牛熊證之補充上市文件 (380.44KB, PDF)
02/10/2013 17:37	63751	恒指法巴四二熊U	債券及結構性產品 - [牛熊證之補充上市文件] 關於BNP PARIBAS ARBITRAGE ISSUANCE B.V.所發行之股份代號 63751 之牛熊證之補充上市文件 (380.44KB, PDF)
02/10/2013 17:37	63860	恒指瑞信四一熊O	債券及結構性產品 - [牛熊證發行公告] 由 Credit Suisse AG 發行之可贖回牛熊證之公佈 (180.36KB, PDF)
02/10/2013 17:37	63863	恒指瑞信四一熊P	債券及結構性產品 - [牛熊證發行公告] 由 Credit Suisse AG 發行之可贖回牛熊證之公佈 (180.36KB, PDF)
02/10/2013 17:37	63865	恒指瑞信四一熊Q	債券及結構性產品 - [牛熊證發行公告] 由 Credit Suisse AG 發行之可贖回牛熊證之公佈 (180.36KB, PDF)

Figure 23 : Exemple de page Web économique parallèle

V.2.2 Nettoyage et filtrage

À partir des sites Web liés à l'échange d'actions, environ 1000 pages bilingues ont été explorées et traitées (le Tableau 43 présente les statistiques de cet ensemble de données).

⁶³ <http://www.boursedeparis.fr>

⁶⁴ <http://www.tmx.com>

⁶⁵ <http://abnnewswire.net>

Tableau 43 : Statistiques des pages Web collectées

Paire de langues	Nombre de pages	Taille
chinois-anglais	761	39,4Mo
anglais-français	250	12,5Mo

Pour obtenir un corpus parallèle propre, il est nécessaire de nettoyer les pages Web récupérées. Ce nettoyage est une opération informatique qui demande un paramétrage souvent poussé.

Nous avons présenté plus haut les problèmes à résoudre pour « filtrer » des corpus bilingues provenant de la faire du Web. Dans le cas des pages du site Web « Bourse de Hong Kong », une page Web en anglais contient le texte entouré de balises HTML, avec une valeur particulière de l'attribut *class*. Dans la Figure 24, le texte anglais et le texte chinois sont alignés, ils ont différentes valeurs de « *class* ». Nous choisissons d'abord les balises associées (<div class= '...'), puis nous comparons le nom de « *class* », et ajoutons un suffixe “_c” pour la partie chinoise. Nous extrayons enfin le texte anglais et le texte chinois. Nous trouvons qu'il y a un commentaire mélangé avec le texte, et nous ne pouvons pas être sûr que ce texte est bien traduit, donc nous l'éliminons.

```

677 <div class="hkex-market-top5shares-disclaimer">HKEx disseminates HSI, HSCEI, VHSI,
... S&amp;P/HKEx Large Cap, S&amp;P/HKEx GEM, CSI 300, CES 120, CES A80 and CES HKMI. Except
... CSI 300 index, CES 120, CES A80 and CES HKMI which are rounded to 2 decimal places, the
... other indices are disseminated on an "as is" basis. HKEx does not accept liability for any
... loss or damage arising from any inaccuracy or omission.
678 <BR>HSI, HSCEI and VHSI are compiled and maintained by Hang Seng Indexes Company Ltd. CSI
... 300 is provided by China Securities Index Company Ltd. CES 120, CES A80 and CES HKMI are
... provided by China Exchanges Services Company Ltd.
679 <!--
680 HKEx disseminates HSI, HSCEI, S&amp;P/HKEx Large Cap, S&amp;P/HKEx GEM and CSI 300. Except
... CSI 300 index which is rounded up to 2 decimal places, the other indices are disseminated
... on an "as is" basis. HKEx does not accept liability for any loss or damage arising from
... any inaccuracy or omission.
681 <br>HSI and HSCEI are compiled and maintained by Hang Seng Indexes Company Ltd and CSI 300
... is provided by China Securities Index Company Ltd.
682 -->
683 </div>
672 }
673 // -->
674 </script>
675 <div class="hkex-market-top5shares-disclaimer_c">
676 香港交易所將恒生指數、恒生中國企業指數、恆指波幅指數、標準普爾／香港交易所大型股指數、
677 標準普爾／香港交易所創業板指數、滬深300指數、中環120、中環A80及中環香港內地指數資料發布，
678 除滬深300指數、中環120、中環A80及中環香港內地指數四捨五入至小數點後兩位外，其他指數皆按既得之資料發布。
679 香港交易所對於任何因資料不確或遺漏而引致之損失或損害概不負責。
680 <BR>恒生指數、恒生中國企業指數及恆指波幅指數屬恒生指數有限公司擁有並由其提供，
681 滬深300指數由中證指數有限公司提供，中環120、中環A80及中環香港內地指數由中環證券交易服務有限公司提供。
682 <!--
683 香港交易所將恒生指數、恒生中國企業指數、標準普爾／香港交易所大型股指數標準、普爾／香港交易所創業板指數及滬
... 其他指數皆按既得之資料發布。香港交易所對於任何因資料不確或遺漏而引致之損失或損害概不負責。<br>
684 恒生指數系列及恒生中國企業指數系列屬恒生指數有限公司擁有並由其提供，滬深300指數由中證指數有限公司提供。
685 -->
686 </div>
687 </td>

```

Figure 24 : Exemple d'une page Web du site de "Bourse de Hong Kong" en format html

À la fin de ce processus, nous avons obtenu 3000 bisegments « propres » extraits des sites Web boursiers, en anglais-chinois. Les segments chinois étaient en caractères traditionnels ; nous les avons convertis en caractères simplifiés, grâce au système de conversion intelligente des caractères chinois simplifiés et traditionnels (Shi et al., 2013)⁶⁶.

Tableau 44 : Exemple de conversion des caractères chinois du traditionnel vers le simplifié

Traditionnel	董事會整體負責確保集團的會計及財務匯報制度健全，並設有適當的內部監控及風險管理系統。
Simplifié	董事会整体负责确保集团的会计及财务汇报制度健全，并设有适当的内部监控及风险管理系统。
Traditionnel	財務匯報 董事會承諾以平衡及清晰可明的方式向股東及其他權益人評估集團的表現、財務狀況及前景。
Simplifié	财务汇报 董事会承诺以平衡及清晰可明的方式向股东及其他权益人评估集团的表现、财务状况及前景。
Traditionnel	這承諾包括所有發布資料，包括但不限於財務報表、按監管要求發出的公告和其他公司通訊。
Simplifié	这承诺包括所有发布资料，包括但不限于财务报表、按监管要求发出的公告和其他公司通讯。

⁶⁶ Précision : 99,991% (<http://jf.cloudtranslation.cc/>)

Nous les avons mis en format TXT et TMX. Voici un exemple des segments extraits au format TMX dans le Figure 25.

```
<?xml version="1.0" encoding="utf-8" ?>
<!DOCTYPE tmx SYSTEM "tmx14.dtd">
<tmx version="1.4">
  <header
    datatype="text"
    encoding="UTF-8"
    segtype="sentence"
    adminlang="ZH"
    srclang="ZH"
    o-tmf="TW4Win 2.0 Format"
  >
  </header>
  <body>
    <tu creationdate="20140116T095746Z" creationid="eco-hkex"><prop
      type="Txt::Note">hkex-zh-en</prop>
    <tuv xml:lang="ZH"><seg>董事会整体负责确保集团的会计及财务汇报制度健全，并设有适当的内部监控及风险管理系统。</seg></tuv>
    <tuv xml:lang="EN"><seg>The Board has overall responsibility for ensuring the integrity of
      the Group's accounting and financial reporting systems, and that appropriate systems of
      internal control and risk management are in place.</seg></tuv> </tu>

    <tu creationdate="20140116T095746Z" creationid="eco-hkex"><prop
      type="Txt::Note">hkex-zh-en</prop>
    <tuv xml:lang="ZH"><seg>财务汇报 董事会承诺以平衡及清晰可明的方式向股东及其他权益人评估集团的表现、财务状况及前景。</seg></tuv>
    <tuv xml:lang="EN"><seg>Financial reporting The Board is committed to presenting a
      balanced, clear and comprehensible assessment of the Group's performance, financial
      position and prospects to our shareholders and other stakeholders.</seg></tuv> </tu>

    <tu creationdate="20140116T095746Z" creationid="eco-hkex"><prop
      type="Txt::Note">hkex-zh-en</prop>
    <tuv xml:lang="ZH"><seg>这承诺包括所有发布资料，包括但不限于财务报表、按监管要求发出的公告和其他公司通讯。</seg></tuv>
    <tuv xml:lang="EN"><seg>This commitment encompasses all published information, including
      but not limited to the financial statements, regulatory announcements and other corporate
      communications.</seg></tuv> </tu>
```

Figure 25 : Exemple de segments chinois-anglais extraits à partir de pages Web

Nous avons utilisé la même méthode pour extraire des segments parallèles anglais-français à partir du site Web « BOURSE DE PARIS ». Nous avons obtenu 2000 segments français-anglais (voir le Tableau 45).

V.2.3 TA par GT, puis PE (production d'un corpus parallèle)

Nous avons donc pu créer un corpus chinois-anglais par extraction de bisegments à partir des sites Web boursiers, et la traduction anglaise est d'assez une bonne qualité. Nous avons donc pu les utiliser pour entraîner le système de TA.

Pour la paire de langues français-chinois, nous n'avons pas de corpus parallèle, bien que nous ayons trouvé des sites Web boursiers et économiques disponibles en français et en anglais. C'est le cas du site « TMX » de Toronto, où le texte français est traduit par TA sans PE en anglais, mais pas en chinois.

Remarque : les utilisateurs semblent bien « couler » en anglais, mais il y a de nombreux contresens.

Pour construire le corpus français-chinois, nous créons une iMAG pour post-éditer les segments du français vers le chinois. Nous choisissons le site Web « BOURSE DE PARIS⁶⁷ » comme la ressource française. 2 étudiantes chinoises nous ont aidé à faire la post-édition dans le cadre de leurs stages TER de M1-info.

⁶⁷ <http://www.boursedeparis.fr>



Figure 26 : Capture d'écran de l'iMAG "Bourse de Paris" en chinois

Pour assurer la qualité des segments, il faut choisir quelques segments à post-éditer. Le site Web « Bourse de Paris » contient un centre d'apprentissage, et il présente les règles du marché, des définitions, un glossaire de bourse, etc. Nous avons post-édité les résultats de TA fournis par GT.

Dans la Figure 27, nous montrons un exemple de la traduction de GT et la post-édition faite dans SECTRA_W. GT traduit "Action" comme "Ce que l'on fait." (行动), mais ici, il faut le traduire par « stock » (股票).

Texte Français original: Google
实用信息

Actions

[Proposer une meilleure traduction](#)

行动

A股代表所有权的公司股份。如果你是股东，你拥有一定比例的股份公司的资产，这是可能的，你是根据社会的利益作为支付股息。

有股两大类：普通股和优先股。普通股给他们所有人的权利给公司：知情权，有权在股东和财务权的股东大会上，如以收取股息投票。优先股可提供经济利益，如优先股息，但没有投票权。附着到优先股权利取决于公司和在不同国家的规定的状态。



Figure 27 : Comparaison de la traduction de GT et de la post-édition humaine

Les deux post-éditeurs étaient des étudiantes chinoises en informatique à l'Université Joseph Fourier, sans connaissances particulières en économie. Cependant, on voit dans la Figure 27 que les segments post-édités sont très compréhensibles et fidèles.

V.3 Construction de systèmes de TA

V.3.1 Construction de systèmes Moses "ligne de base"

Pour construire des systèmes de TA français-chinois adaptés au domaine économique et boursier, nous avons extrait des segments français-anglais, anglais-chinois à partir de sites Web boursiers et économique. Nous avons ainsi obtenu 2000 segments français-chinois via la post-édition (voir le Tableau 45).

Tableau 45 : Statistiques sur la ressource économique et boursière

Paire de langues	Nb de segments	Nb de mots (source)	Nb de mots (cible)	Nb de p_std (source)	Nb de p_std (cible)
anglais-chinois	3000	104K	128K (caractère chinois)	416	321
anglais-français	2000	60K	54K	241	215
français-chinois (PE)	2000	44K	46K (caractère chinois)	176	115

Dans le cadre du projet TABE-FC, nous avons d'abord construit deux systèmes de TA (anglais→français et anglais→chinois). Ensuite nous avons traduit les segments anglais extraits de « Bourse de Paris » vers le chinois, et les segments anglais extraits de « Shenzhen Stock Exchange (SZSE)⁶⁸ » vers le français. Enfin nous avons construit le système de TA français→chinois en utilisant les segments traduits par les systèmes de TA et les segments post-édités du français vers le chinois.

⁶⁸ <http://www.szse.cn/main/>

V.3.1.1 Système de TA anglais→français

Nous sommes parti des 2000 segments parallèles anglais-français obtenus à partir du site Web « BOURSE DE PARIS », concernant le domaine boursier et économique. Comme ces données ne sont pas suffisantes pour entraîner un système de TA français→anglais avec MOSES, nous avons ajouté aux données d'entraînement 2M bisegments du corpus MultiUN anglais-français.

Le système de TA anglais→français a été utilisé pour traduire les segments anglais (3000 segments) extraits du site Web « SZSE » vers le français. Enfin, nous avons produit 3000 bisegments français-chinois, et la partie française a été traduite par le système de TA anglais→français.

V.3.1.2 Système de TA anglais→chinois

Pour construire le système de TA anglais→chinois, nous avons pris 3000 segments parallèles anglais→chinois extraits du site Web « SZSE », et nous avons ajouté aux données d'entraînement 2M bisegments du corpus MULTIUN anglais-chinois.

Le système de TA anglais→chinois a été utilisé pour traduire les segments anglais extraits du site Web « SZSE » vers le français. Enfin, nous avons produit 2000 bisegments français-chinois, et la partie chinoise a été traduite par notre système de TA anglais→chinois de bonne qualité.

V.3.1.3 Système de TA français→chinois

Nous avons traduit 3000 segments d'anglais en français, et 2000 segments d'anglais en chinois. Plus de 2000 segments ont été post-édités du français vers le chinois. Au total, nous avons obtenu 7000 bisegments français-chinois de bonne qualité.

V.3.1.4 Système de TA chinois→français

Pour l'instant, nous n'avons pas de corpus parallèle de bonne qualité pour entraîner un système de TA du chinois vers le français adapté au domaine économique et boursier, et nous n'avons pas non plus de MT post-éditée du chinois vers le français. Pour tester la traduction français→chinois, nous avons essayé d'utiliser la MT français→chinois « à l'envers ». Mais les résultats, comme on pourrait s'y attendre, sont assez désastreux.

V.3.2 Avancement de l'expérimentation

Pour obtenir plus de segments post-édités, nous avons construit deux nouvelles iMAG pour les sites Web « *SIX Swiss Exchange (SIX)*⁶⁹ » (français→chinois), et « *Hong Kong Stock Exchange (HKEX)*⁷⁰ » (anglais→chinois).

Cette expérience est toujours en cours au moment de la rédaction. Nous utilisons nos systèmes de TA et post-éditons leurs résultats sous SECTRA_W. Les systèmes de TA sont installés sur un serveur du laboratoire LIG, et on peut y accéder via TRADOH.

V.3.3 Résultats provisoires

Nous avons commencé à utiliser notre système pour post-éditer les pages Web boursières et économiques des sites « SIX », « HKEX », et « BOURSE DE PARIS ». Pour l'instant, nous n'avons pas encore assez de données pour présenter des statistiques et surtout une évaluation

⁶⁹ <http://www.six-swiss-exchange.com/index.html>

⁷⁰ <http://www.hkex.com.hk/eng/index.htm>

significatives. Nous espérons que nous pourrons fournir cela au moment de la soutenance, et que nous pourrons répondre à trois questions importantes dans ce contexte :

1. La qualité d'usage est-elle suffisante ou non sans intervention humaine ?
2. Le passage par l'anglais est-il une bonne chose (avec ou sans PE de l'anglais) ?
3. Peut-on réduire significativement le temps de recompilation d'un système MOSES avec de nouvelles données (bisegments) obtenue par PE.

Chapitre VI Démonstration de l'intérêt de l'apprentissage incrémental en TA statistique

Introduction

Dans l'objectif permanent d'améliorer la qualité de notre système de TA fr-zh, nous avons accumulé des bisegments de bonne qualité par post-édition sur la plate-forme SECTRA_W/MAG. Étant donné la controverse qu'il y avait vers 2012 sur l'utilité de l'apprentissage incrémental pour des systèmes de type MOSES, nous avons cherché à monter une expérience pour déterminer les conditions favorables à l'utilisation de cette technique. Cette expérience a été menée en 3 phases. (1) J'ai d'abord construit un système appelé `Moses_fr-zh_v0` à partir de 10K bisegments post-édités, et j'ai travaillé seulement sur la MT du LIG-LAB, en post-éditant des segments choisis parmi les 6000 non post-édités, et en faisant des mesures sur des ensembles croissants de segments post-édités. Après 10 itérations d'apprentissage incrémental, les résultats étaient encourageants mais pas encore concluants. (2) Dans le cadre d'un stage d'été, deux étudiants chinois ont travaillé sur la post-édition en chinois d'articles en français et sur des supports de divers cours de Master 1 environ 13000 segments. (3) Enfin, disposant d'environ 30K bisegments de qualité (en fr-zh), nous avons recommencé l'expérience en améliorant son organisation et en l'automatisant.

Au total, nous avons prouvé que l'apprentissage incrémental peut améliorer la qualité de TA fr-zh, en tout cas si l'on vise un « sous-langage », et si de temps en temps, on réentraîne le système en utilisant toutes les données d'entraînement.

VI.1 Contexte

VI.1.1 Motivations

Pour obtenir une bonne qualité avec un modèle de type Moses sur un "sous-langage" précis, on n'a pas besoin de très grandes quantités de données d'entraînement (comme les 200M de bisegments du système généraliste de l'UE), mais il en faut quand même un certain nombre, en fonction de la taille et de la complexité du sous-langage en question. Une expérience précédente a montré que, à l'extrême, on pouvait se contenter de 1000 bisegments dans le cas d'un tout petit sous-langage (SMS de petites annonces en occasion automobile (Daoud, 2007), (Hajlaoui and Boitet, 2008)). Dans le cas de notes techniques et de la paire français-anglais, la société L&M est arrivée à des taux de BLEU de l'ordre de 70% avec 300K exemples.

Si l'on s'attaque à un nouveau sous-langage, on n'a pas encore de corpus de bisegments de taille suffisante. On commence donc par utiliser un système généraliste existant, et à post-éditer ses résultats, jusqu'à obtenir une taille suffisante. On entraîne alors son propre système sur un corpus d'entraînement petit mais de bonne qualité, en le mélangeant, si les résultats sont meilleurs, avec une certaine quantité de corpus général, et on met ce système en service. Au début, la qualité d'usage, mesurée par exemple par le temps de post-édition, n'est en général pas supérieure à celle d'un système généraliste (comme GT). Mais, si l'on produit de nouvelles versions de « son » système de TA en intégrant aux données d'entraînement les post-éditions qu'il produit sur de nouveaux segments, on y arrive.

Malheureusement, la création d'une nouvelle version consiste à "tout recompiler" et est extrêmement longue (plusieurs dizaines d'heures en temps partagé). Par exemple, on a fait une expérience avec 1,2M bisegments fr-zh sur un serveur à quatre cœurs (Intel i7-3770 CPU 3.40GHz), qui a pris 17~20h pour l'entraînement.

C'est un obstacle majeur à l'accroissement régulier de la qualité. Or, il est très important que les utilisateurs aient l'impression que « le système apprend », et, par exemple, que les contributions qu'ils ont faites dans leur activité de post-édition sont prises en compte

Fin 2011, une possibilité d'apprentissage incrémental venait d'être introduite dans MOSES, et L&M m'a demandé de l'étudier, ce que j'ai fait. C'était dans le cadre de la préparation d'une réponse à un appel d'offres, qui n'a pas été obtenu. L&M a alors provisoirement arrêté le travail sur ce point. Pour ma part, je continuais d'être intéressé par cette possibilité, et cela d'autant plus que mes expériences sur les corpus de L&M semblaient prouver que cette technique était prometteuse, même si les gains de qualité étaient nettement plus faibles que si l'on "recompilait tout", et aussi parce qu'il y avait dans le domaine diverses polémiques, l'une sur l'utilité potentielle de l'apprentissage incrémental dans MOSES, et l'autre sur l'intérêt pour les traducteurs de post-éditer des résultats de TA. J'avais en particulier lu un article mettant en doute l'utilité de cette technique (Mirkin, 2014).

Je me suis donc demandé comment je pourrais démontrer l'utilité de cette technique, et ce thème est devenu assez important dans ma recherche. J'ai expérimenté cette technique pour la première fois dans le cadre de la création d'un corpus bilingue français-chinois dans le domaine de l'énergie (voir IV.3.3).

J'avais aussi constaté que les améliorations obtenues en ajoutant des « mini-batch » et en exécutant l'AI tendaient à diminuer avec le temps, et étaient toujours bien inférieures aux améliorations obtenues par un réentraînement complet.

Comme L&M a abandonné ce projet, n'ayant pas obtenu le contrat espéré avec EDF, j'ai poursuivi cette recherche dans le cadre du laboratoire, en prenant comme base expérimentale le site du LIG, pour lequel nous avons déjà une MT 8000 segments dont 2000 post-édités.

VI.1.2 Expérience sur le site du LIG

VI.1.2.1 Motivations et conditions expérimentales

Pourquoi le site du LIG ? Le site Web LIG-Lab⁷¹ est le site Web officiel du laboratoire LIG (Laboratoire d'Informatique de Grenoble). Le LIG rassemble près de 500 chercheurs, enseignants-chercheurs, doctorants et personnels en support à la recherche. C'est aussi le tout premier site Web pour lequel une iMAG a été créée (voir I.1.2.1.2). Un certain nombre de chercheurs, doctorants, ou stagiaires sont chinois, et ils souhaitent accéder notre site Web traduit en chinois.

Sous-langage envisagé. Le site Web du LIG contient la présentation du laboratoire, la présentation de ses 22 équipes de recherche, ses contributions du développement des aspects fondamentaux de la discipline (modèles, langages, méthodes, algorithmes) et les événements du laboratoire. Tous les contenus tournent autour de l'informatique. Une instance d'iMAG existe depuis fin 2008 et est dédiée à ce site, ou plutôt à son *sous-langage*.

Pré-condition de l'expérience. La Figure 28 montre la page d'accueil de cette iMAG dédiée (LIG-LAB) accédée en chinois. La langue source est le français, et il y a 34 langues d'accès possibles. Notre expérience porte sur la paire de langues français-chinois. L'iMAG du LIG contenait environ 8K segments français (en février 2013), et environ 2K segments avaient été post-édités.

⁷¹ Site Web du laboratoire LIG : <https://www.liglab.fr>

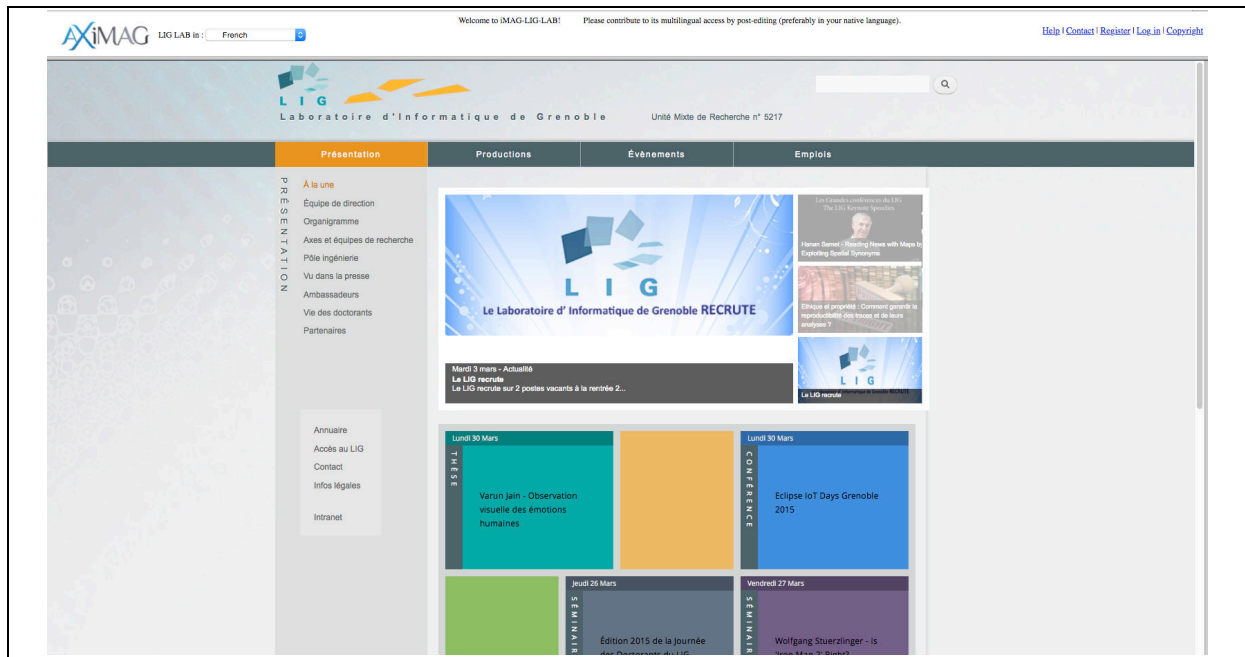


Figure 28 : Site du LIG vu en chinois à travers une iMAG

VI.2 Expérimentation

VI.2.1 Phase 1 (2-6/2013)

VI.2.1.1 Expérience sur le site du LIG

Notre corpus d'apprentissage initial a été construit en choisissant dans SECTRA deux parties tirées de deux MT différentes : une partie provenait de la MT de LIG-LAB⁷², qui avait au total 8000 segments en français, dont 2000 étaient déjà post-édités de français en chinois. L'autre partie provenait de la MT DEMO⁷³, qui avait été créée en 2012, et qui contenait 8000 segments post-édités de français en chinois. Ces 10K bisegments post-édités ont été utilisés pour créer la version initiale de notre système Moses français-chinois dédié au sous-langage du LIG. Nous l'avons appelée MOSES-LIG-FR-ZH-V₀ (Tableau 46).

Tableau 46 : Statistiques sur les données d'entraînement de la phase 1

	Nb de bisegments	Nb de mots (source)	Pages standard	Caractères chinois	Pages standard
Moses V ₀ (initial)	10724	160K	644	182K	455

Dans l'expérience, nous avons travaillé seulement sur la MT du LIG-LAB, en post-éditant des segments choisis parmi les 6000 non post-édités, et en faisant les mesures sur les ensembles croissants de segments post-édités.

Nous utilisons les termes et les notations suivants.

- **Page standard** : texte contenant 250 mots (1400 caractères) en français, ou 400 caractères en chinois.
- **Segment** : unité de post-édition (phrase ou un titre).
- **Page logique** : page Web générée par SECTRA_W pour l'interface de PE et contenant N segments. La valeur par défaut est N=20.

⁷² LIG-LAB : un projet d'iMAG. LIG-LAB est l'iMAG dédiée au site du LIG.

⁷³ Demo2 : c'est une MT dédiée au site du premier auteur.

- **MT_LIG** : MT de LIG-LAB.
- **MT_INC** : nouvelle partie de MT à extraire pour l'apprentissage incrémental.
- **NewPE** : ensemble des nouveaux segments post-édités.
- **Moses-LIG-fr-zh-V_N** : N-ième version de Moses-LIG-fr-zh (après N mises à jour de Moses-LIG).
- **Moses_INC** : opération d'apprentissage incrémental.
- **SEG[i]** : segment i.
- **SEGPE[i]** : post-édition du segment i.
- **ApplyPE (Moses-LIG_N, SEG[i])** : post-édition du segment i à partir de sa TA par Moses-LIG_N.

Dans cette expérience, j'ai post-édité les segments non post-édités, par pages logiques. SECTra note le temps de post-édition primaire (Tpe_1)⁷⁴ pour chaque segment. Lorsque 10 pages logiques (200 segments) ont été post-éditées, on les utilise pour mettre à jour Moses-LIG-fr-zh (de V_N à V_{N+1}). Nous répétons cela 10 fois de suite (2000 segments), et mesurons le temps moyen de post-édition. Enfin, nous utilisons les segments source du site LIG, leurs post-éditions à partir de la MT de LIG-LAB, et les traductions par Moses-LIG-fr-zh de V₁ à V₁₀ pour évaluer la performance de chaque Moses-LIG-fr-zh-V_N par BLEU (Papineni, Roukos et al., 2002), NIST (Doddington, 2002) et par TER (Snover et al., 2009).

Il y a finalement 3 opérations essentielles

Opération 1 : post-édition des segments non post-édités, traduits par Moses-LIG-fr-zh-V_N.

`SEGPE[i] := ApplyPE (Moses-LIG-fr-zh-VN, SEG[i]) si SEGPE[i] ∉ MT_LIG_PE`

Opération 2 : apprentissage incrémental de Moses-LIG_{N+1}.

`MT_INCR := Extraire (NewPE, Niveau, Score) NewPE ⊆ MT_LIG`

`Moses-LIG-fr-zh-VN+1 := Moses_INC (Moses-LIG-fr-zh-VN, MT_INCR);`

Opération 3 : mesures (Voir VI.2.1.3).

VI.2.1.2 Processus d'AI

À l'itération N du processus, tous les segments de la MT sont traduits par la version N du TA (Moses-V_N), y compris, de façon continue, les nouveaux segments créés par le site. Certains segments non encore post-édités sont post-édités à cette itération (de façon opportuniste ou organisée).

Quand on a un certain nombre de nouveaux segments post-édités, jugé "bons" (200 dans notre expérience), on les traite (séparation des mots, nettoyage, traitement de la casse (*truecasing*), alignement, etc.), et on les ajoute à la table de traductions.

Moses-V_N est mis à jour vers Moses-V_{N+1}. À la fin de l'itération N, on met en service Moses-V_{N+1}, et on lui fait traduire tous les segments, post-édités ou non. Cela permet de mesurer la différence de qualité sur la partie déjà post-éditée. On passe alors à l'itération N+1. On obtient ainsi une suite de versions du STA (Moses-V₁, ..., Moses-V_N, Moses-V_{N+1}, ...), qui ne s'arrête (en usage) que quand tous les segments sont post-édités et qu'il n'en reste pas assez de "bons" pour procéder à l'itération suivante.

⁷⁴ Temps passé dans la zone de PE.

VI.2.1.3 Résultats

VI.2.1.3.1 Évolution du temps de post-édition

Ces temps sont ceux mesurés par SECTRA_W. Le temps moyen de post-édition (PE), pour chaque nouvel ensemble de 200 segments, diminue à chaque itération. Après la dixième, le temps moyen de PE par page standard (de 250 mots) a été réduit de 3,8 minutes sur 30,7 au départ, soit 12,4%. Si on le compare avec le temps de PE en partant de résultats de GT, on voit qu'il s'en rapproche. La forme de la courbe faisait penser (à la fin de cette phase 1 de l'expérience) qu'il passerait dessous après environ 30 itérations.

Tableau 47 : Évaluation du temps de post-édition (2-6/2013)

Fois	Segments	Mots/seg	page standard	temps de PE (min)	min/page standard
1	200	4,93	3,9	121,0	30,7
2	200	5,15	4,1	125,7	30,5
3	200	4,95	4,0	119,2	30,1
4	200	5,03	4,0	119,4	29,7
5	200	5,29	4,2	123,5	29,2
6	200	5,15	4,1	118,2	28,7
7	200	5,13	4,1	115,2	28,1
8	200	5,25	4,2	116,8	27,8
9	200	5,15	4,1	112,8	27,4
10	200	5,34	4,3	115,0	26,9
Total	2000		41,1	1186,8	

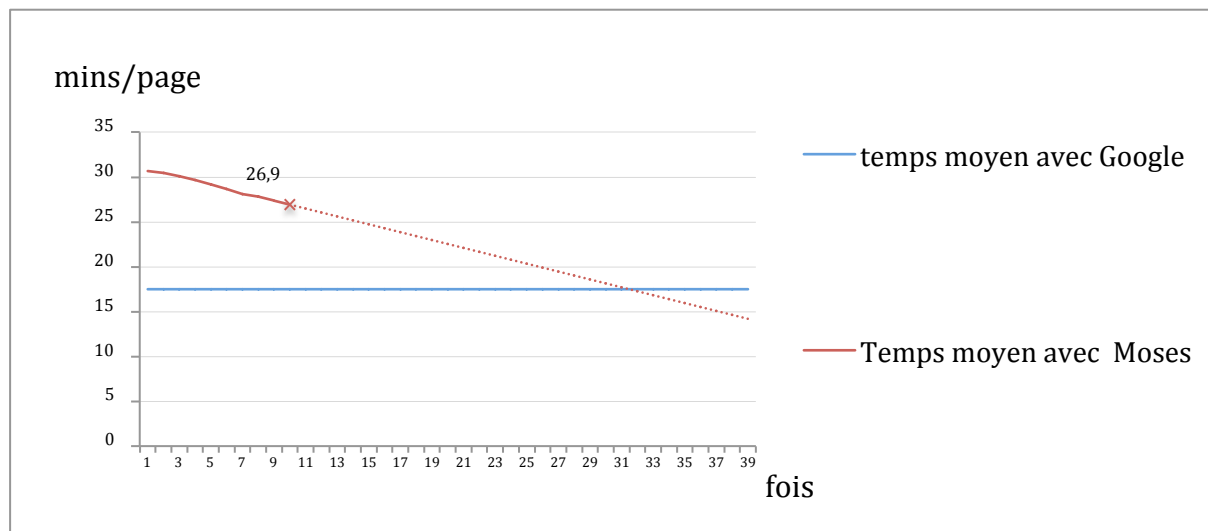


Figure 29 : Diminution de temps moyen de PE (par page standard) avec AI dans la phase 1 de l'expérience

VI.2.1.3.2 Évaluations basées sur des références

On a M segments post-édités, et on a la nouvelle version N de Moses-LIG. On traduit tous les segments avec elle, et on mesure BLEU, NIST et TER sur :

- **Origine (OR)** : les M segments post-édités.
- **Nouveau (NV)** : les 200 segments nouvellement post-édités.
- **Tout (TT)** : tout ce qui a été post-édité.

Tableau 48 : Évaluations basées sur des références (BLEU, NIST, TER)

Version de Moses-LIG	BLEU			NIST			TER		
	OR	NV	TT	OR	NV	TT	OR	NV	TT
0	22,10	19,78	21,58	5,8782	5,5122	5,7655	64,32	65,51	65,32
1	22,12	23,31	21,79	5,8823	5,7356	5,8018	64,28	64,89	64,73
2	22,23	24,64	21,87	5,9125	5,8267	5,8230	64,14	64,62	64,21
3	22,37	20,82	21,53	5,9532	5,4112	5,6547	64,02	66,74	64,88
4	22,46	20,08	21,61	5,9708	5,6349	5,7743	63,98	65,63	64,43
5	22,59	21,33	21,65	5,9821	6,1183	5,8312	63,90	64,37	64,06
6	22,64	20,23	21,67	6,0246	5,4721	5,8414	63,85	65,32	64,52
7	22,71	18,92	21,43	6,0412	5,0537	5,5752	63,65	67,58	65,03
8	22,79	22,41	22,57	6,0856	5,8269	5,8763	63,59	64,73	64,31
9	22,84	20,23	21,41	6,0954	5,6412	5,8838	63,46	66,08	65,27
10	22,93			6,1139			63,34		

À la 10^{ème} itération, il n'y avait pas de segments post-édités (assez bons). Au début de la N^{ème} itération, $M=2000+200*N$, $TT=2000+200*(N+1)$.

Conclusion

Nous avons fait une première expérience sur l'évaluation et l'amélioration incrémentale de la TA, réalisée avec le système MOSES, basé sur SECTRA_W et MOSES-LIG. Les résultats des mesures montrent que la méthode d'apprentissage incrémental permet de réduire le temps de post-édition de 12,4% en 10 itérations, et d'améliorer les mesures basées sur les références (BLEU : \uparrow 0,83, NIST : \uparrow 0,2357, TER : \downarrow 0,98).

VI.2.2 Phase 2 (7-9/2013)

VI.2.2.1 PE par deux étudiants chinois dans le cadre d'un stage d'été (Wu Jiang et Chen Shi)

Wu Jiang et Chen Shi étaient deux étudiants chinois de l'UJF qui, en juin 2013, étaient en fin de leur 1ère année de master en informatique. Ils ont fait un stage d'été en binôme au laboratoire, durant lequel ils ont fait de la post-édition sur des textes et sur des supports de divers cours de Master 1, du français vers le chinois.

Ils ont travaillé sur 2 IMAG dédiées. (1) « *Corpus par jour* »⁷⁵ est un site Web que j'avais créé pour récupérer des ressources français-chinois. J'ai sélectionné des textes en français à partir de sites de journaux LE MONDE et 20 MINUTES. Je les ai mis en format HTML, et les ai ajoutés dans mon site Web. La Figure 30 montre une capture d'écran de post-édition sur l'IMAG associée.

⁷⁵ <http://service.aximag.fr/xwiki/bin/view/imag/Corpus-lx>

{_中瑞达成自贸协定_}

{ 正值中国这个世界第二大经济体与欧盟在贸易领域精力之时，中国与瑞士之间签署了一项自贸协定。这是继今年四月的冰岛之后，第二个同中国达成自贸协定的欧洲国家。 }

{ 中国商务部部长高虎城和瑞士联邦委员会经济部长约翰·施奈德-阿曼在北京的中国商务部举行了签约仪式。 }

{ 仪式后，阿曼表示：“此项自贸协定对两国关系具有重要意义。” }

{ 中国是继欧盟与美国之后瑞士的第三次贸易伙伴。 }

{ 双边贸易额达265亿美元，其中瑞士向中国出口额为228亿美元。 }

{ 作为世界银行评估的全球第19大经济体，瑞士是少数几个能与中国在贸易领域维持积极平衡姿态的西方国家。 }

{ 中瑞间的自贸协定还需通过瑞士联邦议会通过后方能付诸生效。 }

{_皮诺向中国归还两件青铜兽首_}

{ 法国收藏家及商界掌门弗朗索瓦·皮诺先生今日向中国政府归还两件青铜兽首，这两件兽首是1860年10月第二次鸦片战争期间，英法联军劫掠圆明园时被掠夺的。 }

{ 二者一为鼠首一为兔首。 }

{ 弗朗索瓦·皮诺的儿子、开云集团董事长弗朗索瓦·亨利·皮诺在北京天安门广场的国家博物馆所举行的捐赠仪式上说：“这两件兽首带回中国的时候，我的家族也恪守了承诺，保护一国之尊严。” }

{ 弗朗索瓦·皮诺与法国前总理马林·勒克莱尔一道在媒体面前揭开并归还两件青铜兽首的面纱。 }

{_我们的学生将成为具有影响力的领袖中的一员_}

{ 此举实为百年奇遇。 }

{ 福布斯杂志上排名第63位的富豪、黑石（私募股权）投资公司的创始人苏世民先生捐赠一亿美元设立奖学金，以资助赴中国高校（留学）的项目。 }

{ 通过增加在华留学生的人数，此项目旨在增进中国与外部世界的相互理解。 }

{ 1902年，美国人塞西尔·罗德爵士出于类似的目的，用自己的遗产设立了一个以其名字命名的奖学金，这项奖学金资助着以美国学子为主的年轻学生赴英国牛津大学继续深造。 }

{ 只是当时，罗德先生考虑的是要紧密两国之间的关系。 }

{ 时至今日，世界的目光已转向了以中国为首的亚洲地区。 }

{ 这一套精英精英教育旨在北京的清华大学付诸实践，该校又被称为“中国的麻省理工”，现任及前任中国国家主席习近平与胡锦涛均在此校度过了青春岁月。 }

{ 围绕此项目还成立了一个名流汇聚的顾问委员会，其中包括三任美国前国务卿——亨利·基辛格、科林·鲍威尔和康多莉扎·赖斯；此外还包括有（法国前总统）尼古拉·萨科齐以及英国前首相（布莱尔）、加拿大前总理（马尔罗尼）和澳大利亚前总理（陆克文）。 }

{ 苏世民奖学金预计总额为3亿美元。 }

{ 其中教育、美国石油、卡特彼勒、通用电气以及摩根大通将注资一亿美元。 }

{ 剩下的一亿美元资金则由待发放。 }

{ 该项目每年将选出4名幸运的法国学生（赴清华留学深造），但目前还没有找到有意愿资助的会主。 }

{_继法国葡萄酒后，中国又将矛头指向了德国乘用车_}

{ 中欧之间的贸易大战将会愈演愈烈吗？（继葡萄酒之后，以德国车为主的欧洲高档乘用车进入了中国当局的视野。） }

{ 据《回声报》报道，中国拟对进口的欧洲高档乘用车征收新关税。 }

{ 中国本土汽车制造商可能已向中国商务部投诉其欧洲同行（的不当行为）。 }

{ 此次征税主要针对排量大于等于2.0的轿车。 }

{ 而这款车更多的是以德国车作为行业领头羊所制造的一系列豪车。 }

Accord de libre-échange Chine/Suisse

La Chine et la Suisse ont signé un accord de libre-échange, le deuxième de Pékin avec un pays d'Europe - après l'Islande, en avril dernier -, en plein bras de fer commercial entre la deuxième économie mondiale et l'Union européenne.

Le ministre chinois du Commerce, Gao Hucheng, et le ministre suisse de l'Économie, Johann Schneider-Ammann, ont apposé leur signature sur l'accord, au cours d'une cérémonie qui s'est déroulée au ministère du Commerce à Pékin.

"Cet accord de libre-échange a une signification importante pour les relations entre nos deux pays", a déclaré le ministre suisse de l'Économie, après la signature.

La Chine est le troisième partenaire commercial de la Suisse, après l'Europe et les États-Unis.

Les échanges s'élevaient à 26,3 milliards de dollars, avec 22,8 milliards de dollars pour les exportations suisses en Chine.

La Suisse, 19e économie mondiale, selon la Banque mondiale, est ainsi l'un des rares pays occidentaux à avoir une balance commerciale positive avec le géant chinois.

L'accord entre la Suisse et la Chine doit encore être approuvé par le Parlement suisse pour pouvoir entrer en vigueur.

Pinault rend deux sculptures à la Chine

La famille de l'homme d'affaires et collectionneur français François Pinault a remis aujourd'hui à la Chine deux sculptures de bronze volées lors du sac du palais d'Été par les troupes franco-britanniques en octobre 1860, lors de la Seconde Guerre de l'opium.

Ces sculptures représentent l'une une tête de rat, l'autre une tête de lapin.

"En rapportant ces deux merveilles en Chine, ma famille respecte sa promesse de préserver l'héritage national et la création artistique", a déclaré François-Henri Pinault, fils de François Pinault et PDG du groupe Kering, lors d'une cérémonie au Musée national de Chine, place Tiananmen à Pékin.

François Pinault et le vice-Premier ministre chinois, Liu Yangong, ont dévoilé les deux statues devant la presse.

Nos étudiants auront vocation à être des leaders d'influence

Une parallèle initiative n'arrive qu'une fois par siècle.

Le fondateur du fonds d'investissement Blackstone, Stephen Schwarzman- 63e homme le plus riche des États-Unis selon le magazine Forbes- vient de donner 100 millions de dollars de sa poche afin de créer une bourse d'études associée à un nouveau programme universitaire en Chine.

Objectif, favoriser la compréhension entre la Chine et le reste du monde en immergeant des étudiants étrangers en Chine.

Le même dessein poursuivi par le Britannique Cecil Rhodes lorsqu'il légua sa fortune, en 1902, à une fondation éponyme finançant les études à Oxford de jeunes diplômés, américains surtout.

Sauf qu'à l'époque, le philanthrope voulait renforcer les liens entre l'Amérique et l'Europe.

Désormais, c'est vers l'Asie que les regards se tournent, et vers la Chine en particulier.

Très élitiste, le nouveau cursus sera abrégé par l'université Tsinghua, à Pékin - le «MIT chinois» -, dont le président chinois Xi Jinping et son prédécesseur Hu Jintao ont usé les bords.

Tout aussi huppé le conseil consultatif mis en place autour de ce programme, voit trois ex-secrétaires d'État américains, Henry Kissinger, Colin Powell, Condoleezza Rice, côtoyer Nicolas Sarkozy ou encore des anciens premiers ministres britannique, canadien, australien.

Au total, la Schwarzman Scholarship représente un budget de 300 millions de dollars.

Boeing, BP, Caterpillar, GE ou encore JPMorgan Chase apportent 100 millions de dollars complémentaires.

Reste encore 100 millions de dollars à trouver.

Comme Stephen Schwarzman l'explique au Figaro, manque notamment à l'appel un donateur désireux de financer 4 heureux élus français par an.

Figure 30 : Capture d'écran de l'iMAG « Corpus par jour »

À la fin de cette période (04/07/2013-13/09/2013), nous avons post-édité 21 articles français. La statistique est présentée dans le Tableau 49.

Tableau 49 : Statistiques de post-édition sur 21 articles français 4/7-13/9/2013

Nb de doc	Segments	Mots	Mots/seg	Page_std	Temps en sec.	Temps/page_std en sec.
21	4775	115564	24,20	462,26	159408	344,85

(2) « MACAU »⁷⁶, Cette iMAG a été créée dans le cadre du projet MACAU-OFI. Elle permet aux étudiants étrangers de l'UJF d'accéder dans leur langue à des supports pédagogiques, par l'amélioration contributive des "pré-traductions" obtenue par TA (Figure 31). (Voir la présentation de MACAU au III.2.3).

⁷⁶ <http://46.105.41.94/macau/chamilo-macau/index.php>



Figure 31 : Capture d'écran de Chamilo affichant le lien AXiMAG

En septembre 2013, nous avons post-édité les supports de cours de l'UJF comme « Complexité en M1 », « IHM en M1 », « *Logique en L2* », etc. Le Tableau 50 donne la statistique de post-édition sur les supports de cours.

Tableau 50 : Statistiques de post-édition sur les supports de cours

Nb de doc	Segments	Mots	Mots/seg	page_std	Temps en sec.	Temps/page_std (calculé)
236	16069	136573	8,50	546,29	412784	755,61s =12,6 mn

VI.2.2.2 Processus d'AI

Dans cette phase, nous utilisons la méthode présentée au IV.3.4 pour construire notre système. Tout d'abord, nous utilisons le système de TA français→chinois, qui est entraîné dans la phase 1, comme système de base (« *baseline* »). Les nouveaux segments post-édités sont utilisés pour entraîner la nouvelle partie de TA. Les paramètres sont les mêmes que ceux de la phase I.

Après la construction du système de TA français→chinois avec tous les segments post-édités (Phase I + Phase II)⁷⁷, nous avons fait une expérimentation pour comparer le temps de post-édition parmi les différents systèmes de TA dans le cadre d'un stage M1 TER.

VI.2.2.3 Résultats

Pour évaluer notre système de TA, nous avons choisi 50 segments à partir du site Web du LIG en langue française, et un étudiant chinois a ensuite post-édité les prétraductions produites par GT et MOSESLIG français-chinois (100 observations) dans SECTRA_W. En fait, dans cette situation, nous ne pouvons pas analyser directement le Tpe_{total} , parce que, pour chacune des prétraductions, le segment source est le même, et le temps de recherche lexicale (l'essentiel de

⁷⁷ Dans la phase II, on n'a pas de façon pour accéder le réseau intranet du labo, donc nous faisons un fois d'apprentissage incrémental.

Tpe_2) est souvent passé sur la première prétraduction. Néanmoins, nous pouvons analyser le Tpe_1 qui est fortement corrélé avec le Tpe_{total} .

VI.2.2.3.1 Évolution du temps de post-édition

Dans le Tableau 51, le nombre moyen de mots par segment des prétraductions de MOSESLIG est presque le même qu'avec GT. Nous pouvons constater que, pour post-éditer un segment prétraduit, le Tpe_1 moyen de MOSESLIG est inférieur à celui de GT.

Tableau 51 : Résultat de l'expérimentation (en français-chinois)

TA	Nb moyen de mots par segment (source)	Nb moyen de mots par segment (cible)	Tpe_1 moyen par segment	Tpe_1 moyen par page_std ⁷⁸
MosesLIG	26	25,4	23,4 s	6,1 mn
Google	26	27,1	25,3 s	6,3 mn

VI.2.2.3.2 Conclusions provisoires

Dans la phase 2, nous n'avons pas fait l'AI comme dans la phase 1, mais nous avons réentraîné notre système de TA français-chinois avec une plus grosse MT. Notre expérimentation a produit un bon résultat. Du point de vue du temps de post-édition, le résultat est claire : notre système est maintenant un peu meilleur que GT.

VI.2.3 Phase 3 (9-12/14 et 7-11/15)

VI.2.3.1 Motivation

Dans la phase 1, nous avons fait le premier essaie de l'AI, et nous avons obtenu un résultat encourageant. Mais notre système de TA était limité par la quantité de segments post-édités, et la longueur moyenne du segment (7 ou 9 mots).

Dans la phase 2, nous n'avons pas pu faire d'AI. Nous avons réentraîné notre système avec une grande MT français-chinois, mais cela ne prouve pas que notre méthode (AI) peut améliorer le système de TA, mais ne permet pas non plus de trouver le point de « réentraînement du système ». Nous voulions donc refaire notre expérimentation pour répondre nos « doutes ».

VI.2.3.2 Poursuite de l'expérience en français-chinois sur le site du LIG

Dans cette phase, nous avons exporté tous les segments post-édités de SECTRA_W. Dans le Tableau 52, nous montrons la quantité de segments post-édités.

Tableau 52 : Nombre de segments dans chaque MT

Nom du corpus / de la MT	Quantité de segments PE
Demo2	28K
MT_INC	7K
MT_Macau	16K
MT_TED	14K
MT_énergie	9K

Préparation des données. Nous avons mélangé toutes les données, environ 80K bisegments, puis nous les avons nettoyées et les avons séparées en 16 parties. Chaque partie contient 5000 segments. Les statistiques sur les données nettoyées sont données dans le Tableau 53.

⁷⁸ Une page standard contient 250 mots français ou 400 caractères chinois

Conception de l'expérience. Nous avons fait 16 itérations. Le processus de l'AI a été présenté dans le VI.2.2.2. Le système MOSES-LIG V₀ a été construit avec 100K bisegments extraits à partir du corpus MULTIUN.

Tableau 53 : Statistiques sur les données pour l'AI (phase 3 de l'expérience)

Incrémental	Longueur moyenne de segment (source)	Longueur moyenne de segment (cible)	Nombre de mots (source)	Nombre de caractères chinois (cible)	Pages standard (source)	Pages standard (cible)
Base (MultiUN)	52,7	41,8	263 K	209 K	1 054	522,50
1	42,9	27,3	214 K	136 K	858	341,25
2	43,3	26,7	216 K	133 K	866	333,75
3	39,6	25,3	198 K	126 K	792	316,25
4	34,5	23,3	172 K	116 K	690	291,25
5	31,7	22,3	158 K	111 K	634	278,75
6	30,1	21,9	150 K	109 K	602	273,75
7	28,8	21,5	144 K	107 K	576	268,75
8	27,7	21,1	138 K	105 K	554	263,75
9	27,1	20,9	135 K	104 K	542	261,25
10	26,5	20,7	132 K	103 K	530	258,75
11	25,9	20,5	129 K	102 K	518	256,25
12	25,5	20,4	127 K	102 K	510	255,00
13	25,1	20,2	125 K	101 K	502	252,50
14	24,8	20,1	124 K	100 K	496	251,25
15	23,6	19,2	118 K	96 K	472	240,00
16	22,3	19,1	111 K	95 K	446	238,75
Au total	Longueur moyenne de segment	Longueur moyenne de segment	Nombre de mots (source)	Nombre de caractère chinois (cible)	Pages standard (source)	Pages standard (cible)
	31,3	23,1	2 660 K	1 961 K	10 642,00	4 903,75

VI.2.3.3 Résultats

VI.2.3.3.1 Évolution du temps de post-édition

Pour évaluer le temps de post-édition, nous continuons à post-éditer les nouveaux segments sur l'iMAG LIG-LAB comme nous l'avons présenté au VI.2.1.3.1. Les segments sont sélectionnés dans la MT LIG-LAB. Nous avons filtré sur la longueur, et supprimé les segments très courts.

Tableau 54 : Évaluation du temps de post-édition (9-12/2014)

Fois	Segments	Mots/seg	Page standard	Tpe (min)	Min/page standard
0	200	16,7	13,36	60,4	45,21
1	200	17,5	14	39,9	27,52
2	200	18,6	14,88	40,4	26,83
3	200	21,2	16,96	45,1	26,62
4	200	24,6	19,68	49,0	24,89
5	200	16,8	13,44	32,3	24,03
6	200	18,2	14,56	34,4	23,62
7	200	16,5	13,2	29,2	22,12
8	200	16,1	12,88	27,8	21,56
9	200	14,3	11,44	23,7	20,72
10	200	20,5	16,4	32,9	20,06
11	200	27,1	21,68	41,1	18,98
12	200	13,8	11,04	18,8	17,04
13	200	15,7	12,56	21,0	16,73
14	200	19,1	15,28	24,5	16,06
15	200	13,7	10,96	16,7	15,24
16	200	18,6	14,88	22,5	15,11

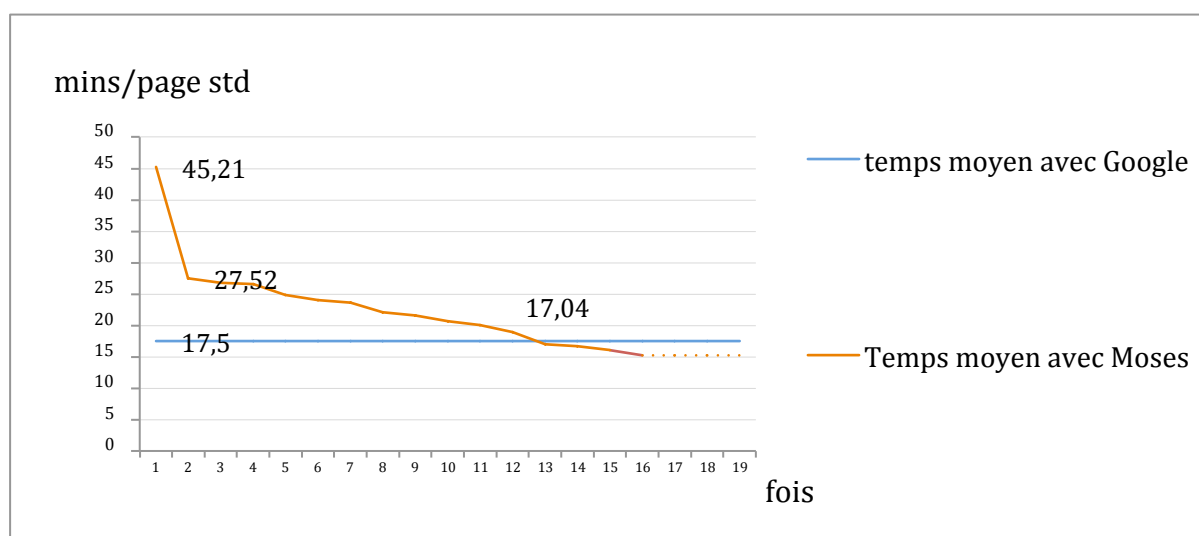


Figure 32 : Diminution du temps moyen de PE (par page) avec AI dans la phase 3 de l'expérience

VI.2.3.3.2 Évaluations basées sur des références

Pour évaluer notre système de TA V_N , nous avons choisi 500 bisegments à partir des données d'entraînement V_{N+1} comme données de test. Par exemple, nous avons fait l'AI avec les données « *Incrémental 10* », et choisi 500 bisegments des données d'entraînement « *Incrémental 11* » pour évaluer notre système. Les statistiques sur les données de test et le score BLEU sont données dans le Tableau 55.

Tableau 55 : Données de test et scores BLEU

	Longueur moyenne de segment (source)	Longueur moyenne de segment (cible)	Longueur moyenne de segment (résultat)	Nombre de mots (source)	Nombre de caractère (cible)	Nombre de caractère (résultat)	Pages standard (source)	Pages standard (ref)	Pages standard (résultat)	BLEU
Eval1	37,2	24,8	23,1	18612	12406	11544	74,4	31,0	28,9	0,1382
Eval2	40,6	25,8	27,6	20309	12912	13817	81,2	32,3	34,5	0,1428
Eval3	20,2	18,5	15,9	10098	9244	7974	40,4	23,1	19,9	0,1603
Eval4	20,7	18,3	16,3	10336	9141	8141	41,3	22,9	20,4	0,2357
Eval5	25,4	22,9	20,3	12687	11462	10165	50,7	28,7	25,4	0,2823
Eval6	19,1	17,3	15,4	9541	8646	7692	38,2	21,6	19,2	0,3346
Eval7	23,1	21,0	18,7	11543	10518	9356	46,2	26,3	23,4	0,3770
Eval8	20,7	19,3	17,2	10364	9651	8582	41,5	24,1	21,5	0,4285
Eval9	19,5	17,5	15,8	9725	8766	7901	38,9	21,9	19,8	0,4331
Eval10	21,8	19,8	17,8	10899	9921	8897	43,6	24,8	22,2	0,4597
Eval11	20,5	18,4	16,9	10248	9194	8459	41,0	23,0	21,1	0,4682
Eval12	23,0	20,5	18,5	11501	10251	9260	46,0	25,6	23,2	0,4694
Eval13	21,9	20,0	17,9	10973	10004	8964	43,9	25,0	22,4	0,4734
Eval14	14,5	13,0	11,3	7242	6481	5655	29,0	16,2	14,1	0,4792
Eval15	16,2	17,4	17,2	8107	8734	8622	32,4	21,8	21,6	0,4831

VI.3 Analyse des résultats

Après 16 itérations, les résultats des mesures montrent que la méthode d'apprentissage incrémental permet de réduire le temps de post-édition et d'améliorer les mesures basées sur les références (BLEU).

La V_0 de notre système avait été entraînée sur le corpus MULTUN, et ce système trop généraliste ne pouvait pas nous donner une bonne qualité de traduction. Avec 45mn/page_std, on avait en effet $Q=10\%=2/20$.

Cependant en l'utilisant comme point de départ, grâce à l'AI, nous sommes arrivé assez rapidement (16 itérations et environ 90h de calcul pour l'AI, sans jamais tout recompiler) à un système de bonne qualité (avec 15mn/page_std de PE en moyenne, notre formule donne $Q=70\%=14/20$).

Les temps de PE primaire (Tpe_1) et les scores BLEU sont donnés dans le Tableau 54 et dans le Tableau 55.

Partie C Contribution d'outils et de ressources

Introduction

Pendant ma thèse, j'ai participé au projet MUMIA⁷⁹ (*Multilingual and Multifaceted Interactive Information Access*). Tout d'abord, j'ai travaillé sur le corpus de brevets CLEF-IP 2011⁸⁰. J'ai construit les mémoires de traductions à partir de ce corpus pour les paires des langues allemand↔anglais, anglais↔français, et français↔allemand. Ensuite, j'ai construit plusieurs systèmes de TA avec ces mémoires de traductions. Enfin, j'ai extrait les segments monolingues dans ce corpus. Avec l'aide de Huanan SUN, étudiant de M1, j'ai utilisé les MAG dédiées à trois sites Web pour post-éditer ces segments monolingues à l'aide de nos systèmes de TA.

Dans le cadre de ce projet, on a demandé de mettre à disposition des participants à MUMIA, dont une bonne partie travaille sur la RI translingue sur les brevets.

Ce travail m'a montré l'intérêt de telles mises à disposition, non seulement pour leurs utilisateurs potentiels, mais aussi du point de vue de l'ingénierie linguistique. J'ai donc essayé de mettre à disposition le plus possible d'outils et de ressources dérivables de mon travail de thèse. Ce faisant, j'ai pu dégager un certain nombre de problèmes liés à ce type de tâche, et leur trouver des solutions assez génériques. Cette dernière partie est consacrée à cet aspect.

⁷⁹ <http://www.mumia-network.eu/index.php/the-action/objectives> et http://www.cost.eu/COST_Actions/ict/Actions/IC1002

⁸⁰ <http://ifs.tuwien.ac.at/~clef-ip/>

Chapitre VII Construction de systèmes de TA et support à la RI multilingue pour MUMIA

Résumé

Le corpus CLEF-IP 2011 contient 3,5 millions de documents de brevet ; il a été distribué dans le cadre du projet MUMIA. Nous en avons extrait les segments parallèles, et puis nous avons construit les mémoires de traductions correspondantes. Pour augmenter la quantité de segments multilingues, tout d'abord, nous avons filtré les fichiers XML pour supprimer les segments traduits, ensuite nous avons transformé les fichiers XML en format HTML, et les avons importé dans notre plate-forme SECTRA_w/IMAG. Enfin, nous avons fait la post-édition pour l'élargissement à d'autres langues.

VII.1 Contexte et motivations

VII.1.1 Description du projet MUMIA et du WG2

Le projet MUMIA (EU COST ACTION) a pour but de construire et mettre en réseau une communauté de chercheurs travaillant dans le domaine de l'accès à l'information multilingue, multimodale, et multifacette, et plus particulièrement aux documents de brevet. Il encourage la recherche et la communication de la technologie dans les domaines généraux de la traduction automatique, de la recherche d'information, et de l'accès multimodèle, multilingues à l'information interactive, visant à faciliter le développement de systèmes de recherche de prochaine génération.

Pour augmenter le synergie entre les communautés de RI et de TA, une série d'activités ont été organisées (par exemple réunion interdisciplinaire MUMIA GT à Tallinn et à Amsterdam, réunion WG/atelier ouvert à Hissar) pour discuter des ressources, et des outils linguistiques pour la recherche de développement des systèmes, et aussi pour soutenir la production des prochains systèmes de recherche multilingue.

Du point de vue de l'ACTION COST MUMIA, la principale motivation de ces activités était d'accroître la communication entre les scientifiques RI/TAL/TA sur les problèmes et les défis principaux suivant : (1) le développement/l'aide/l'intégration des ressources et des outils linguistiques pour le développement des systèmes de RI, et (2) la clarification des ressources et des solutions technologiques pour la mise en œuvre de RI multilingue dans les systèmes de recherche professionnels modernes.

Par exemple, un progrès technologique est le travail sur la production de corpus parallèles de haute qualité et sur l'évaluation permanente basée sur les tâches de plusieurs systèmes de TA pour les brevets. Une vitrine des corpus parallèles développés en utilisant la MT de CLEF-IP est disponible sur la page web de MUMIA. L'autre exemple est l'intégration de divers services de TA pour soutenir la recherche multilingue dans un système fédératif de RI sur les brevets (PERFEDPAT).

L'ACTION MUMIA COST contient cinq groupes de travail (GT) (*Integrating and Managing Language Resources, Processing Infrastructures for IR and MT, User Centered Aspects of MUMIA, Semantic Search and Faceted Search, Visualization, Distributed and Social Search*). J'ai participé au WG2. L'objectif de ce groupe de travail est l'étude de nouvelles infrastructures pour la recherche plus efficace dans les grands environnements numériques.

VII.1.2 PerFedPat et Khresmoi

Le système PERFEDPAT (Salampasis and Hanbury, 2014), basé sur le cadre EZDL, fournit des services basiques de RI en utilisant une méthode qui fédère l'accès à de multiples systèmes (en ligne) de RI concernant les brevets (actuellement ESPACENET⁸¹, GOOGLE PATENTS, PATENTSCOPE⁸² et MAREC⁸³), et qui repose sur un accès parallèle à de multiples sources de brevet.

PERFEDPAT cache la complexité à l'utilisateur qui utilise un outil unique pour interroger tous les ensembles de données de brevets. PERFEDPAT fournit des services comme la recherche par champ, la fusion, le regroupement et le filtrage des résultats. Il offre un soutien pour l'historique de la requête en cours et des sessions de recherche. La deuxième caractéristique innovante de PERFEDPAT est qu'il a une architecture enfichable (par plugins) et extensible.

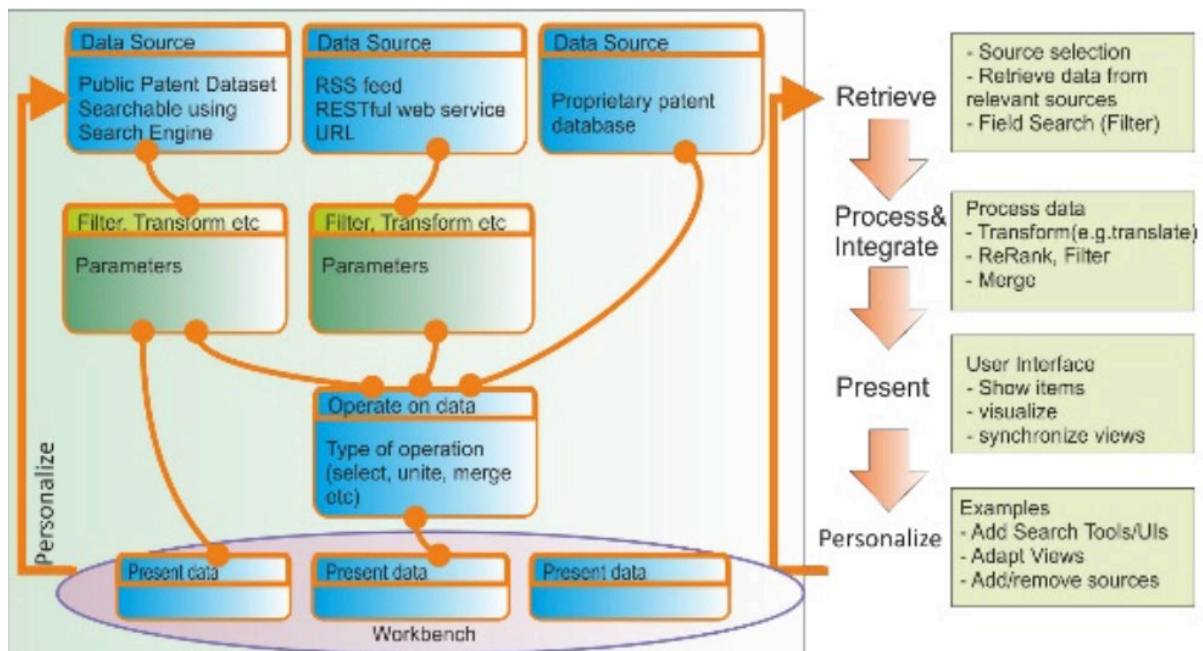


Figure 33 : Architecture de PerFedPat⁸⁴

Actuellement les outils intégrés concernent la recherche CIB (Classification internationale des brevets), la navigation à facettes, le regroupement des résultats de la recherche et la TA des requêtes. Le système est disponible en téléchargement à partir d'Internet⁸⁵.

KHRESMOI⁸⁶ (Hanbury et al., 2011) est un système de recherche multilingue/multimodèle, et d'accès d'information pour l'information biomédicale, construit dans le cadre du projet éponyme KHRESMOI. KHRESMOI utilise des données disponible en ligne, comme les sites de santé certifiés par *Health on the Net*⁸⁷ et des revues en accès ouvert.

⁸¹ <http://worldwide.espacenet.com>

⁸² <http://www.wipo.int/patentscope/fr/>

⁸³ <http://www.ifs.tuwien.ac.at/imp/marec.shtml>

⁸⁴ L'image est à partir du transparent de « PerFedPat patent search system »

⁸⁵ <http://www.perfedpat.eu/index.php/download-perfedpat>

⁸⁶ <http://www.khresmoi.eu>

⁸⁷ <https://www.healthonnet.org>



Figure 34 : KHRESMOI

VII.1.3 Objectif poursuivi

Nous travaillons pour supporter la recherche en RI multilingue sur les brevet. Les systèmes de recherche existant, comme PERFEDPAT et KHRESMOI, utilisent des systèmes de TA pour réaliser la fonction de recherche en contexte multilingue. La qualité du système de TA détermine la qualité du résultat de la recherche. Par exemple, PERFEDPAT utilise GT pour traduire les requêtes, mais le résultat n'est pas satisfaisant, et il est impossible d'améliorer la qualité de traduction.

Un document de brevet contient des « segments multilingues » au sens défini plus haut (p.22) : un élément XML correspondant à un segment contient ses versions en plusieurs langues). Si nous pouvions construire des systèmes de TA (STA) avec ces segments parallèles, nous pourrions améliorer la qualité de la RI translingue. C'est ce que nous avons fait.

VII.2 Construction de MT et de STA à partir des corpus CLEF-IP 2011

VII.2.1 Description du corpus CLEF-IP

La collection CLEF-IP contient des brevets, physiquement stockés comme une collection de fichiers XML encodant des documents de brevets. Un document de brevet peut être un document de demande de brevet, un rapport de recherche, ou un document de brevet accordé. À chaque document de brevet est attribué un code de type, qui apparaît comme un suffixe attaché à l'identifiant du brevet (par exemple, EP-nnnnnnn-A1, WO-nnnnnnnnn-A2). Dans le cas de l'OEB (EPO), les documents de demande de brevet qui incluent un rapport de recherche portent le code A1, les rapports de recherche de demandes de brevet portent le code A3, les documents de brevets accordés portent le code B1, etc.⁸⁸.

La collection de données CLEF-IP 2011 est basée sur les données de 2010, et est extraite à partir du corpus de données MAREC. La collection CLEF-IP contient principalement des documents brevets publiés par l'OEB (EPO).

Deux ajouts importants ont été apportés au corpus de la collection 2010. Le premier a été d'inclure dans le corpus distribué certains brevets publiés par l'Organisation mondiale de brevets intellectuels OMPI (WIPO). Un pourcentage élevé des brevets de l'OEB (EPO) figurant dans le corpus CLEF-IP sont des demandes de brevet déposées à l'échelle internationale dans le cadre du Traité de coopération en matière de brevets (PCT), auquel cas l'OEB ne republie pas la demande de brevet en totalité, mais seulement une entrée bibliographique faisant référence à la demande originale. Pour ces brevets, CLEF-IP 2011 a ajouté leur équivalent de l'OMPI à la

⁸⁸ <https://register.epo.org/help?topic=kindcodes>

collection, afin de fournir une collection qui soit à la fois plus grande et plus réaliste. Le deuxième ajout à la collection CLEF-IP concerne l'une des nouvelles tâches basées sur les images contenues dans les brevets, à savoir la tâche de recherche basée sur les images des brevets. Pour cette tâche, nous avons ajouté à la collection cible CLEF-IP les images des brevets pour trois classes d'ICP: A43B, A61B, et H01L.

Dans nombre de documents, ajouter les documents de brevets de l'OMPI (WIPO) au corpus de la collection l'a augmentée de 1,2 million de documents de brevets, arrivant à un nombre final d'environ 3,5 millions de documents XML, correspondant à environ 1,5 million de brevets. Les images correspondant aux 47 000 documents XML dans les trois classes de la CIB (ICP) dans la tâche IMG PAC occupent 5,4 Go, pour 290 880 fichiers TIFF.

Comme dans les années précédentes, le corpus de la collection de test a été livré aux participants sans fusionner les documents relatifs à un même brevet en un seul document. Chaque brevet est identifié par un numéro de brevet unique, une chaîne identifiant le bureau éditeur («EP» pour l'OEB et «WO» pour l'OMPI), suivie par une série de chiffres. Pour chaque brevet, il y a un répertoire contenant les documents XML représentant les documents de brevet liés à ce brevet. Pour les brevets de l'EP, le format du nom est 00000n/nn/nn/nn/*.xml, où la séquence de chiffres dans le nom du répertoire correspond à celle du nom du brevet. Pour les brevets WO, le format du nom est 00000n/nn/nn/nn/*.xml, où les 4 premiers chiffres (après « 00 ») représentent l'année de publication du document.

Par exemple, le brevet EP 09 81 2 01 correspond au répertoire contenant les fichiers EP-0981201-A2.xml, EP-0981201-A3.xml, et EP-0981201-B1.xml:

EP/000000/98/12/01

Au brevet WO 1994030029 correspond le répertoire contenant le fichier WO-1994030029-A1.xml:

WO/001994/03/00/29 → WO-1994030029-A1.xml

Les fichiers d'images des brevets sont stockés comme des fichiers TIF dans un dossier séparé, la correspondance entre le fichier d'images et son fichier XML est établie par un petit ensemble de règles.

Tous les documents textuels de la collection CLEF-IP contiennent les principaux champs XML suivants: données bibliographiques, des résumés, résumé et revendications. Tous les documents n'ont pas nécessairement de contenu dans ces champs. Le contenu des différents champs XML peut être de l'anglais, de l'allemand ou du français. Certains champs peuvent apparaître plus d'une fois, chaque fois avec une langue différente. Un fichier de brevet XML a également une langue du document (anglais, allemand ou français), ce qui n'exclut pas que ses sous-champs apparaissent avec un attribut de langue différent de la langue du document. Par exemple, les documents de brevet EP (EP-nnnnnn-Bn.xml) accordé doivent contenir les revendications en trois langues (anglais, allemand et français).

La dernière version de la collection est la même que celle utilisée dans le CLEF-IP 2011 LAB (le corpus de données utilisé en 2012 et 2013 est le même que celui utilisé en 2011), de sorte que notre travail est basé sur la collection CLEF-IP 2011.

```

▼<patent-document ucid="EP-0071719-B1" country="EP" doc-number="0071719" kind="B1" lang="EN" date="19881102" family-id="23112365"
date-produced="20100220" status="new">
  ▼<bibliographic-data>
    ▼<publication-reference fvid="19385621" ucid="EP-0071719-B1" status="new">
      ▼<document-id status="new" format="original">
        <country status="new">EP</country>
        <doc-number>0071719</doc-number>
        <kind>B1</kind>
        <date>19881102</date>
        <lang>EN</lang>
      </document-id>
    </publication-reference>
    ▼<application-reference mxw-id="PAPP16307477" ucid="EP-82105214-A" load-source="docdb" status="new" is-representative="NO">
      ▼<document-id format="epo" status="new">
        <country status="new">EP</country>
        <doc-number>82105214</doc-number>
        <kind>A</kind>
        <date>19820615</date>
        <lang>EN</lang>
      </document-id>
    </application-reference>
    ▼<priority-claims status="new">
      ▼<priority-claim mxw-id="PPC19696867" ucid="US-28963181-A" status="new">
        ▼<document-id format="epo" status="new">
          <country status="new">US</country>
          <doc-number>28963181</doc-number>
          <kind>A</kind>
          <date>19810803</date>
        </document-id>
      </priority-claim>
    </priority-claims>
  </bibliographic-data>
</patent-document>

```

Figure 35 : Exemple de fichier XML Dans CLEF-IP

VII.2.2 Extraction de MT à partir de CLEF-IP 2011

VII.2.2.1 Analyse du corpus

Nous avons analysé les brevets en ce qui concerne la structure de leurs champs XML. Nous avons constaté que quatre grands champs peuvent avoir des segments parallèles: *<invention-title>*, *<abstract>*, *<description>*, et *<claims>*. Chaque champ peut avoir des sous-champs, par exemple, un champ *<claims>* peut contenir 6 sous-champs *<claim>* dans *EP-0260000-B1.xml* (Figure 36).

```

▼<claims load-source="patent-office" status="new" mxw-id="PCLM9874066" lang="DE">
  ▶<claim num="1">...</claim>
  ▶<claim num="2">...</claim>
  ▶<claim num="3">...</claim>
  ▶<claim num="4">...</claim>
  ▶<claim num="5">...</claim>
  ▶<claim num="6">...</claim>
</claims>

```

Figure 36 : Exemple de champ *<claims>* contenant 6 sous-champs *<claim>* dans *EP-0260000-B1.xml*

Nous commençons par ces champs, en cherchant des champs qui apparaissent plus d'une fois dans le document de brevets et tels que chaque champ ait un attribut de langue différente. Par exemple, Figure 37 montre un champ *<invention-title>* avec 3 attributs de langue différents (*lang = "DE"*, *lang = "EN"*, et *lang = "FR"*). Chaque champ contient également du contenu, dans la langue qui correspond à son attribut.

```

▼<invention-title lang="DE" load-source="ep" status="new">
  Verfahren zur Herstellung von Polyimidestern der Trimellitsäure
</invention-title>
▼<invention-title lang="EN" load-source="ep" status="new">
  PROCESS FOR THE FABRICATION OF POLYIMIDE-ESTERS OF TRIMELLITIC ACID
</invention-title>
▼<invention-title lang="FR" load-source="ep" status="new">
  Procédé pour la fabrication de polyimide-esters de l'acide trimellitique
</invention-title>

```

Figure 37 : Exemple d'un champ `<invention-title>` avec 3 attributs de langue différents et les contenus correspondants en 3 langues différentes

Il est bien connu que les corpus parallèles ont aussi un problème avec la direction de la traduction, car la relation de traduction est symétrique au niveau des termes, mais pas des phrases. Si une MT contient des directions de traduction différentes, cela affecte directement notre travail. Lorsque nous traitons la collection CLEF-IP 2011, nous devons donc également tenir compte de la langue source de chaque document de brevet. Chaque document XML de CLEF-IP 2011 comporte une indication de sa langue de départ. Par exemple, la Figure 38 montre un document de brevet en langue anglaise. Au cours de notre processus d'extraction, nous considérons la langue du document comme état la langue source de tous ses segments.

```

<patent-document ucid="EP-0071719-B1" country="EP" doc-
number="0071719" kind="B1" lang="EN" date="19881102" family-
id="23112365" date-produced="20100220" status="new">

```

Figure 38 : Un champ `<patent-document>` avec attribut `lang = "EN"`

VII.2.2.2 Traitement du corpus

(Utiyama and Isahara, 2007) ont utilisé les parties du champ de description "description détaillée des modes de réalisation préférés" (Detailed Description of the Preferred Embodiments) et "Contexte de l'invention" de chaque brevet pour trouver des segments parallèles (japonais-anglais), car ils ont constaté que ces deux parties ont plus de traductions que d'autres. Parce qu'ils avaient moins de paires de brevets, (Lu et al., 2009) ont utilisé toutes les parties des documents de brevet pour trouver des segments parallèles (chinois-anglais). Dans notre travail, nous avons extrait toutes les parties des documents de brevet, mais dans le but d'assurer la qualité du corpus parallèle, nous avons rendu le champ `<invention-title>` et les parties `<claims>` disponibles dans la première version du corpus parallèle CLEF; les autres parties du corpus parallèle seront disponibles dans la prochaine version.

Notre travail est basé sur 3,5 millions de documents de brevet (fichiers XML), et nous voulons en extraire autant de segments parallèles utiles que possible. Tout d'abord, nous parcourons chaque document de brevet. Pour chaque document de brevet, nous sélectionnons la langue source à partir du champ `<patent-document>`, selon l'attribut de langue de ce champ. Deuxièmement, nous recherchons les segments parallèles contenus dans les quatre champs principaux (`<invention-title>`, `<abstract>`, `<description>`, et `<claims>`). Parfois, certains champs ont un attribut de langue différent de la langue du document. Par exemple, dans `EP-0260000-B1.xml`, l'anglais est la langue du document, mais `<claims>` segments ne existent pas en anglais, seules les versions allemandes et françaises sont disponibles. Même s'il est toujours souhaitable de collecter autant de texte que possible, il est encore plus important de veiller à la qualité des textes, de sorte que, dans ce cas, nous ne stockons pas les parties en allemand et en français comme un segment parallèle, parce que nous ne savons pas laquelle est la source.

Tous les champs qui apparaissent plus d'une fois dans un document de brevet et qui ont différents attributs de langue sont traités comme une collection. En général, un document de brevet OEB (IPO) a un maximum de 3 langues (anglais, français et allemand). Nous avons choisi comme segment source un segment dont l'attribut langue est compatible avec la langue source, puis avons ensuite extrait le segment parallèle cible à partir des autres champs. Par exemple, dans EP-0301015-B1.xml, la langue source est l'anglais, et `<revendications>` champ apparaît 3 fois. Donc, nous utilisons la partie anglaise des champs de revendications comme segments source, et nous considérons les parties en français et en allemand comme les segments cibles. Le segment source et les segments cibles sont ensuite stockés séparément dans des fichiers différents. Dans l'exemple ci-dessus, le segment source a été stocké dans *CLEF_claims_en-fr.en* et *CLEF_claims_en-de.en*, et les segments cibles dans *CLEF_claims_en-fr.fr* et *CLEF_claims_en-de.de*, respectivement. Afin de réduire le bruit dans les données, nous ne gardons que le texte extrait, et enlevons toutes les balises.

Toutes les données extraites ne sont pas entièrement adaptées à une utilisation directe pour les applications de TAL (NLP). Nous devons nettoyer les données extraites et éliminer un peu de bruit. Pour l'alignement, nous avons utilisé LF ALIGNER⁸⁹, un outil open-source basé sur HUNALIGN (Varga et al., 2007) développé par András Farkas, qui, surtout, a la couverture linguistique la plus large (un total de 32 langues), et permet la génération automatique des dictionnaires dans une combinaison quelconque de ces langues. Les segments alignés sont préparés de façon bilingue pour 4 types (titre, résumé, description et revendications), et toutes les 6 paires de langues (de_en, de_fr, en_de, en_fr, fr_de, fr_en).

Le Tableau 56 montre le nombre de segments et de mots qui sont extraits à partir des champs de titre et de revendications en source et en cible après l'alignement des segments. Toutes les phrases parallèles extraites sont enregistrées dans les formats TMX et TXT, et peuvent être trouvées à <http://membres-liglab.imag.fr/wang/downloads>.

Tableau 56 : Nombre de segments extraits comme source et cible après l'alignement de segments dans les champs `<title>` et `<claims>`

Paires de langues		Titre		Revendications	
		Segments	Mots	Segments	Mots
de-en	de	311,298	2,038,785	1,696,498	62 M
	en		2,582,703		71 M
de-fr	de	311,184	2,036,112	1,661,419	79 M
	fr		2,482,257		86 M
en-de	en	884,759	6,661,481	5,218,024	332 M
	de		5,508,289		296 M
en-fr	en	884,727	6,661,322	5,373,452	330 M
	fr		8,538,012		380 M
fr-de	fr	106,211	963,508	572,356	36 M
	de		1,204,439		37 M
fr-en	fr	106,246	1,285,467	586,498	38 M
	en		1,048,374		37 M

⁸⁹ <http://sourceforge.net/projects/aligner/>

VII.2.3 Construction des systèmes de TA

Nous avons utilisé notre corpus parallèle extrait (le titre et les champs de revendication) pour construire des systèmes de TA avec MOSES. Tout d'abord, pour la préparation de l'ensemble de développement et de l'ensemble de test, nous avons extrait 2000 phrases pour la le réglage des poids de fonctionnalités de Moses, et extrait 1000 phrases pour les tests. Ensuite, nous avons utilisé le reste pour former les modèles de traduction de Moses. Nous avons en fait construit des systèmes de TA pour seulement 3 directions: de-fr, fr-de, et en-fr.

Les systèmes comprennent également des modèles de langage 5-grammes formés sur le côté cible correspondant des textes parallèles à l'aide de IRSTLM. Les poids de fonctionnalités requises par le décodeur de MOSES ont encore été déterminés avec MERT en optimisant les scores BLEU sur l'ensemble de développement (1000 phrases). Les ensembles de test ont été traduits par les systèmes qui en résultent et ensuite utilisés pour évaluer les systèmes en termes de scores BLEU, comme indiqué dans Tableau 57.

Tableau 57 : Scores BLEU des systèmes de TA tirés de CLEF-IP

Paires de langues	Jeu de développement	Jeu de test
de-fr	35.41	28.72
en-fr	42.59	38.82
fr-de	34.85	30.14

VII.3 Expérimentation et élargissement à d'autres langues

VII.3.1 Reconstruction de trois sites Web de brevets monolingues

Le corpus de brevets CLEF-IP contient certain nombre segments monolingues, par exemple, dans la Figure 39, ce fichier XML est en anglais (dans la balise <document>, l'attribut de langue est lang="EN"), et son titre anglais est traduit en français et en allemand, mais le champ <claims> n'est pas traduit. CLEF-IP contient environ 500K fichiers monolingues.

```
<patent-document ucid="EP-0000004-B1" country="EP" doc-number="0000004" kind="B1" lang="FR" date="19800903" family-id="9191604" date-produced="20100220" status="new">
  <bibliographic-data>
    <publication-reference fvid="19066349" ucid="EP-0000004-B1" status="new" /> </publication-reference>
    <application-reference mw-id="PAPP16167232" ucid="EP-78100037-A" load-source="docdb" status="new" is-representative="NO" /> </application-reference>
    <priority-claims status="new" /> </priority-claims>
    <dates-of-public-availability status="new" /> </dates-of-public-availability>
  </bibliographic-data>
  <technical-data status="new">
    <classification-ipc status="new" /> </classification-ipc>
    <classifications-ipcr /> </classifications-ipcr>
    <classification-ecla status="new" /> </classification-ecla>
    <invention-title mw-id="PT72447844" lang="FR" load-source="patent-office" status="new">
      Dispositif de filtration centrifuge continue pour séparation, lavage ou époussement, et application à une machine à café automatique
    </invention-title>
  </technical-data>
  <parties /> </parties>
  <international-convention-data /> </international-convention-data>
</bibliographic-data>
<copyright>
  User acknowledges that the Information Retrieval Facility (IRF) and its third party providers retain all right, title and interest in and to this xml under applicat
  copyright laws. User acquires no ownership rights to this xml including but not limited to its format. User hereby accepts the terms and conditions of the Licence
  Agreement set forth at http://www.ir-facility.org/legal/marec/data_licence
</copyright>
</patent-document>
```

Figure 39 : Exemple de fichier XML monolingue

L'expression des phrases du document de brevet est très similaire dans un document, et beaucoup de mots-clés sont présentés plusieurs fois. Par exemple, nous voyons les revendications du fichier XML *EP0203923B1.xml*. Chaque revendication dans le champ `<claims>` change quelques mots entre eux, mais la structure de phrase est presque la même.

CLAVIER A TOUCHES CAPACITIVES SELON LA R1 CARACTERISE EN CE QUE LE CHOIX DE FREQUENCES ELEVEES (SUPERIEURES A 500 KHZ) POUR LES OSCILLATEURS DE TOUCHE PERMET UN FONCTIONNEMENT FIABLE JUSQU'A UNE EPAISSEUR DE LA PAROI DIELECTRIQUE (2) ATTEIGNANT 30 MM DE VERRE COURANT.
CLAVIER A TOUCHES CAPACITIVES SELON LES R1, R2, R3 CARACTERISE EN CE QU'UN REGLAGE IMMEDIAT ET SIMPLE PERMET D'ADAPTER LE FONCTIONNEMENT DU CLAVIER A L'EPAISSEUR DE LA PAROI DIELECTRIQUE (2).
CLAVIER A TOUCHES CAPACITIVES SELON LES R1 ET R2 CARACTERISE EN CE QUE LA CIRCUITERIE ASSOCIEE AUX TOUCHES COMPREND UN COMPTEUR D'ADRESSE, UN DECODEUR ET LES OSCILLATEURS, ET PEUT ETRE RELIEE AU MICROPROCESSEUR PAR UN CABLE LIMITE A 5 FILS ET POUVANT ATTEINDRE 10 METRES DE LONGUEUR.

Figure 40 : Exemple de revendication dans le fichier EP0203923B1.xml

Un système de TA entraîné avec notre MT pourra fournir des traductions de haute qualité, ce qui facilitera beaucoup la post-édition des segments ne contenant pas leur version dans la langue cible. Nous avons post-édité les segments monolingues et bilingues (nous avons 3 langues au total) pour agrandir la MT.

D'abord, nous filtrons les fichiers XML du corpus CLEF-IP, et utilisons *JSoup*⁹⁰ pour analyser la langue source. Nous trouvons la langue dans le champ `<document>`. Ensuite, nous cherchons les balises qui contiennent l'attribut `lang`, et où la valeur de l'attribut `lang` est différente de la langue source du document. Enfin, nous supprimons les balises et leurs contenus. Par exemple, dans le fichier EP-0000004-B1.xml, l'attribut `lang` dans le champ `<document>` est l'anglais. Les trois champs `<title>` sont en anglais, en français, et en allemand, et le champ `<claims>` (revendications) est en anglais. Donc nous supprimons le champ `<title>` en français et en allemand, et conservons le champ `<claims>`.

Après le prétraitement, il faut adapter nos fichiers XML monolingues à SECTRA_w/IMAG. Le fichier XML doit avoir l'aide de l'outil pour la visualisation de contexte. Nous le transformons au format HTML. Pour conserver la structure XML du document de brevet, nous reconstruirons les balises du fichier XML, et nous utilisons les balises HTML/CSS pour décorer les balises XML (avec les couleurs et le style). Dans le Figure 41, on donne un exemple de fichier XML monolingue présenté en format HTML.

⁹⁰ <http://jsoup.org>


```

<patent-document uci="EP-0000004-B1" country="EP" doc-number="0000004" kind="B1" lang="FR" date="19800903" family-id="9191604" date-produced="20100220"
status="new">
  <bibliographic-data>
    <publication-reference fvid="19066349" uci="EP-0000004-B1" status="new" /> </publication-reference>
    <application-reference mw-id="PAPP16167232" uci="EP-78100037-A" load-source="docdb" status="new" is-representative="NO" /> </application-reference>
    <priority-claims status="new" /> </priority-claims>
    <dates-of-public-availability status="new" /> </dates-of-public-availability>
  </bibliographic-data>
  <technical-data status="new">
    <classification-ipc status="new" /> </classification-ipc>
    <classifications-ipc /> </classifications-ipc>
    <classification-ecla status="new" /> </classification-ecla>
    <invention-title mw-id="PT72447844" lang="FR" load-source="patent-office" status="new">
      Dispositif de filtration centrifuge continue pour séparation, lavage ou épuisement, et application à une machine à café automatique
    </invention-title>
  </technical-data>
  <parties /> </parties>
  <international-convention-data /> </international-convention-data>
</bibliographic-data>
<copyright>
  User acknowledges that the Information Retrieval Facility (IRF) and its third party providers retain all right, title and interest in and to this xml under applicat
  copyright laws. User acquires no ownership rights to this xml including but not limited to its format. User hereby accepts the terms and conditions of the Licence
  Agreement set forth at http://www.ir-facility.org/legal/marec/data_licence
</copyright>
</patent-document>

```

Figure 41 : Exemple de fichier HTML décoré

VII.3.2 Accès multilingue en utilisant les systèmes de TA créés

Nous avons construit 3 sites Web pour les fichiers monolingues, et chaque site Web contient tous les fichiers monolingues traités. Nous avons aussi créé 3 iMAG dédiées à ces 3 sites Web monolingues. Le site Web anglais est présenté dans AXIMAG. Ces 3 iMAG dédiées partagent la MT CLEF-IP. Nous importons les segments parallèles dans la MT CLEF-IP. Quand une iMAG dédiée à un site Web monolingue de CLEF-IP (une page Web monolingue) est accédée par l'utilisateur, l'iMAG d'abord cherche d'abord sa traduction dans la MT CLEF-IP. Si elle existe, il l'affiche dans la page Web traduite, sinon elle appelle nos systèmes de TA via Tradoh.

Nous avons entraîné 3 systèmes de TA (anglais→français, français→anglais, et allemand→français) en utilisant les MT extraites du corpus CLEF-IP. Pour fournir le service de traduction, nous avons intégré le plugin MOSES dans TRADOH, puis traduit les segments monolingues.

VII.3.3 Accès multilingue utilisant d'autres systèmes et pour d'autres langues

Sur SECTRA, nous pouvons ajouter des différents systèmes de TA pour traduire les segments source. Il nous aide à évaluer la qualité des différents systèmes de TA.

VII.3.3.1 Comparaison sur les mêmes langues

Pour comparer la qualité du système de TA sur les mêmes langues, nous utilisons GT à retraduire les 3 sites Web monolingues. Les segments traduits sont sauvegardés dans la MT CLEF-IP, et ils sont présentés dans SECTRA.

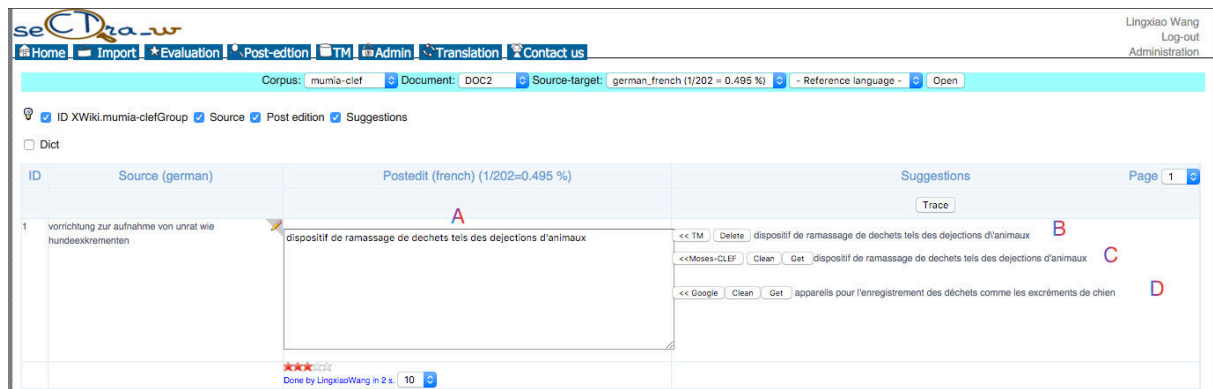


Figure 44 : PE en mode avancé, avec pseudo-trace montrant les différences entre les sorties de TA, la post-édition (utilisée comme référence), et la MT.

Dans le carré rouge de la Figure 44, il y a 4 segments. Le premier segment (A) est la post-édition de l'utilisateur, le deuxième segment (B) est la MT de CLEF-IP, le troisième (C) est la traduction du système de TA entraîné avec la MT CLEF-IP, et le dernier (D) est la traduction de GT. Quand on clique sur le bouton « Trace », la distance entre les segments et la référence peut être affiché dans la colonne « Suggestion ».

On remarque que le résultat de traduction de GT ne correspondant pas au segment source (). Il est très différent de la référence. Si l'utilisateur fait la post-édition sur ce résultat, il aura besoin de plus de temps.

VII.3.3.2 Utilisation pour d'autres langues

Nous ajoutons les nouvelles langues dans la MT de CLEF-IP. Nous utilisons GT et les systèmes de TA pour traduire les segments monolingues vers d'autres langues.

Les pages Web sont segmentées, puis leurs segments sont sauvegardés dans la MT. Pour ajouter un nouveau système de TA ou ajouter une nouvelle langue, nous retraduisons tous les segments sous SECTRA. Par exemple, pour ajouter la traduction chinoise, on coche notre système de TA français→chinois MOESLIG dans la colonne « MTs ». La Figure 45 montre comment nous retraduisons les segments du français vers le chinois dans le pseudo-document DOC6.

seCTra_w

Home Import Evaluation Post-edition TM Admin Translation Contact us

Corpus name: mumia-clef

All	MTs	Pivot language	Target language	Translation selection
<input type="checkbox"/>	<input type="checkbox"/> Google <input type="checkbox"/> Systran <input checked="" type="checkbox"/> MosesLIG	No pivot language	-select-	Start

	No	Document name	Source language	Pivot language	Target language	Operate
<input type="checkbox"/>	1	DOC1	german	No pivot language	-select-	Start
<input type="checkbox"/>	2	DOC2	german	No pivot language	-select-	Start
<input type="checkbox"/>	3	DOC3	german	No pivot language	-select-	Start
<input type="checkbox"/>	4	DOC4	german	No pivot language	-select-	Start
<input type="checkbox"/>	5	DOC5	german	No pivot language	-select-	Start
<input checked="" type="checkbox"/>	6	DOC6	french	No pivot language	Chinese	Start
<input type="checkbox"/>	7	DOC7	german	No pivot language	-select-	Start
<input type="checkbox"/>	8	DOC8	german	No pivot language	-select-	Start
<input type="checkbox"/>	9	DOC9	english	No pivot language	-select-	Start

Figure 45 : Retraduction des segments du français vers le chinois pour DOC6 avec le système de TA français→chinois MosesLIG

Chapitre VIII Mise à disposition de ressources

Résumé

Dans le cadre de ma thèse, je me suis aussi attaché à mettre à disposition de la communauté des chercheurs les ressources et outils que j'ai développés, sous des formes « statiques » ou « dynamiques ».

Les ressources et outils « statiques » sont mis à disposition sous forme de fichiers téléchargeables.

La mise à disposition dynamique consiste à proposer des serveurs Web d'utilisation de systèmes de TA et de post-édition collaborative de MT associées à des sites Web, via les iMAG associées.

Introduction

Nos ressources « statiques » sont particulièrement intéressantes à cause de leur taille, et parce qu'elles sont utilisables pour construire ou expérimenter des systèmes de TA.

Nous avons mis à disposition des corpus multilingues post-édités dans des formats usuels comme TXT et XML. Ces formats sont bien adaptés aux systèmes existants : les segments source et les segments post-édités sont alignés, et sauvegardés dans deux fichiers TXT. Ils sont utilisables directement par les boîtes à outils comme MOSES et JOSHUA pour entraîner des systèmes de TA. Le format XML peut être importé dans plusieurs plates-formes comme DEJA VU, MEMOQ, SDL TRADOS, etc.

Une grande partie de ces ressources statiques a été construite en utilisant la plate-forme SECTra_w/iMAG. Les segments post-édités sont munis de l'information associée, comme le temps de PE, le système de TA initial, la référence, etc. L'utilisateur en a besoin pour certains travaux, l'évaluation de systèmes de TA, par exemple.

Les MT tirées de la collection de brevets CLEF-IP sont les plus grosses (17,5M segments en fr-de-en). Suivent les MT fr-zh : 9000 pour le sous-langage d'EDF, 15000 pour celui du LIG, 23000 pour les documents pédagogiques de MACAU-UJF, et 10000 pour le roman de Jules Verne « Voyage au centre de la Terre ».

Les systèmes de TA téléchargeables sont ceux déduits de CLEF-IP (fr-de, de-fr, fr-en) et quelques variantes de notre système fr-zh MOSESLIG-FR-ZH.

Les services Web en ligne sont ceux des quelques iMAG parmi lesquelles deux servent actuellement à construire de bonnes MT du français vers le somali et vers le comorien.

VIII.1 Contribution de ressources statiques sous forme de MT

VIII.1.1 Formats choisis

Nous avons choisi 2 types de format pour sauvegarder et partager nos mémoires de traduction, nous avons choisi 2 formats documentaires, TXT et XML.

TXT est le format plus utilisable pour l'apprentissage du système de TA. La boîte à l'outil, comme MOSES et JOSHUA, ils supportent des données en format txt (*raw text*) à entraîner les modèles. Figure 46 présente un exemple de la MT aligné CLEF-IP anglais-français. Voir l'exemple avec 100 bisegments anglais-français dans l'Annexe 12.

en-fr.en	en-fr.fr
1 An additive for a lubricating oil or hydrocarbon fuel obtainable by a process which comprises	1 Additif pour une huile lubrifiante ou un combustible hydrocarboné, pouvant être obtenu par un
2 An additive as claimed in Claim 1, wherein R ₄ is alkylene of 2 or 3 carbon atoms.	2 Additif suivant la revendication 1, dans lequel R ₄ représente un groupe alkylène ayant 2 ou 3
3 An additive as claimed in Claim 1 or 2, wherein R _s is hydrogen or alkyl of from 1 to 10 carbon	3 Additif suivant la revendication 1 ou 2, dans lequel R _s représente l'hydrogène ou un groupe al
4 An additive as claimed in Claim 1, 2 or 3, wherein W and X are both oxygen.	4 Additif suivant la revendication 1, 2 ou 3, dans lequel W et X représentent l'un et l'autre l'
5 An additive as claimed in Claim 1, 2 or 3, wherein W is sulfur and X is oxygen.	5 Additif suivant la revendication 1, 2 ou 3, dans lequel W représente le soufre et X représente
6 An additive as claimed in Claim 1, 2 or 3, wherein W and X are both sulfur.	6 Additif suivant la revendication 1, 2 ou 3, dans lequel W et X représentent l'un et l'autre le
7 An additive as claimed in any preceding claim, wherein the reaction is conducted at from 0°C t	7 Additif suivant l'une quelconque des revendications précédentes, dans lequel la réaction est c
8 An additive as claimed in any preceding claim, wherein the molar charge of the compound of For	8 Additif suivant l'une quelconque des revendications précédentes, dans lequel le rapport molaire
9 A lubricating oil composition comprising an oil of lubricating viscosity and 0.2 to 10 percent	9 Composition d'huile lubrifiante, comprenant une huile de viscosité propre à la lubrification e
10 A fuel composition comprising a hydrocarbon boiling in the gasoline or diesel range and an add	10 Composition de combustible comprenant un hydrocarbure bouillant dans la plage de l'essence o
11 A fuel composition as claimed in Claim 10, wherein the additive is present in an amount of fro	11 Composition de combustible suivant la revendication 10 dans laquelle l'additif est présent en
12 An output apparatus for outputting information to a recording medium of the kind which compris	12 Appareil périphérique de sortie destiné à délivrer en sortie une information vers un support d
13 An output apparatus according to claim 1, wherein said first transfer means comprises a clutch	13 Appareil périphérique de sortie selon la revendication 1, dans lequel lesdits premiers moyens
14 An output apparatus according to claim 1, wherein said third feed mechanism also is operable t	14 Appareil périphérique de sortie selon la revendication 1, dans lequel ledit troisième mécanism
15 An apparatus according to claim 1, further comprising key means (5) for inputting data and con	15 Appareil selon la revendication 1, comportant en outre des moyens à touche (5) destinés à intr
16 An apparatus according to claim 4, wherein said control means (201) controls said motor (19) t	16 Appareil selon la revendication 4, dans lequel lesdits moyens de commande (201) commandent led
17 An apparatus according to claim 4, wherein said control means (201) is divided into a keyboard	17 Appareil selon la revendication 4, dans lequel lesdits moyens de commande (201) sont divisés e
18 An apparatus according to claim 1, wherein said first transfer means comprises a clutch mechan	18 Appareil selon la revendication 1, dans lequel lesdits premiers moyens de transfert comprennent
19 An apparatus according to claim 5, wherein the predetermined amount of angle of rotation is ab	19 Appareil selon la revendication 5, dans lequel la valeur prédéterminée de l'angle de rotation
20 An apparatus according to claim 5, wherein the direction of rotation of said motor (19) is a d	20 Appareil selon la revendication 5, dans lequel le sens de la rotation dudit moteur (19) est u
21 An apparatus according to claim 1, wherein said apparatus can erase one character by hitting i	21 Appareil selon la revendication 1, dans lequel ledit appareil peut effacer un caractère en le
22 An apparatus according to claim 1, wherein said first feed mechanism (1, 24, 33) feeds said re	22 Appareil selon la revendication 1, dans lequel ledit premier mécanisme d'alimentation (1, 24,
23 An apparatus according to any previous claim wherein said motor (19) is a stepper motor.	23 Appareil selon l'une quelconque des revendications précédentes, dans lequel ledit moteur (19)
24 Apparatus for generating from an input analog signal having first and second signal elements t	24 Appareil destiné à produire, à partir d'un signal analogique d'entrée possédant des premier et
25 Apparatus as claimed in claim 1 wherein said storing means (14) is a random access memory and	25 Appareil selon la revendication 1, où ledit moyen d'empaasage (14) est une mémoire vive et
26 A semipermeable microcompartment which is artificially prepared by reassembly of proteinaceous	26 Microcompartment semipermeable, qui est préparé artificiellement par réassemblage de macromol
27 A microcompartment according to claim 1 wherein substantially all of said macromolecules are d	27 Microcompartment selon la revendication 1 ou 2, caractérisé en ce qu'il présente la forme gén
28 A microcompartment according to claim 1 or claim 2 which has the shape generally of a sphere.	28 Microcompartment selon la revendication 1 ou 2, caractérisé en ce qu'il présente la forme gén
29 A microcompartment according to claim 1 or claim 2 which has the shape generally of a sphere i	29 Microcompartment selon la revendication 1 ou 2, caractérisé en ce qu'il présente la forme gén
30 A microcompartment according to any of claims 3 to 5 which has an overall diameter in the rang	30 Microcompartment selon l'une quelconque des revendications 3 à 5, caractérisé en ce qu'il pré

Figure 46 : Exemple de données en format TXT (MT CLEF-IP anglais-français)

Pour adapter aux système de l'aide de traduction comme SDL TRADOS et MEMOQ, nous avons choisi le format XML (fichier TMX ou fichier XML avec fichier accompagne) pour partager ces données (Figure 47). Voir l'exemple avec 20 bisegments anglais-français en format TMX dans l'Annexe 14.

```

<!DOCTYPE tmx SYSTEM "tmx14.dtd">
<tmx version="1.4">
  <header creationtool="LF Aligner" creationtoolversion="3.11"
    datatype="patent document" segtype="sentence" adminlang="EN"
    srclang="EN"
    o-tmf="TW4Win 2.0 Format"
  >
  </header>
  <body>
    <tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop
      type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>
    <tuv xml:lang="EN"><seg>An additive for a lubricating oil or hydrocarbon fuel obtainable
      by a process which comprises reacting a polyamino alkenyl or alkyl succinimide with a
      compound of the general formula: wherein W is oxygen or sulfur; X is oxygen or sulfur; R4
      is an alkylene group of 2 or 3 carbon atoms optionally substituted by from 1 to 3 alkyl
      groups of 1 or 2 carbon atoms each; and Rs is hydrogen or alkyl of from 1 to 20 carbon
      atoms.</seg></tuv>
    <tuv xml:lang="FR"><seg>Additif pour une huile lubrifiante ou un combustible hydrocarboné,
      pouvant être obtenu par un procédé qui consiste à faire réagir un polyamino-alcényl- ou
      alkyl-succinimide avec un composé de formule générale dans laquelle W représente
      l'oxygène ou le soufre ; X représente l'oxygène ou le soufre ; R4 représente un
      groupe alkylène ayant 2 ou 3 atomes de carbone, facultativement substitué avec 1 à 3
      groupes alkyle ayant chacun 1 ou 2 atomes de carbone ; et Rs représente l'hydrogène
      ou un groupe alkyle ayant 1 à 20 atomes de carbone.</seg></tuv> </tu>

    <tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop
      type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>
    <tuv xml:lang="EN"><seg>An additive as claimed in Claim 1, wherein R4 is alkylene of 2 or
      3 carbon atoms.</seg></tuv>
    <tuv xml:lang="FR"><seg>Additif suivant la revendication 1, dans lequel R4 représente un
      groupe alkylène ayant 2 ou 3 atomes de carbone.</seg></tuv> </tu>

    <tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop
      type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>
    <tuv xml:lang="EN"><seg>An additive as claimed in Claim 1 or 2, wherein Rs is hydrogen or
      alkyl of from 1 to 10 carbon atoms.</seg></tuv>
    <tuv xml:lang="FR"><seg>Additif suivant la revendication 1 ou 2, dans lequel Rs représente
      l'hydrogène ou un groupe alkyle ayant 1 à 10 atomes de carbone.</seg></tuv> </tu>
  
```

Figure 47 : Exemple de données en format TMX (MT CLEF-IP anglais-français)

VIII.1.2 Méthode de création

Nous avons construit les ressources statiques en utilisant 2 méthodes. Grâce à SECTra_w, qui offre un système d'annotation de chaque traduction ou post-édition d'un segment par un niveau de fiabilité (de * à *****) et un score (de 0 à 20), il est possible d'extraire de la mémoire de traductions, associée à un site Web, une sous-MT vérifiant n'importe quel prédicat basé sur les niveaux et les scores. Le résultat de post-édition peut être exporté en format XML et TXT. Nous avons choisi les segments post-édités de haute qualité (par les critères de SECTra).

L'exemple suivant montre une extraction à partir de la partie français-chinois de la MT-demo2 associée à la IMAG-DOC_PAR_JOUR. Le prédicat est tout simplement [Level = 3 & score >= 13], et ses paramètres peuvent être modifiés directement via l'interface graphique.

No	Pseudo Doc	Source	Cible	Stars	Notes
1	DOC16	la salle du haut conseil située au 9ème étage de l'Institut du monde arabe, dans le 5ème.	位于巴黎5区的阿拉伯世界博物馆10楼高级理事会议厅。	3	20
2	DOC16	le 24ème étage de la tour zamansky, l'université pierre et marie curie, sur le campus de jussieu, dans le 5ème.	位于巴黎5区的皮埃尔和玛丽·居里大学加希耶校区的扎曼斯基大楼25楼。	3	20
3	DOC16	le 18ème étage de la bibliothèque François Mitterrand, dans le 13ème.	位于巴黎13区的法国国家图书馆密特朗首19楼。	3	20
4	DOC16	le 6ème étage de l'hôtel industriel de Dominique Perrault, dans le 13ème.	位于巴黎13区的多米尼克·佩罗工业馆7楼。	3	20
5	DOC16	la nuit blanche 2012 permettra aux visiteurs de découvrir la ville lumière d'en haut grâce à 15 belvédères normalement fermés au public.	2012巴黎不眠夜将使参观者可以从15个平时不对外开放的平台发现欣赏巴黎这座光影之城。	3	20
6	DOC16	jk rowling	jk 罗琳	3	20
7	DOC16	en effet, jk rowling, qui a créé notre sorcier à lunettes, pourrait se replonger dans l'univers d'harry potter.	事实上, "创造了"我们那位眼镜魔法师的jk 罗琳, 有可能会写哈利波特的魔法世界里的续集。	3	20
8	DOC16	c'est en 2011 que la saga de nos célèbres sorciers s'est achevée, à l'issue du septième livre intitulé « harry potter et les reliques de la mort » qui a été divisé en deux parties au cinéma.	2011年, 第七本《哈利波特和死亡圣器》分为上下两部电影, 上映结束之后, 我们这著名的哈利波特系列魔法小说完结了。	3	20
9	DOC16	cinq ans après le dernier opus de cette série à succès, jk rowling revient avec une nouvelle œuvre, pour adultes cette fois.	在这成功的系列小说最后一章完结的5年之后, jk 罗琳带着她的新作品回归了。这次是给大人看的小说。	3	20
10	DOC16	et la star des librairies a su entretenir le mystère.	这位图书史上的明星会将这光芒延续下去。	3	20
11	DOC14	la joie	欢乐	3	20
12	DOC16	« une bourgade apparemment idyllique mais qui va faire face aux tourments les plus violents ».	"一个表面看起来非常美好的田园小镇, 但是它将面临最猛烈的动荡。"	3	20

Figure 48 : Extraction d'une "bonne" MT de la MT produite par post-édition "naturelle"

Avec l'autre méthode, nous récupérons d'abord les documents multilingues comme le corpus CLEF-IP 2011, et analysons les documents.

Ensuite, nous extrayons les segments parallèles, et les nettoyons.

Dans un troisième temps, nous les alignons et les mettons en format TXT et TMX.

VIII.1.3 Résultats

Après les 4 premières années d'utilisation (09/2011-09/2015), il y avait environ 80 sites Web accédés par des IMAG. Il y en a maintenant plus de 200. Les sites visités ont plus de 8 langues comme langue source, et plus de 10 langues comme langue cible avec de la post-édition. Il y a plus de 820.000 segments, et environ 45% (370.000+) des segments ont été post-édités par des contributeurs (payés ou non, organisés ou occasionnels). La plupart des segments parallèles sont en anglais-français, en anglais-chinois, et en français-chinois.

Voici quelques données chiffrées.

Dans le Tableau 58, on donne la longueur du texte source des segments post-édités, le balisage XML personnalisé est supprimé, et on compte le nombre de caractères chinois pour le chinois.

Tableau 58 : Segments post-édités dans SECTra_w à partir de 3 langues source

Langue	Phrase	Mots
anglais	350 174	5 252 610
français	55 642	834 621
chinois	270 108	3 241 296

Le Tableau 59 présente les segments post-édités avec la direction de post-édition et la taille de la MT.

Tableau 59 : Segments parallèles obtenus à partir des MT (mêmes remarques)

Paire de langues (L1→L2)	Phrase	Mots L1 (ou caractères)	Mots L2 (ou caractères)	Taille L1	Taille L2
anglais→français	121 074	2 542 731 9 685,9 pages	2 613 351 9862,1 pages	10,1Mo	10,4Mo
anglais→chinois	208 106	4 370 530 16 648,5 pages	<u>6 063 942</u> 151 159,8 pages	19,1Mo	17,6Mo
français→anglais	29 079	627 661 2 326,3 pages	610 098 2 210,1 pages	4Mo	3,9Mo
français→chinois	10 890	228 703 871,2 pages	<u>317 322</u> 793,3 pages	1,5Mo	1,25Mo
chinois→anglais	2 013	<u>58 656</u> 146,6 pages	42 275 161,1 pages	240Ko	263Ko
chinois →français	10 062	<u>291 192</u> 727,9 page	211 185 804,7 page	874Ko	1Mo

VIII.1.3.1 Français-chinois pour l'énergie

Dans le cadre de ma thèse, j'ai travaillé sur la construction d'un système de TA français→chinois pour un client potentiel de L&M. J'ai collecté la ressource française concernant le domaine de l'énergie, et puis nous l'avons post-éditée en chinois à l'aide de GT. La Figure 49 affiche 2 fichiers TXT exportés de SECTra. Le fichier en haut est celui des segments source, le fichier du bas est celui des segments post-édités.

Répondre à la demande croissante d'énergie tout en parant aux risques climatiques et à la raréfaction des ressources : c'est le défi lancé aux énergéticiens. C'est pourquoi le développement durable est au cœur de notre stratégie. Notre politique traduit la volonté du Groupe de « changer l'énergie ensemble » en apportant des solutions réalistes. La lutte contre le changement climatique commence par la maîtrise de nos impacts environnementaux.

Nous nous sommes fixé 9 engagements pour répondre à 3 enjeux prioritaires

La lutte contre le changement climatique et la préservation de la biodiversité :

Un enjeu environnemental

rester, en tant que Groupe, le moins émetteur de CO des grands énergéticiens européens,

1. adapter notre parc de production et nos offres au changement climatique,
2. réduire notre impact environnemental, notamment sur la biodiversité.
3. Faciliter l'accès à l'énergie et développer des liens de proximité avec les territoires :

Un enjeu sociétal

favoriser l'accès à l'énergie et l'éco-efficacité énergétique,

4. développer dans la durée la proximité avec les territoires où nous opérons,
5. contribuer à l'effort éducatif sur les questions liées à l'énergie.
6. Contribuer au débat sur le développement durable par le dialogue, l'information et la communication :

Un enjeu de gouvernance

poursuivre le développement des politiques et le partage des valeurs au sein du Groupe, en relation avec les parties prenantes,

7. communiquer et rendre compte des activités et résultats du Groupe en matière de développement durable,
8. participer au débat sur le développement durable au niveau national et international.

既能满足日益增长的能源需求，又要应对气候变化风险和资源枯竭问题，这就是能源企业面临的挑战。为此，可持续发展成为法国电力集团的战略核心。“共同改变能源结构”体现了法国电力集团的意愿，为此，我们将提供切实可行的技术方案。

应对气候变化首先应从控制环境影响做起。

我们制定了九项承诺，优先应对以下三方面的挑战：

应对气候变化，保护生物多样性：

环境的挑战

法国电力集团将继续保持欧洲大型能源公司二氧化碳排放最少企业的地位；

1. 电力生产设施和服务项目适应气候变化的要求；
2. 减少对环境，特别是对生物多样性的影响。
3. 为普及能源使用提供方便，与地方合作，推行就近供电政策：

社会的挑战

为普及能源使用提供便利，发挥能源项目的生态效益；

4. 与地方合作，在集团经营区域内长期实施就近供电政策；
5. 支持能源教育的各种活动。
6. 支持关于可持续发展的辩论，可采用对话、信息与交流等各种形式：

管理方面的挑战

与有关方面合作，在集团内部继续制定相关政策，共享企业价值观；

7. 交流和通报集团在可持续发展方面开展的工作和取得的成果；
8. 参与法国和国际有关可持续发展的讨论。

Figure 49 : Segments post-édités pour la ressource énergie

Il y a environ 10K segments français, que nous divisons en 500 pages Web. Après la post-édition, nous avons récupéré 9000 segments français→chinois (Tableau 61).

Tableau 61 : Statistique des données pour la ressource énergie

	Nb de segments	Longueur moyenne de segment	Nb de mot (ou caractère chinois)	Page_std
français	9000	41	368K	263
chinois	9000	36	324K	810

VIII.2 Contribution sous forme de systèmes de TA

Il s'agit soit de systèmes téléchargeables, soit de système utilisables comme des services Web.

VIII.2.1 Systèmes de TA téléchargeables

Nous les avons mis à l'url <http://...> Il s'agit des systèmes créés en français-chinois et des systèmes créés à partir de la collection de brevets CLEF-IP 2011.

Tableau 62 : Systèmes de TA téléchargeables

Nom du système de TA	Langues	Ressources utilisées	Commentaire
Moses-L&M	fr-zh	MT de 9000 segments créée par L&M	Qualité assez basse BLEU TER Tpe estimé
Moses-LIG	fr-zh	MT de 80K segments PE + 100K MultiUN	
Moses-CLEF-IP	fr-de fr-en de-fr		Très bien sur les parties de brevets BLEU :

VIII.2.2 Systèmes de TA utilisables comme des services Web

Certains des systèmes précédents sont des services Web utilisables directement via l'API de TRADOH, ou la plate-forme JIANDAN-EVAL. Certains sites Web boursiers et économiques sont aussi utilisables via une iMAG, ce qui permet l'amélioration par AI.

Tableau 63 : Systèmes de TA utilisables comme des services Web

Nom du système de TA	Langues	Ressources utilisées	Commentaire
Moses-L&M	fr-zh	MT de 9000 segments créée par L&M	Qualité assez basse
Moses-LIG	fr-zh	MT de 80K segments PE + 100K MultiUN	
Moses-CLEF-IP	fr-de fr-en de-fr		Très bien sur les parties de brevets

VIII.3 Passerelles iMAG vers des sites Web statiques ou dynamiques

Une ressource de ce type est simplement une iMAG-S dédiée à un site Web S. La MT associée peut être propre ou partagée, cela n'a pas d'importance pour les visiteurs. Nous distinguons les sites accédés S selon qu'ils sont statiques ou dynamiques.

VIII.3.1 Passerelles iMAG pour des sites statiques

Il s'agit essentiellement de sites contenant des documents mis sous forme HTML, et pour lesquels on vise à obtenir une assez bonne qualité de traduction.

Tableau 64 : Passerelles iMAG pour des sites statiques

iMAG	Documents accédés	Commentaires
BemBook	Le BemBook (de livre 700 p.) Site : http://www.bembook.ibpsa.us/	Post-édité vers le français de façon occasionnelle Expérimentation organisée pour PE du chapitre 21 en 7~8 langues.
Corpus par jour	« Voyage au centre de la terre » de Jules Verne	Post-édition de ce roman du français vers le chinois
Manchanda	« IITB : Monastery, Sanctuary, Laboratory »	Préparé pour la PE vers français et hindi
Macau	Support de cours	Pour les étudiants étrangers
Powers	« The book of me » de Powers	Expérimentation de L. Besacier
CLEF-IP	Corpus de brevets « CLEF-IP 2011 »	PE inutile pour les langues de la collection, mais possible vers le chinois

VIII.3.2 Passerelles iMAG pour des sites dynamiques

Il s'agit ici de sites dont certaines pages changent régulièrement, par exemple celles des nouvelles dans des sites Web d'organismes ou de firmes.

Il y a aussi des sites de journaux, utilisés pour la construction de ressources (corpus parallèles, dictionnaires bilingues) utilisables pour construire ensuite des systèmes de TA (empiriques ou experts ou les deux) vers des langues peu dotées. C'est le cas des iMAG créées pour la Nation de Djibouti (thèse sur le somali) et pour La Gazette des Comores et Al-Watwan (thèse sur le comorien et plus particulièrement le shingazidja) visant à créer un système de lecture active de textes administratifs et juridiques écrits en français (appris à l'école) à l'intention des Comoriens qui parlent presque toujours le comorien en fait.

Tableau 65 : iMAG pour des sites dynamiques

iMAG	sites accédés	Commentaires
LIG-LAB	Site Web du LIG	Post-édité vers le chinois (avec expérimentation) et l'anglais. Un peu vers l'allemand.
lanationdedjibouti	Site Web du journal La Nation de Djibouti. Environ 500 nouveaux segments par jour, 3 dernières années en ligne.	Construction par PE (à partir des résultats de GT) d'une bonne MT français-somali. But : créer un système de TA Moses meilleure que GT pour d'accéder à ce journal.
lametro	Site Web de La Métro à Grenoble	Site de La Métro. Gros travail (fr→zh/en) en 2010.
alwatwan	Sites Web des 2 journaux Al-Watwan et La Gazette des Comores	Construction par PE à partir des traductions en wahili de GT d'une bonne MT français-comorien

VIII.3.3 Structure d'une contribution « dynamique » par iMAG

Il s'agit ici d'iMAG dédiées à des sites Web dynamiques. Les utilisateurs post-éditent le contexte sur SECTRA_w/iMAG. Les segments post-édités sont sauvegardés dans les MT.

Il y a 2 types de sites Web dynamiques. (1) Un site Web (par exemple le site Web *lamétre*⁹¹) met à jour son contexte en temps réel. Pendant que nous accédons son iMAG dédiée dans une langue cible, nous choisissons les phrases intéressantes, et les post-éditons. Un segment peut être post-édité (vers une ou plusieurs langues) par un ou plusieurs post-éditeurs, donc chaque segment source peut avoir plusieurs versions de segments post-édités. Quand nous avons un certain nombre de segments post-édités, nous les exportons pour construire la MT. (2) Nous collectons des ressources spécifiques pour construire une MT adaptée au sous-langage. Tout d'abord, nous récupérons les textes, et les nettoyons. Ensuite nous les mettons en format HTML (Figure 50). Enfin, nous les ajoutons dans SECTRA_w/iMAG.

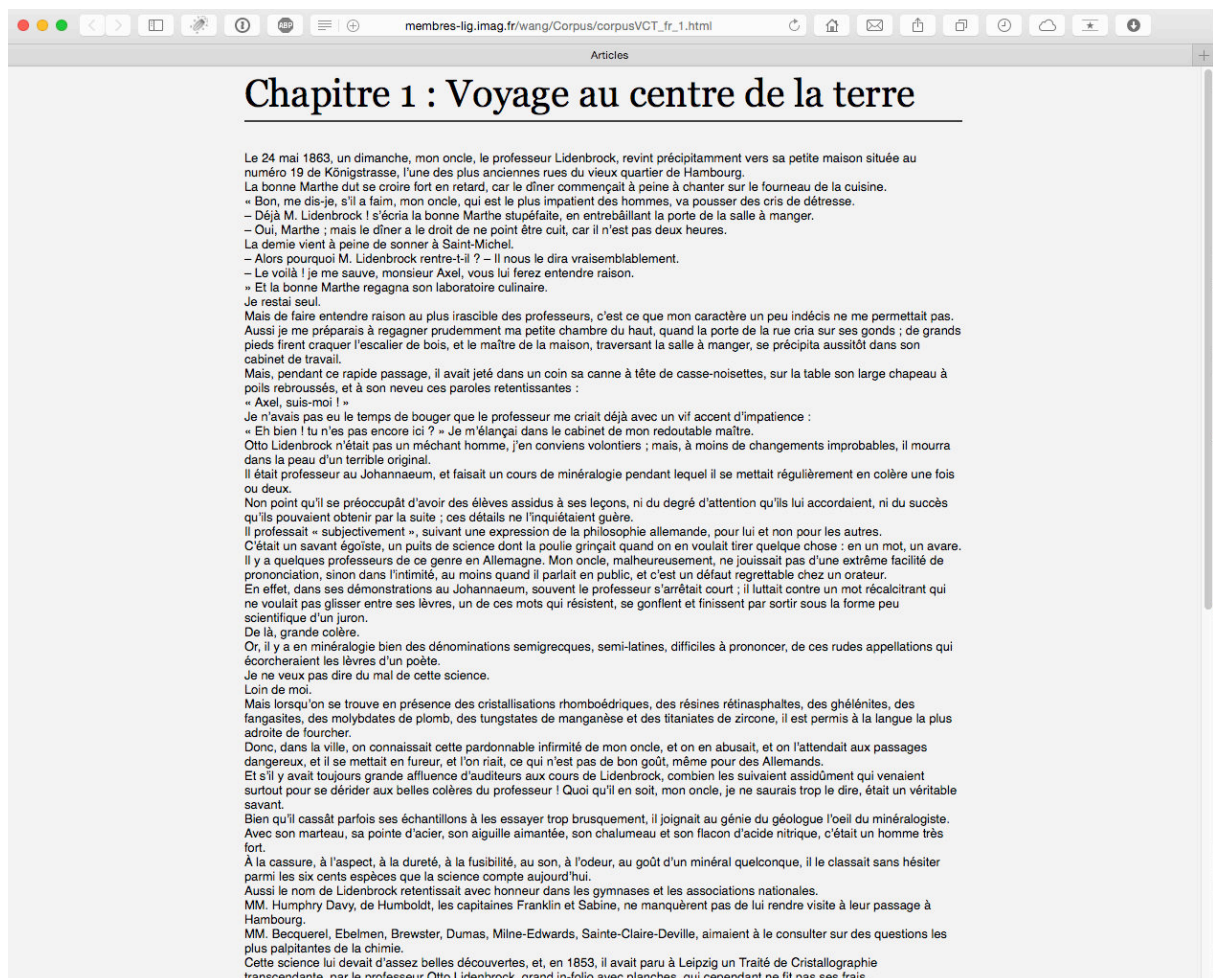


Figure 50 : Exemple de contribution au format HTML (Chapitre 1 : Voyage au centre de la Terre)

⁹¹ <http://www.lametro.fr>

VIII.3.4 Remarques sur la création de certains des sites « contribués »

Certaines de ces contributions ont demandé des travaux de préparation spécifiques.

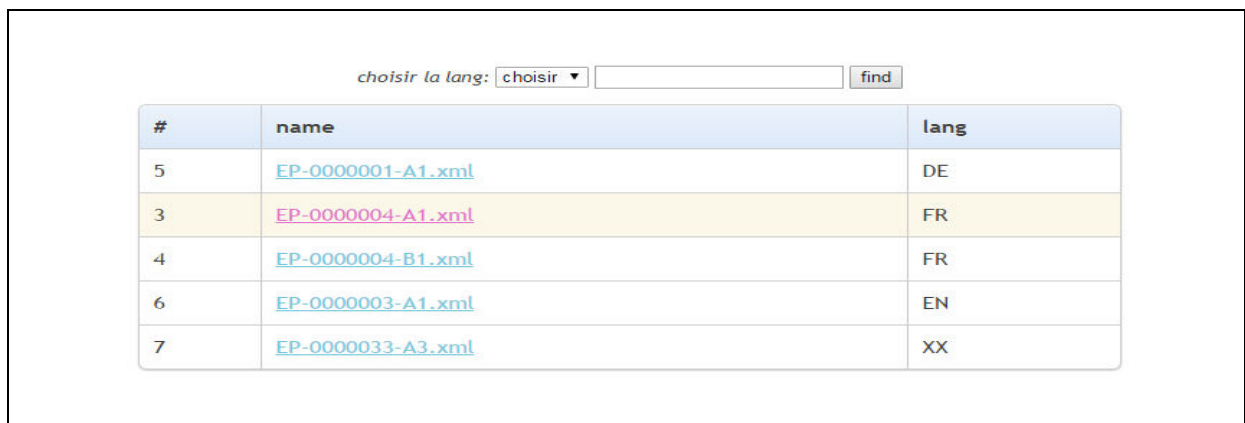
VIII.3.4.1 La Métro

Le site Web de La Métro de Grenoble contient 2500 pages Web, soit environ 30000 segments. Plus de la moitié ont été post-édités en chinois pour l'Exposition universelle de Shanghai en 2010. Nous avons créé un système français-chinois spécifique en choisissant une partie "assez bonne" de cette MT du 26/06/2015 au 30/06/2015 par la sélection:

```
TM-lametro-extract-good = TM_select (lametro, [level=3 & score >=13 |  
level=4 & score >=13 | level=5 & score >=11.5], [langage_pair=fr-zh |  
langage_pair=fr-en | langage_pair=zh-fr], [Begin_date : 20150626,  
End_date : 20150630]).
```

VIII.3.4.2 MUMIA-CLEF

Avec Huanan Sun (M1-TER), nous avons construit 3 sites Web monolingues pour le corpus CLEF-IP. Chaque site Web contient environ 500K pages Web. Chaque page Web représente un fichier XML de CLEF-IP (Figure 51). J'ai créé 3 iMAG dédiées pour accéder à ces sites en utilisant les systèmes que nous avons créés pour les couples fr→en, de→fr, fr→de, et les services de TA disponibles pour les autres couples.

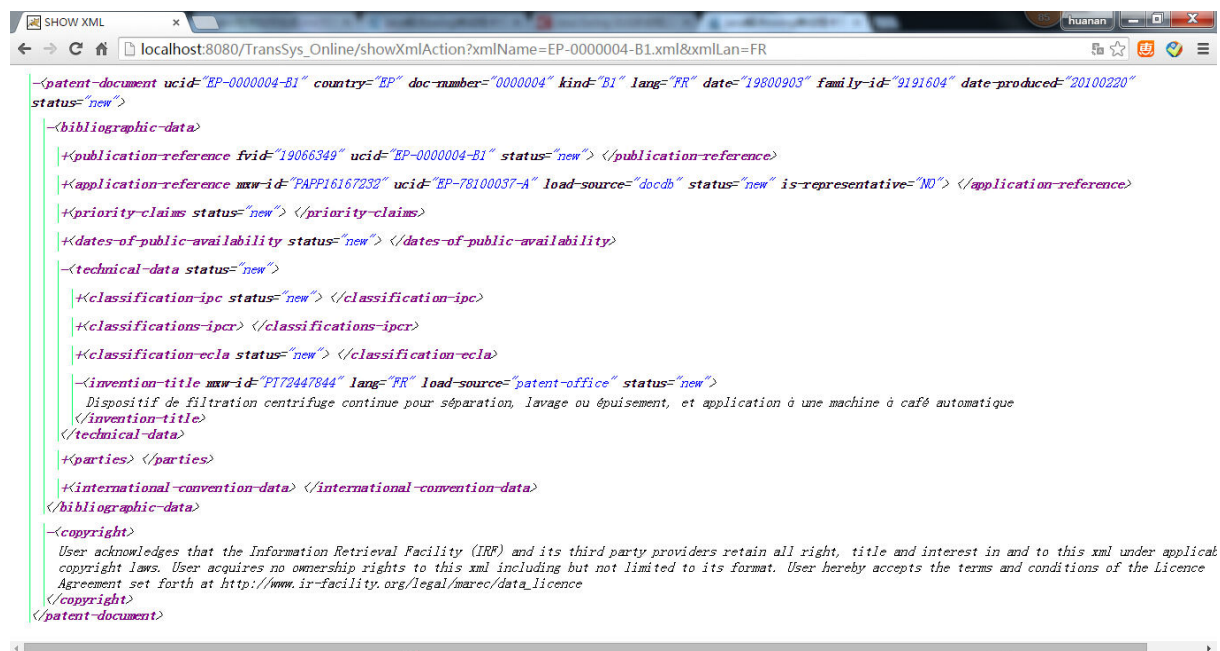


The screenshot shows a web interface for selecting a language. At the top, there is a label "choisir la lang:" followed by a dropdown menu currently set to "choisir", an empty text input field, and a "find" button. Below this is a table with three columns: "#", "name", and "lang". The table contains five rows of data, each representing an XML file.

#	name	lang
5	EP-0000001-A1.xml	DE
3	EP-0000004-A1.xml	FR
4	EP-0000004-B1.xml	FR
6	EP-0000003-A1.xml	EN
7	EP-0000033-A3.xml	XX

Figure 51 : Capture d'écran du site Web monolingue de CLEF-IP

La Figure 52 montre une capture d'écran de l'iMAG CLEF-EN utilisée pour accéder dans d'autres langues aux brevets écrits initialement en anglais. Les segments entre parenthèses vertes ont été trouvés dans la MT CLEF-IP, et les segments entre parenthèses rouges ont été traduits par le système de TA anglais→français.



The screenshot shows a web browser window with the address bar displaying "localhost:8080/TransSys_Online/showXmlAction?xmlName=EP-0000004-B1.xml&xmlLan=FR". The main content area displays XML data for a patent document. The XML is color-coded with green tags and red text for values. The document includes bibliographic data, technical data, and a copyright notice.

```
<patent-document ucid="EP-0000004-B1" country="EP" doc-number="0000004" kind="B1" lang="FR" date="19800903" family-id="9191604" date-produced="20100220" status="new">
  <bibliographic-data>
    +<publication-reference fvid="19066349" ucid="EP-0000004-B1" status="new"> </publication-reference>
    +<application-reference mxw-id="PAPP16167232" ucid="EP-78100037-A" load-source="docdb" status="new" is-representative="NO"> </application-reference>
    +<priority-claims status="new"> </priority-claims>
    +<dates-of-public-availability status="new"> </dates-of-public-availability>
  </bibliographic-data>
  <technical-data status="new">
    +<classification-ipc status="new"> </classification-ipc>
    +<classifications-ipc> </classifications-ipc>
    +<classification-ecla status="new"> </classification-ecla>
    <invention-title mxw-id="PT72447844" lang="FR" load-source="patent-office" status="new">
      Dispositif de filtration centrifuge continue pour séparation, lavage ou époussement, et application à une machine à café automatique
    </invention-title>
  </technical-data>
  +<parties> </parties>
  +<international-convention-data> </international-convention-data>
</bibliographic-data>
<copyright>
  User acknowledges that the Information Retrieval Facility (IRF) and its third party providers retain all right, title and interest in and to this xml under applicat
  copyright laws. User acquires no ownership rights to this xml including but not limited to its format. User hereby accepts the terms and conditions of the Licence
  Agreement set forth at http://www.ir-facility.org/legal/marec/data_licence
</copyright>
</patent-document>
```

Figure 52 : Capture d'écran de l'iMAG dédiée CLEF-IP

Chapitre IX Vers une plate-forme de construction, déploiement et évaluation de systèmes de TA: JianDan-eval

Résumé

Cette partie présente la construction en cours de la plate-forme JIANDAN-EVAL. On présente d'abord le cahier des charges, la spécification externe et l'architecture logicielle. On décrit ensuite brièvement l'implémentation : outils et environnements de base, composants et bibliothèques intégrés, et composants développés. On donne enfin un exemple d'utilisation, la construction d'un système de TA fondé sur MOSES par un utilisateur de JIANDAN-EVAL. Ce système a déjà servi à un étudiant de M2R pour construire et évaluer des systèmes MOSES.

IX.1 Cahier des charges, architecture, et spécifications externes

IX.1.1 Cahier des charges

IX.1.1.1 Problèmes actuels

Les systèmes de TA sont beaucoup utilisés dans le domaine scientifique et dans le domaine industriel. Malheureusement, il arrive fréquemment qu'ils ne soient pas adaptés à la situation traductionnelle précise. Souvent, un traducteur professionnel voudrait construire son système de traduction automatique (TA) en utilisant sa propre mémoire de traductions (MT). Comme nous l'avons vu dans le chapitre précédent, il existe plusieurs boîtes à outils pour fabriquer un système personnalisé, mais l'utilisateur doit avoir de bonnes connaissances en informatique pour manipuler les commandes. C'est un vrai défi pour les non-informaticiens. C'est pourquoi nous avons voulu proposer une plate-forme en utilisation libre, JIANDAN-EVAL.

IX.1.1.2 Objectifs du projet

JIANDAN-EVAL est définie comme une plate-forme de création, d'évaluation, d'utilisation et d'amélioration de systèmes de TA. Notre objectif est d'aider les non-informaticiens à construire leurs systèmes de TA sans connaissances sur la boîte à outils MOSES, à mener des campagnes d'évaluation ponctuelles ou permanentes, et, de façon plus générale, à créer une infrastructure pour l'évaluation humaine (objective par PE, et subjective).

IX.1.1.3 Description des modules principaux

La première version de JIANDAN-EVAL contient 4 modules réalisant les 4 fonctionnalités principales : construction, évaluation, déploiement, et amélioration de systèmes de TA.

IX.1.1.3.1 Construction de systèmes de TA

JIANDAN-EVAL est construit sur la boîte à outils MOSES, et intègre aussi des outils pour aider à construire des systèmes de TA. JIANDAN-EVAL contient tous les services nécessaires, du prétraitement jusqu'à l'entraînement de systèmes de TA.

La conception de la fonction « construction de systèmes de TA » a été inspirée par KANTANMT⁹². Comme dans KANTANMT, cette partie est limitée aux système de type MOSES.

L'utilisateur télécharge les données d'entraînement via l'interface de JIANDAN-EVAL, puis choisit les paramètres pour la construction d'un système de TA. Il peut choisir différents types de modèles de traduction. JIANDAN-EVAL permet de télécharger les données pour la construction du modèle de langue cible, l'optimisation des poids de traduction, et l'évaluation du système.

⁹² KantanMT : <http://www.kantanmt.com>

Pour les données d'entraînement, il ne s'agit pas seulement de phrases parallèles, mais aussi de documents multilingues parallèles alignables (comme les brevets), et JIANDAN-EVAL intègre des outils pour aider à les aligner. Pour améliorer la qualité du système de TA en construction, nous avons aussi ajouté la plupart des outils de filtrage existants.

Comme il utilise la boîte à outil MOSES, JIANDAN-EVAL dispose du segmenteur de mots intégré à MOSES, qui est bon mais limité à certaines langues, comme l'anglais, le français, l'allemand, etc. Pour pouvoir créer des systèmes de TA français-chinois, nous avons intégré à JIANDAN-EVAL d'autres segmenteurs, comme STANFORD SEGMENTOR, NEON SEGMENTOR, etc.

JIANDAN-EVAL a un moniteur permettant de surveiller l'avancement de l'entraînement. Ce moniteur affiche toute l'information pour chaque étape de l'entraînement, comme le nom de l'étape, les temps d'exécution, etc.

IX.1.1.3.2 Évaluation des systèmes de TA

JIANDAN-EVAL est conçu pour supporter des campagnes d'évaluation ponctuelles (comme celles de IWST) ou permanentes.

Paramétrabilité des systèmes TA et des données. L'utilisateur peut évaluer un ou plusieurs systèmes de TA. Il peut s'agir non seulement de différentes versions de systèmes de TA entraînés par l'utilisateur, mais aussi de systèmes comme GOOGLE TRANSLATE, SYSTRAN et BING. L'utilisateur peut télécharger les données de test ou choisir une petite partie des données d'entraînement (par exemple, 1000 bisegments) pour évaluer les systèmes.

Paramétrabilité des méthodes « objectives ». JIANDAN-EVAL supporte l'évaluation « objective » (non liée à des jugements humains) comme BLEU, NIST, TER, notre distance de post-édition, etc. L'utilisateur utilise les données de test pour évaluer son système de TA. Les scores (comme BLEU) sont calculés par des scripts intégrés à JIANDAN-EVAL.

Paramétrabilité des méthodes « subjectives » (liées à des jugements humains). L'organisateur de l'évaluation définit les mesures comme fluidité, adéquation, valeur d'usage, fiabilité, fidélité, grammaticalité, qualité terminologique, etc., et le système associe un score à chaque segment traduit par TA ou post-édité, pour chaque mesure.

IX.1.1.3.3 Déploiement de systèmes de TA

JIANDAN-EVAL permet de déployer uniquement des systèmes de type MOSES. D'une part, les systèmes de TA entraînés sous JIANDAN-EVAL sont déployés directement pour la traduction et l'évaluation. D'autre part, si un utilisateur a construit son système sur une machine locale, il peut le téléverser sur JIANDAN-EVAL, puis l'utiliser quelques heures plus tard.

IX.1.1.3.4 Amélioration du système de TA

JIANDAN-EVAL permet d'améliorer les systèmes de TA de type MOSES par post-édition contributive. L'utilisateur utilise son système pour produire les « prétraductions », puis il les post-édite. Ensuite, il peut filtrer et choisir les segments post-édités à partir de sa MT (cette fonction est réalisée via la communication avec SPECTRA/IMAG), et il les ajoute aux données d'entraînement. Enfin, l'utilisateur peut relancer la procédure d'entraînement de TA pour améliorer la qualité de son système. Quand on entraîne une nouvelle version de TA, la version courante reste utilisable.

IX.1.2 Spécifications externes

IX.1.2.1 Introduction

Cette partie a pour but de spécifier l'aspect externe de la plate-forme JIANDAN-EVAL, ce qui est utilisable par les différents types d'utilisateur au travers de l'utilisation des fonctionnalités auxquelles ils ont accès. Elle contient quatre sections :

- les rôles d'utilisateur,
- les scénarios d'utilisation associés,
- les fichiers d'entrée et les fichiers de sortie,
- les interfaces du logiciel et les modèles de tâches.

IX.1.2.2 Acteurs du système

Utilisateur. L'acteur « utilisateur » représente tout utilisateur physique souhaitant se servir de JIANDAN-EVAL. Il devra s'identifier dans le système, afin d'avoir par la suite les droits et privilèges correspondant à son profil. Un utilisateur normal a le droit de fabriquer un système de TA avec ses données, mais la taille des données d'entraînement est limitée à 50Mo. Si l'utilisateur veut utiliser de nouvelles données monolingues (différentes des données d'entraînement cible) pour fabriquer le modèle de langue, la taille est limitée à 50Mo. Lors de chaque session, un utilisateur peut fabriquer de 1 à 3 systèmes de TA. Il peut toujours choisir les données pour évaluer son système par BLEU/NIST, TER, etc.

Organisateur du projet. L'acteur « organisateur du projet » est un utilisateur qui peut proposer un projet d'évaluation de systèmes de TA. Il peut distribuer les rôles d'évaluateur et développer les mesures d'évaluation.

Évaluateur. L'acteur « évaluateur » peut participer au projet et évaluer la qualité de TA.

Administrateur. L'acteur « administrateur » gère les utilisateurs, et leurs privilèges.

Technicien. Un technicien surveille le processus de la création des systèmes de TA, et résout les problèmes, en particulier ceux liés au service Web. Il peut aussi utiliser un sous-système de communication avec les utilisateurs.

IX.1.2.3 Scénarios

IX.1.2.3.1 Création de systèmes de TA

La première tâche principale est le support de la création de systèmes de TA en ligne. JIANDAN-EVAL intègre les outils de MOSES, et permet d'entraîner les modèles de traduction en utilisant les données fournies par l'utilisateur ou des données publiques (EUROPARL, MULTIUN, EUBOOKSHOP, etc.).

a. Scénario 1 : Soumission des données d'entraînement

1. Un utilisateur crée un système de TA du français vers l'anglais en utilisant une MT qui est sauvegardée au format TMX. Il entre un nom pour ce système.
2. JIANDAN-EVAL demande que le fichier de données soit compressé au format .zip ou .tar.gz, et que la taille du fichier compressé soit inférieure ou égale à 50Mo (pour plus, il faudra acquérir des droits supplémentaires).
3. Cet utilisateur soumet ses données d'entraînement à JIANDAN-EVAL. Il choisit la paire de langues français-anglais et soumet le fichier zip dans l'interface de JIANDAN-EVAL.
4. Il clique sur le bouton « Upload », et une fenêtre de sélection apparaît dans la page Web. Il choisit le fichier sur son ordinateur local. Il fait un double-clic sur le fichier cible (ou bien il clique sur le fichier cible, puis sur le bouton « Confirm ») pour soumettre son fichier zip.
5. La progression du téléversement est représentée par une barre de progression.
6. Après le téléversement, le bouton « Next » redevient valide pour cet utilisateur. Quand il clique sur « Next », il entre dans l'interface d'entraînement.

b. Scénario 2 : Configuration d'un système de TA

1. D'abord, l'utilisateur choisit son système préféré, par exemple MOSES, sur l'interface de configuration des systèmes de TA. Les paramètres modifiables de MOSES

apparaissent. L'utilisateur choisit MGIZA++ pour l'alignement des mots, et il coche la case « Parallel » pour créer les modèles dans les deux sens (français↔anglais).

2. Ensuite, il choisit IRSTLM pour l'apprentissage des modèles de langue, et PHRASE-BASED MODEL pour le modèle de MOSES. Il garde les valeurs par défaut des autres paramètres.
3. Enfin, pour faire le « Tuning », JIANDAN-EVAL demande à l'utilisateur de téléverser ses données. Ce dernier choisit les données de Tuning déjà préparées avec l'aide de JIANDAN-EVAL, ou ignore cette étape, et les données de Tuning sont extraites à partir des données d'entraînement (premières lignes des données d'entraînement). Pour évaluer la qualité du système de TA obtenu, l'utilisateur téléverse les données de test sur le serveur, et coche les cases BLEU, NIST, et TER.
4. JIANDAN-EVAL envoie un mail de confirmation de l'entraînement de MOSES à l'utilisateur.

c. *Scénario 3 : Configuration simple d'un système de TA*

1. Après le téléchargement des données (au format .txt), l'utilisateur choisit la configuration simple de Moses. JIANDAN-EVAL demande de soumettre les données de test.
2. L'utilisateur soumet les données de test, et clique sur le bouton « Confirm ».
3. Le serveur SMTP envoie un mail de confirmation de l'entraînement de MOSES à l'utilisateur.

d. *Scénario 4 : Surveillance de la procédure d'entraînement de système de TA*

1. L'utilisateur clique sur le bouton « Monitor » pour entrer dans l'interface de surveillance de l'état de TA.
2. L'utilisateur consulte l'étape actuelle de l'entraînement de système de TA. Il voit que la tâche d'entraînement est dans la deuxième étape « Clean data ». Il clique sur le bouton vert pour exporter un fichier LOG, qui contient l'information plus détaillée sur le déroulement de l'entraînement de TA. Il trouve que les données ne sont pas correctes, et il veut arrêter cette tâche.
3. Il clique alors sur le bouton « Cancel job », puis sur « Yes » dans une fenêtre de confirmation.
4. Sa tâche s'arrête, et un mail d'annulation de tâche lui est envoyé par le serveur SMTP.
5. La page Web retourne à la page d'accueil.

e. *Scénario 5 : réussite de la création d'un système de TA*

1. Après la configuration de la création de son système de TA sur JIANDAN-EVAL, l'utilisateur se déconnecte du site JIANDAN-EVAL.
2. Il reçoit un mail d'annonce de fin d'entraînement de TA 10 heures plus tard.
3. Il se connecte alors sur le site JianDan-eval, et voit que son système de TA a été mis en place dans sa page de TA personnelle. Il commence à l'utiliser.

IX.1.2.3.2 Utilisation d'un système de TA

L'utilisateur peut utiliser un ou plusieurs systèmes de TA pour la traduction ou pour l'évaluation de la qualité de traduction. Les systèmes sont de 2 sortes : les systèmes créés par l'utilisateur, et les systèmes publics. Il s'agit donc de systèmes de TA déjà créés et disponibles dans JIANDAN-EVAL, ou bien de systèmes de TA comme GOOGLE TRANSLATE, BING, SYSTRAN, etc.

JianDan-eval fournit 2 façons d'utiliser les systèmes de TA personnels : (1) JianDan-eval contient une interface pour faire la traduction en ligne, comme GOOGLE TRANSLATE, et (2) l'API de JIANDAN-EVAL peut être intégrée dans le système de l'utilisateur.

a. *Scénario 1 : Traduction en ligne*

1. L'utilisateur clique sur le lien « My MT » pour entrer dans l'interface de traduction de JIANDAN-EVAL.
2. Il choisit son système de TA « Moses-fr-en » (c'est le nom qu'il a donné à son système) pour la traduction. À gauche, il y a un champ de saisie, qui permet d'entrer les phrases source. L'utilisateur colle un texte, qui contient 5 phrases, dans le champ de saisie. Il clique sur le bouton « Translate », et la traduction apparaît dans la partie droite de l'interface.
3. L'utilisateur clique sur la case à cocher « PDF », puis sur le bouton « export ». Un fichier PDF, qui contient les résultats de traduction du système de TA, est téléchargé sur la machine locale de l'utilisateur.

b. *Scénario 2 : Traduction de document*

1. L'utilisateur clique sur le lien « My MT » pour entrer dans l'interface de traduction de JIANDAN-EVAL.
2. Il clique sur le bouton « Translate a document », et choisit un document .txt (2000 phrases) sur la machine locale. Son fichier est soumis à JIANDAN-EVAL, et un mail est envoyé dans sa boîte aux lettres.
3. Après 5 minutes, il reçoit un mail de confirmation, disant que son document a été traduit. Ce mail contient un lien de téléchargement de fichier.
4. Il retourne à l'interface de traduction de JIANDAN-EVAL, et clique sur l'icône du fichier .txt de traduction. Le fichier de traduction est téléchargé dans la machine locale.

c. *Scénario 3 : Intégration d'API dans le système de l'utilisateur*

1. L'utilisateur a un projet de service Web développé sous ECLIPSE. Il importe le fichier JIANDAN-EVAL.JAR dans son projet.
2. Il ajoute un bandeau dans le code html de sa page d'accueil ; ce bandeau contient les systèmes de TA valables pour cet utilisateur dans JIANDAN-EVAL. Il choisit la langue source/cible pour son projet.
3. Il déploie son projet dans Tomcat, et un bandeau caché est intégré dans sa page d'accueil. Si on place la souris en haut de cette page, le bandeau de JIANDAN-EVAL apparaît.
4. L'utilisateur choisit la langue « English », et sa page d'accueil est traduite de français en anglais.
5. JIANDAN-EVAL est aussi un système d'évaluation, qui permet à l'utilisateur d'évaluer en ligne la qualité de son système.

IX.1.2.3.3 *Évaluation des systèmes de TA*

Après l'entraînement d'un système de TA, l'utilisateur peut évaluer son système avec les données de test, et il peut calculer des scores comme BLEU, NIST, TER, etc.

a. *Scénario 1 : Calcul de BLEU*

1. L'utilisateur clique sur le lien « Evaluation », et entre dans l'interface d'évaluation.
2. Les systèmes de TA qu'il a créés sont présentés dans un champ de sélection. Par exemple, on y trouve les 3 systèmes de TA français-anglais que l'utilisateur a créés.
3. L'utilisateur choisit « All » TA, puis la langue source « fr » et la langue cible « en ». Chacun des 3 systèmes de TA est alors évalué par la mesure BLEU.
4. L'utilisateur clique sur le bouton « Upload » pour télécharger de nouvelles données de test (par exemple, 50 bi-phrases, les phrases source et leurs traductions de référence), et coche les cases « BLEU » et « NIST ». Il clique sur le bouton « Start ». Le calcul des scores des 3 systèmes de TA commence.

5. Après quelques minutes, les scores ont été calculés sur les données de test, chaque phrase source a 3 résultats, et chaque résultat a 2 scores.
- b. *Scénario 2 : évaluation humaine dans un cadre de « Project »*
1. L'utilisateur a créé un projet dans JIANDAN-EVAL. Le projet contient 2 systèmes de TA du français vers le chinois, et 3 évaluateurs.
 2. L'utilisateur propose un article en français pour évaluer la qualité de la TA. D'abord, l'article est segmenté en phrases. Ensuite, les phrases sont traduites de français en chinois par les 2 systèmes de TA.
 3. Les 3 évaluateurs évaluent chaque traduction sur l'interface d'évaluation. Ils donnent un score à chaque traduction, et proposent une meilleure traduction pour chaque traduction. En même temps, ils ajoutent des commentaires aux traductions.

IX.1.2.4 Données d'entrée et données de sortie

Les données d'entrée de JIANDAN-EVAL sont de deux types :

1. des fichiers en format txt ou tmx, qui contiennent les données d'entraînement, les données de test, les données d'évaluation, etc. Ils sont utilisés pour entraîner les systèmes de TA, et les évaluer.
2. un système de TA (fabriqué avec Moses), destiné à être déployé sur JIANDAN-EVAL.

Les données de sortie de JIANDAN-EVAL contiennent les traductions obtenues par TA, les résultats d'évaluation, l'information attachée, ou le système entraîné par l'utilisateur.

IX.1.3 Architecture logicielle

L'architecture de JIANDAN-EVAL, présentée ci-dessus, est composée de 3 parties : un serveur d'applications, un serveur de TA, et un serveur d'entraînement.

IX.1.3.1 Serveur d'applications

Le serveur d'applications

- assure le rôle de passerelle vers la base des données de tous les utilisateurs,
- contient une base de données,
- fournit une interface utilisateur (client léger) pour accéder aux pages Web,
- gère les appels à la TA.

Le serveur d'applications assure un double rôle dans l'architecture de JIANDAN-EVAL.

D'abord, il héberge l'application Web d'édition des disponibilités (client léger). Afin d'éviter un nombre important de problèmes de sécurité, les communications entre le serveur et l'utilisateur pourront être chiffrées.

Le serveur d'applications fournit également les services permettant aux utilisateurs d'accéder aux informations de la base de données. L'avantage est qu'ainsi les interactions (consultation, modification ou suppression) ne se font que depuis un seul et unique point d'accès, ce qui facilite grandement la maintenance. De cette manière, nous ne dépendons pas du choix de la base de données. De plus, cela permet d'assurer un contrôle d'accès plus strict aux informations applicatives, et simplifie notamment les mises à jour du schéma de la base.

IX.1.3.2 Serveur de TA

Il peut y avoir un ou plusieurs serveurs de TA. Un serveur de TA contient les décodeurs et des modèles de traduction, et fournit le service de traduction. Il n'est pas visible par l'utilisateur, mais l'interface de service de traduction est accessible par l'API de TRADOH.

Les décodeurs peuvent être celui de Moses, ou celui de Joshua, ou celui de XMU-AI.

IX.1.3.3 Serveur d'entraînement

Le serveur d'entraînement est utilisé pour la création de systèmes de TA de type MOSES. Il permet de créer les tables d'un système de TA. Tout d'abord, il reçoit les données d'entraînement à partir du serveur d'applications. Ensuite, il prépare les données, et démarre une procédure d'entraînement. Enfin, il renvoie le nouveau modèle au serveur de TA.

IX.1.3.4 Interconnexion avec les systèmes extérieurs

Authentification avec Facebook/Google+/Twitter. Pour l'authentification des utilisateurs, il sera intéressant de proposer des alternatives de connexion rapide via les réseaux sociaux : FACEBOOK, GOOGLE PLUS ou TWITTER.

L'authentification sociale permettra de gérer à la fois l'authentification sur les réseaux, et aussi de récupérer facilement des données depuis ces derniers.

Envoi de courriels avec SMTP. Les envois de courriels se feront par l'intermédiaire du serveur SMTP de chaque machine.

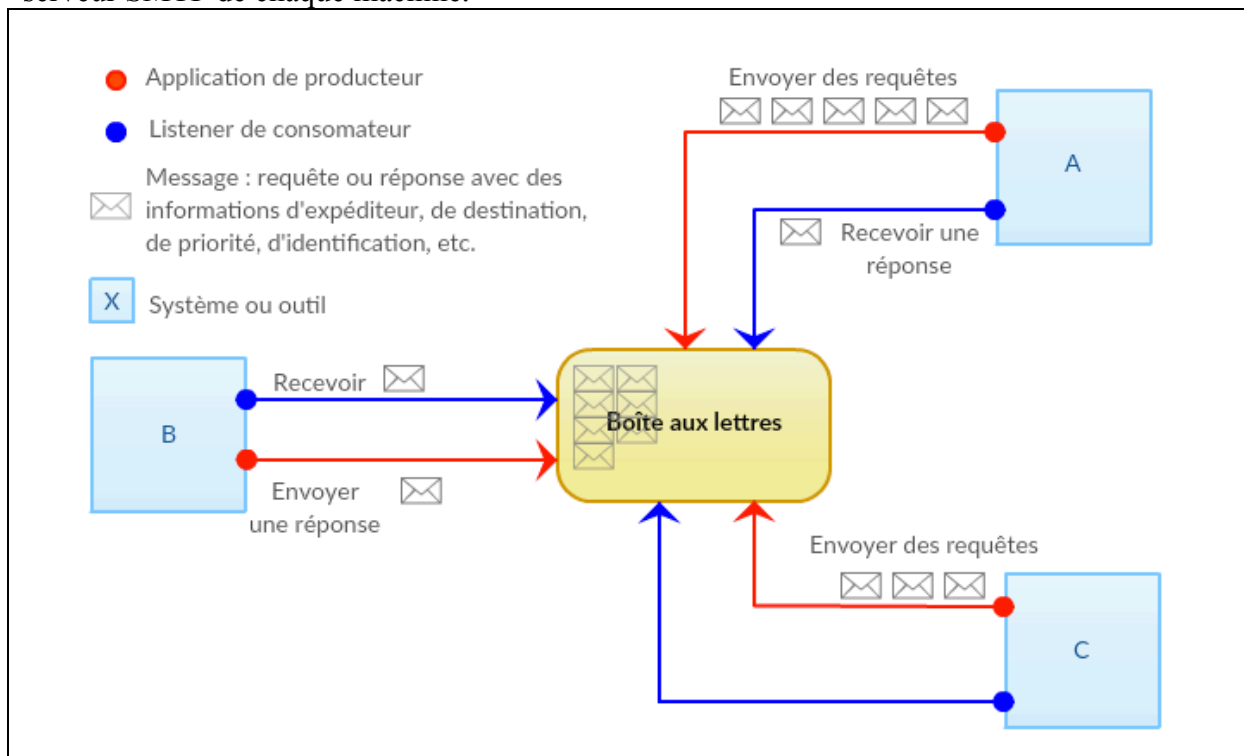


Figure 53 : Architecture initiale de gestion de travaux

Files d'attente & tâches. Après une brève étude de l'existant, nous avons choisi l'intergiciel ACTIVEMQ pour gérer les requêtes. Cet outil permet non seulement de contrôler l'ordre des requêtes, mais aussi de mettre en place une gestion des priorités en donnant une importance particulière à certaines tâches, ou à certains clients. De la même manière, en qualifiant correctement les traitements, il sera facile de mettre en place des serveurs dédiés à certains processus.

IX.2 Implémentation

Pour implémenter la plate-forme JIANDAN-EVAL, nous avons utilisé plusieurs ressources et outils. Cette partie montre les outils utilisés, les composants utilisés, et les composants développés.

IX.2.1 Outils de base utilisés

JIANDAN-EVAL est un projet développé sous NETBEANS 8.0.1⁹³, et nous avons utilisé plusieurs outils et bibliothèques.

- Java version 8,
- JSOUP⁹⁴ pour traiter les documents en format xml/html,
- JUNIT 4.10⁹⁵ pour tester les programmes ; (test unitaires)
- MYSQL Community Server version 5.6⁹⁶,
- Python,
- JSF 2.2⁹⁷ et PrimeFaces 5.2⁹⁸,
- PHP 7⁹⁹.

Les formats des données sont xml, xhtml, txt.

IX.2.2 Composants utilisés

Moses. Le noyau de JIANDAN-EVAL est basé sur Moses et ses scripts. Ils sont utilisés pour la construction de systèmes de TA.

Segmenteurs. Le segmenteur de Moses, et le STANFORD WORD SEGMENTER¹⁰⁰.

Pour compléter les fonctionnalités de JIANDAN-EVAL, nous sommes en train de créer les relations avec les composants suivants.

SECTra : Gestion des MT, et évaluation de type « campagne ».

Tradoh : Intergiciel d'appel à des systèmes de TA, éventuellement via SEGDOC.

Lextoh : Intergiciel permettant d'appeler plusieurs outils d'analyse morphologique et de fusionner leurs résultats.

ActiveMQ : Gestion des tâches et des files d'attente.

Décodeur de XMU : C'est un décodeur statistique différent de celui de MOSES, fabriqué par l'université de Xiamen.

IX.2.3 Composants développés

Construction de systèmes de TA. Ce composant a été développé avec JAVA/JSF et MOSES. Les paramètres sont obtenus via l'interface de JIANDAN-EVAL, qui génère ensuite les commandes de MOSES permettant de construire le système de TA correspondant.

Évaluation de systèmes de TA. Ce composant a été développé par Haozhou WANG dans le cadre de son stage final de Master 2. Il a intégré à JianDan-eval des scripts PYTHON permettant de calculer les scores comme BLEU/NIST et TER.

Déploiement (en cours). Ce composant permet de mettre un système de TA, qui a été entraîné par l'utilisateur, en exploitation sous JIANDAN-EVAL. Nous utilisons la bibliothèque

⁹³ <https://netbeans.org>

⁹⁴ <http://jsoup.org>

⁹⁵ <http://junit.org>

⁹⁶ <http://dev.mysql.com/downloads/mysql/5.6.html>

⁹⁷ <http://www.oracle.com/technetwork/java/javaee/javaserverfaces-139869.html>

⁹⁸ <http://www.primefaces.org>

⁹⁹ <http://php.net>

¹⁰⁰ <http://nlp.stanford.edu/software/segmenter.shtml>

PrimeFaces pour téléverser un système de type de Moses sur un serveur contrôlé par JIANDAN-EVAL.

Amélioration de systèmes de TA (en cours). Ce composant intègre notre script d'entraînement incrémental. L'utilisateur téléverse les nouvelles données (le nouveau corpus ou les segments post-édités) via l'interface de JIANDAN-EVAL, puis lance l'entraînement incrémental.

IX.3 Exemple : création d'un système de TA

Dans le cadre de son stage final de master, Haozhou WANG a proposé une formule qui permet de calculer un score de *qualité d'usage potentielle* pour chaque prétraduction, et de choisir une meilleure prétraduction parmi différentes prétraductions produites par plusieurs systèmes de TA. Les résultats ont montré que cette formule peut aider les post-éditeurs à augmenter leur vitesse de post-édition.

Il a utilisé JIANDAN-EVAL pour fabriquer plusieurs systèmes de TA français-chinois. Les données d'entraînement étaient extraites de la MT de SECTRA. Nous en avons extrait environ 80K bissegments. Après le filtrage, il y avait 75K bissegments.

Haozhou a téléversé ces données sur JIANDAN-EVAL, puis a utilisé les paramètres par défaut (segmenteur STANFORD, GIZA++, ISRILM, 4 N-gram, modèle phrase-based, et MERT). Il a fallu huit heures pour entraîner le système, qui a été déployé sur le serveur de TA, et aussi intégré à Tradoh.

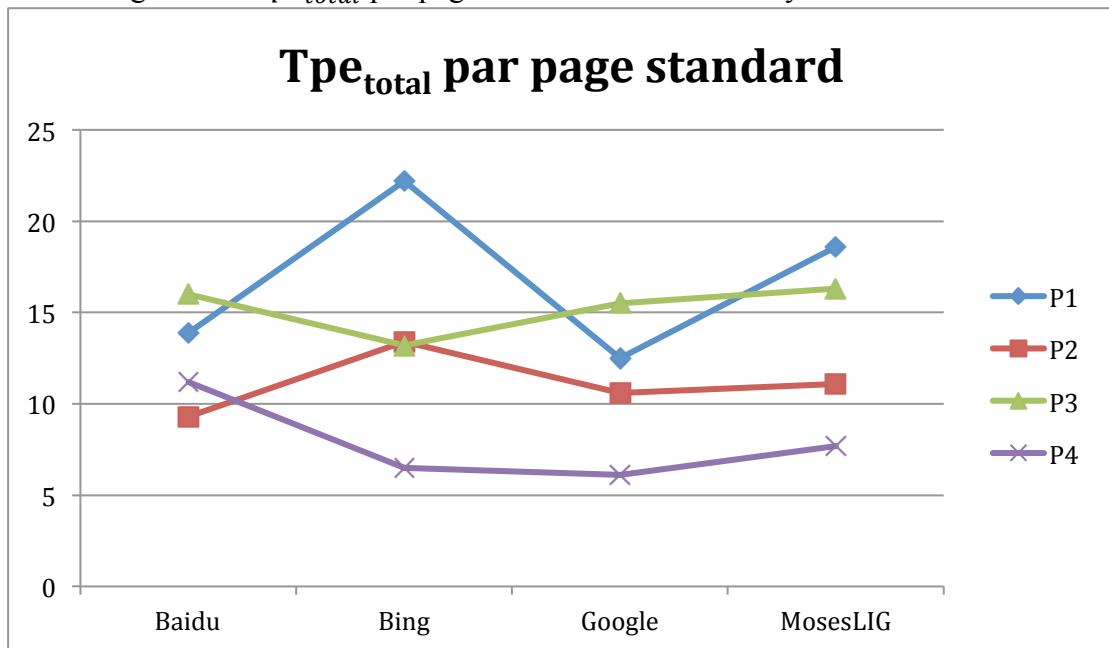
Son corpus de post-édition a été collecté et sélectionné depuis le site du LIG. Il s'agissait d'un corpus de 1000 segments en langue française dans le domaine de l'informatique. Quatre post-éditeurs ont été invités à post-éditer les segments prétraduits par BAIDU FANYI, BING, GOOGLE TRANSLATE, et notre système MOESLIG, dans la plate-forme de SECTRA. Un résumé des données obtenues est présenté dans le Tableau 66.

Tableau 66 : Résumé des données

Post-éditeur	Nombre de segments	Nb de pages standard (source)	Tpe_{total} par page standard (prétraduction de Baidu)	Tpe_{total} par page standard (prétraduction de Bing)	Tpe_{total} par page standard (prétraduction de Google)	Tpe_{total} par page standard (prétraduction de MosesLIG)
P1	600	42.38	13,9 minutes	22,2 minutes	12,5 minutes	18,6 minutes
P2	600	42.38	9,3 minutes	13,4 minutes	10,6 minutes	11,1 minutes
P3	600	42.38	16,0 minutes	13,2 minutes	15,5 minutes	16,3 minutes
P4	400	42.38	11,2 minutes	6,5 minutes	6,1 minutes	7,7minutes

Selon le système de TA utilisé, il semble que l'ordre de productivité des post-éditeurs varie. La Figure 54 suivante l'illustre.

Figure 54 : Tpe_{total} par page standard de différents systèmes de TA



Conclusions et perspectives

Dans le cadre de cette thèse, effectuée dans le cadre d'une bourse CIFRE, et prolongeant un des aspects du projet ANR TRAQUIERO, nous avons d'abord abordé la production, l'extension et l'amélioration de corpus multilingues par traduction automatique (TA) et post-édition contributive (PE). Nous avons apporté des améliorations fonctionnelles et techniques aux logiciels SECTRA_W et IMAG produits lors de thèses antérieures (P.C. Huynh, H.T. Nguyen), et nous pensons avoir progressé vers une définition générique de la structure d'un corpus multilingue, multi-annoté et multimédia, pouvant contenir des documents classiques aussi bien que des pseudo-documents (comme des pages Web) et des méta-segments. Cette partie a été validée par la création de bons corpus bilingues français-chinois, l'un d'eux résultant de la toute première application à la traduction littéraire (un roman de Jules Verne que nous avons traduit en chinois pour améliorer notre connaissance du français).

La seconde partie de cette thèse a initialement été motivée par un besoin industriel. Elle a consisté à construire des systèmes de TA de type Moses, spécialisés à des sous-langages, en français↔chinois, et à étudier la façon de les améliorer dans le cadre d'un usage en continu avec possibilité de PE. Dans le cadre d'un projet interne sur le site du LIG et d'un projet (TABEL-FC) en coopération avec l'université de Xiamen, nous avons pu démontrer l'intérêt de l'apprentissage incrémental en TA statistique, sous certaines conditions, grâce à une expérience qui s'est étalée sur toute la thèse. Lors de la deuxième phase de cette expérience, nous n'avons pas pu faire d'apprentissage incrémental, à cause de l'indisponibilité d'un serveur, et n'avons « recompilé » le système de TA qu'à la fin. Cependant, cette phase a permis de démontrer le gain de qualité obtainable par spécialisation à un sous-langage. La troisième phase, toujours en cours, démontre vraiment l'apport de l'Al, toutes choses égales par ailleurs.

La troisième partie de la thèse a été consacrée à des contributions et mises à disposition de supports informatiques et de ressources. Cet aspect n'est pas seulement intéressant par ses résultats, mais aussi parce que la spécification et surtout l'implémentation de ressources et outils aussi variés et de grande taille pose des problèmes spécifiques, auxquels nous avons trouvé des solutions assez génériques et efficaces.

Nos principales contributions se placent dans le cadre du projet COST MUMIA de l'EU et résultent de l'exploitation de la collection CLEF-IP 2011 de 1,5M de brevets partiellement multilingues. De grosses mémoires de traductions en ont été extraites (17,5 M segments), 3 systèmes de TA en ont été tirés (de-fr, en-fr, fr-de), et un site Web de support à la RI multilingue sur les brevets a été construit.

Enfin, nous avons construit en 2015 deux IMAG utilisées par des doctorants pour construire des MT de bonne qualité en français-somali et français-comorien par post-édition de journaux de Djibouti et des Comores, et ensuite en dériver des systèmes de TA spécialisés aux sous-langages de ces journaux, en extraire du vocabulaire bilingue, etc.

Dans le tout dernier chapitre, nous avons décrit la réalisation en cours de JIANDAN-EVAL, une plate-forme de construction, déploiement et évaluation de systèmes de TA. Ce travail s'est révélé plus lourd que nous ne le pensions initialement, et n'est pas encore achevé. Cependant, une version α a déjà été utilisée depuis début 2015 par H. Wang pour son M2R IdL sur l'estimation a priori (QE) de la qualité (d'usage) des résultats de systèmes de TA variés, toujours en fr-zh.

Les perspectives de cette recherche sont multiples. D'abord, nous voulons terminer l'implémentation de JIANDAN-EVAL, le mettre à disposition de la communauté, et l'évaluer.

Ensuite, nous voudrions continuer notre travail sur la modélisation des corpus de traductions, et arriver à une nouvelle implémentation de SECTRA dans laquelle on pourrait définir un corpus en s'appuyant sur la description de sa structure (*macrostructure*, *microstructure* et *mésosstructure*), écrite dans un langage qui reste encore à définir (sans doute par une DTD ou un schéma XML, puisque nous avons déjà décrit toutes les métadonnées de nos différents corpus en XML).

Enfin, nous souhaitons trouver un cadre approprié pour valoriser notre expertise en TA français→chinois, et pour développer aussi des systèmes chinois→français. Un cadre approprié pourrait être celui de l'accès (de HQ) à des sites boursiers et économiques en français par des sinophones, puis à des sites de même type en chinois par des francophones. Deux projets ont été soumis dans ce sens aux appels des PRC *Cai Yuanpei* et *Xu Guangqi*¹⁰¹.

¹⁰¹ <http://www.campusfrance.org/fr/caiyuanpei>

Bibliographie

- [1] (Alabau, Bonk et al., 2013) Alabau, V., Bonk, R., Buck, C., Carl, M., Casacuberta, F., García-Martínez, M., González, J., Koehn, P., Leiva, L. and Mesa-Lao, B. (2013). CASMACAT: An open source workbench for advanced computer-aided translation. *The Prague Bulletin of Mathematical Linguistics* 100: 101-112.
- [2] (Allen, 1999) Allen, J. (1999). The Postediting Special Interest Group (SIG). *MT News International* 21.
- [3] (Allen, 2003) Allen, J. (2003). Post-editing. *Benjamins Translation Library* 35: 297–318.
- [4] (Arnold, Balkan et al., 1994) Arnold, D., Balkan, L., Meijer, S., Humphreys, R. and Sadler, L. (1994). *Machine translation: An introductory guide*. NCC Blackwell. 200 p.
- [5] (Berment and Boitet, 2012) Berment, V. and Boitet, C. (2012). Heloise—An Ariane-G5 compatible environment for developing expert MT systems online. *Proc. 24th International Conference on Computational Linguistics (COLING)*. Mumbai. 9 p.
- [6] (Bertoldi, Simianer et al., 2014) Bertoldi, N., Simianer, P., Cettolo, M., Wäschle, K., Federico, M. and Riezler, S. (2014). Online adaptation to post-edits for phrase-based statistical machine translation. *Machine Translation* 28(3-4): 309-339.
- [7] (Besacier, 2014) Besacier, L. (2014). Traduction automatisée d'une œuvre littéraire: une étude pilote. *Proc. Traitement Automatique du Langage Naturel (TALN)*. Marseille, France.
- [8] (Blanchon, Boitet et al., 1999) Blanchon, H., Boitet, C. and Caelen, J. (1999). Participation francophone au consortium C-STAR II.
- [9] (Boitet, 1996) Boitet, C. (1996). La synergie entre THAM, réseau et TA comme facteur de progrès théoriques et pratiques en TAO. *Proc. TAL+ AI-96 (le traitement du langage et ses applications industrielles)*. Moncton. 4-6.
- [10] (Boitet, 2009) Boitet, C. (2009). Machine Translation (MT) and Computer-Aided Translation (CAT) — Lecture 1: Linguistic, computational and operational architectures of MT systems. *Proc. NII lecture series*. Tokyo, Japan. Slides available on NII web site. 44 p.
- [11] (Boitet, Bellynck et al., 2008) Boitet, C., Bellynck, V., Mangeot, M. and Ramisch, C. (2008). Towards Higher Quality Internal and External Multilingualization of Web Sites. *Proc. ONII-08 (Summer Workshop on Ontology, NLP, Personalization and IE/IR)* IITB, Mumbai, Inde. 8 p. .
- [12] (Boitet, Bey et al., 2005) Boitet, C., Bey, Y. and Kageura, K. (2005). Main research issues in building web services for mutualized, non-commercial translation. *Proc. the 6th Symposium on Natural Language Processing (SNLP)*. pp. 451–454.
- [13] (Boitet, Blanchon et al., 1994) Boitet, C., Blanchon, H., Seligman, M. and Bellynck, V. (1994). Evolution of MT with the Web. *Proc. Conference "Machine Translation 25 Years On"*. pp. 1-13.
- [14] (Boitet, Boguslavskij et al., 2007) Boitet, C., Boguslavskij, I. M. and Cardeñosa, J. (2007). An evaluation of UNL usability for high quality multilingualization and projections for a future UNL++ language. *Proc. Computational Linguistics and Intelligent Text Processing*. Springer. pp. 361–373.
- [15] (Boitet, Huynh et al., 2010) Boitet, C., Huynh, C.-P., Nguyen, H.-T. and Bellynck, V. (2010). The iMAG concept: multilingual access gateway to an elected Web site with incremental quality increase through collaborative post-edition of MT pretranslations. *Proc. Traitement Automatique du Langage Naturel (TALN)*. Montréal, Canada.
- [16] (Brown, Pietra et al., 1993) Brown, P. F., Pietra, V. J. D., Pietra, S. a. D. and Mercer, R. L. (1993). The mathematics of statistical machine translation: Parameter estimation. *Computational Linguistics* 19(2): pp. 263–311.
- [17] (Callison-Burch, Osborne et al., 2006) Callison-Burch, C., Osborne, M. and Koehn, P. (2006). Re-evaluating the Role of Bleu in Machine Translation Research. *Proc. EACL*. 249-256.
- [18] (Cettolo, Federico et al., 2010) Cettolo, M., Federico, M. and Bertoldi, N. (2010). Mining parallel fragments from comparable texts. *Proc. IWSLT*. 227–234.

- [19] (Chang, Galley et al., 2008) Chang, P.-C., Galley, M. and Manning, C. D. (2008). Optimizing Chinese word segmentation for machine translation performance. Proc. *third workshop on statistical machine translation* Association for Computational Linguistics (ACL-SMT). pp.224-232.
- [20] (Chen, Wang et al., 2014) Chen, Y., Wang, L., Boitet, C. and Shi, X. (2014). On-going Cooperative Research towards Developing Economy-Oriented Chinese-French SMT Systems with a New SMT Framework. Proc. *Traitement Automatique du Langage Naturel (TALN)*. Maresille, France. pp. 401- 408.
- [21] (Chiang, 2012) Chiang, D. (2012). Hope and fear for discriminative training of statistical translation models. *The Journal of Machine Learning Research* 13(1): 1159-1187.
- [22] (Damerau, 1964) Damerau, F. J. (1964). A technique for computer detection and correction of spelling errors. *Communications of the ACM* 7(3): 171-176.
- [23] (Daoud, 2007) Daoud, M. (2007). Towards interactive Multilingual Access Gateways (iMAG). Université Joseph Fourier, Grenoble, France. Rapport de M2R. 60 p.
- [24] (Déjean and Gaussier, 2002) Déjean, H. and Gaussier, E. (2002). Une nouvelle approche à l'extraction de lexiques bilingues à partir de corpus comparables.
- [25] (Delpech, 2013) Delpech, E. (2013). Traduction assistée par ordinateur et corpus comparables: contributions à la traduction compositionnelle. Thèse de doctorat, Université de Nantes: 266 p.
- [26] (Do, 2011) Do, T. N. D. (2011). Extraction de corpus parallèles pour la traduction automatique depuis et vers une langue peu dotée. Thèse de doctorat, Université de Grenoble.
- [27] (Doddington, 2002) Doddington, G. (2002). Automatic evaluation of machine translation quality using n-gram co-occurrence statistics. Morgan Kaufmann Publishers. 138–145.
- [28] (Dong, 1990) Dong, D. Z. (1990). TRANSTAR: a commercial English-Chinese MT system. Proc. *Conf of the Association for Computational Linguistics (ACL)*. Pittsburgh, Pennsylvania, USA. 339–341.
- [29] (Eisele and Chen, 2010) Eisele, A. and Chen, Y. (2010). MultiUN: A Multilingual Corpus from United Nation Documents. Proc. *International conference Language Resources and Evaluation 2011(LREC 2011)*. Istanbul, Turkey. 2868-2872.
- [30] (Esplá-Gomis, 2009) Esplá-Gomis, M. (2009). Bitextor, a free/open-source software to harvest translation memories from multilingual websites. Proc. *Beyond Translation Memories Workshop (MT Summit XII)*. Ottawa, Ontario, Canada.
- [31] (Fafiotte, 2004) Fafiotte, G. (2004). Interprétariat à distance et collecte de dialogues spontanés bilingues, sur une plate-forme générique multifonctionnelle. Proc. *11ème conférence sur le Traitement Automatique des Langues Naturelles (TALN'2004)*, . Fès (Maroc). 10 p.
- [32] (Federico, Bertoldi et al., 2008) Federico, M., Bertoldi, N. and Cettolo, M. (2008). IRSTLM: an open source toolkit for handling large scale language models. Proc. *9th Annual Conference of the International Speech Communication Association 2008 (INTERSPEECH 2008)*. Brisbane, Australia. 1618–1621.
- [33] (Federmann, Giannopoulou et al., 2012) Federmann, C., Giannopoulou, I., Girardi, C., Hamon, O., Mavroeidis, D., Minutoli, S. and Schröder, M. (2012). META-SHARE v2: An Open Network of Repositories for Language Resources including Data and Tools. 3300–3303.
- [34] (Feng, 1981) Feng, Z. (1981). Mémoire pour une tentative de traduction multilingue du chinois en français, anglais, japonais, russe et allemand. Doc. GETA, Grenoble, 40p. + annexes.
- [35] (Forcada, Ginestí-Rosell et al., 2011) Forcada, M. L., Ginestí-Rosell, M., Nordfalk, J., O'regan, J., Ortiz-Rojas, S., Pérez-Ortiz, J. A., Sánchez-Martínez, F., Ramírez-Sánchez, G. and Tyers, F. M. (2011). Apertium: a free/open-source platform for rule-based machine translation. *Machine Translation* 25: 127–144.
- [36] (Ganitkevitch, Cao et al., 2012) Ganitkevitch, J., Cao, Y., Weese, J., Post, M. and Callison-Burch, C. (2012). Joshua 4.0: Packing, PRO, and paraphrases. Proc. *50th Annual Meeting of the Association for Computational Linguistics*. Jeju Island, Korea. 283–291.
- [37] (Gao and Vogel, 2008) Gao, Q. and Vogel, S. (2008). Parallel implementations of word alignment tool. Proc. *Workshop: Software Engineering, Testing, and Quality Assurance for Natural Language Processing*. 49–57.

- [38] (Graff, Kong et al., 2003) Graff, D., Kong, J., Chen, K. and Maeda, K. (2003). English gigaword. Linguistic Data Consortium, Philadelphia.
- [39] (Guilbaud, 1999) Guilbaud, J.-P. (1999). Ariane-G5: Environnement de développement et d'exécution de systèmes (linguiciels) de traduction automatique. GDR I3 ATALA, Paris.
- [40] (Hajlaoui and Boitet, 2008) Hajlaoui, N. and Boitet, C. (2008). TA statistique à petits corpus pour des petits sous-langages. Proc. *TOTH 2008 Conférence sur la Terminologie & Ontologie: Théories et Applications*. 20 p.
- [41] (Hanbury, Boyer et al., 2011) Hanbury, A., Boyer, C., Gschwandtner, M. and Müller, H. (2011). KHRESMOI: towards a multi-lingual search and access system for biomedical information. Med-e-Tel, Luxembourg: 412–416.
- [42] (Hasler, Haddow et al., 2011) Hasler, E., Haddow, B. and Koehn, P. (2011). Margin infused relaxed algorithm for Moses. *The Prague Bulletin of Mathematical Linguistics* 96: 69-78.
- [43] (Hsu and Su, 1995) Hsu, Y.-L. U. and Su, K.-Y. (1995). The New Generation BehaviorTran: Design Philosophy and system architecture. Proc. *ROCLING*. Taiwan. 65-79.
- [44] (Hutchins, 1986) Hutchins, W. J. (1986). Machine translation: past, present, future, Ellis Horwood: Chichester.
- [45] (Huynh, 2010) Huynh, C.-P. (2010). Des suites de test pour la TA à un système d'exploitation de corpus alignés de documents et métadocuments multilingues, multiannotés et multimédia. Thèse de doctorat, Université Joseph Fourier.
- [46] (Kalitvianski, Boitet et al., 2012) Kalitvianski, R., Boitet, C. and Bellynck, V. (2012). Collaborative Computer-Assisted Translation Applied to Pedagogical Documents and Literary Works. *COLING (Demos)*: 255–260.
- [47] (King and Chang, 1963) King, G. W. and Chang, H.-W. (1963). Machine translation of Chinese. *Scientific American* 208: 124-135.
- [48] (Koehn, 2005) Koehn, P. (2005). Europarl: A parallel corpus for statistical machine translation. 79–86.
- [49] (Koehn, 2009) Koehn, P. (2009). *Statistical machine translation*, Cambridge University Press.
- [50] (Koehn, Hoang et al., 2007) Koehn, P., Hoang, H., Birch, A., Callison-Burch, C., Federico, M., Bertoldi, N., Cowan, B., Shen, W., Moran, C., Zens, R. and Others (2007). Moses: Open source toolkit for statistical machine translation. Proc. *Conf of the Association for Computational Linguistics*. 177–180.
- [51] (Kohlschütter, Fankhauser et al., 2010) Kohlschütter, C., Fankhauser, P. and Nejdil, W. (2010). Boilerplate detection using shallow text features. Proc. *third ACM international conference on Web search and data mining (WSDM)*. New York City, NY, USA. 441–450.
- [52] (Levenshtein, 1966) Levenshtein, V. I. (1966). Binary codes capable of correcting deletions, insertions, and reversals. Proc. *Soviet physics doklady*. 707-710.
- [53] (Li, Callison-Burch et al., 2009) Li, Z., Callison-Burch, C., Dyer, C., Ganitkevitch, J., Khudanpur, S., Schwartz, L., Thornton, W. N., Weese, J. and Zaidan, O. F. (2009). Joshua: An open source toolkit for parsing-based machine translation. Proc. *Association for Computational Linguistics*. Singapore. 135–139.
- [54] (Loh and Kong, 1979) Loh, S.-C. and Kong, L. (1979). An interactive on-line machine translation system (Chinese into English). *Translating and the Computer*. North-Holland, Amsterdam: 135-148.
- [55] (Lu, Tsou et al., 2009) Lu, B., Tsou, B. K., Zhu, J., Jiang, T. and Kwong, O. Y. (2009). The construction of a chinese-english patent parallel corpus. *Proceedings of the MT Summit XII*: 17–24.
- [56] (Lü, Huang et al., 2007) Lü, Y., Huang, J. and Liu, Q. (2007). Improving Statistical Machine Translation Performance by Training Data Selection and Optimization. Proc. *EMNLP-CoNLL*. 3-350.
- [57] (Mangeot, 2010) Mangeot, M. (2010). Environnements pour lexicographes et lexicologues: Environnements génériques, centralisés et distribués en ligne de construction collaborative de bases lexicales en contexte multilingue. Thèse, Université Joseph Fourier 280 p.
- [58] (Mirkin, 2014) Mirkin, S. (2014). Incrementally Updating the SMT Reordering Model. Proc. *Proceedings of The 28th Pacific Asia Conference on Language, Information and Computing (PACLIC)*, Phuket, Thailand.

- [59] (Munteanu and Marcu, 2005) Munteanu, D. S. and Marcu, D. (2005). Improving machine translation performance by exploiting non-parallel corpora. *Computational Linguistics* 31(4): 477–504.
- [60] (Munteanu and Marcu, 2006) Munteanu, D. S. and Marcu, D. (2006). Extracting parallel sub-sentential fragments from non-parallel corpora. *Proc. 21st International Conference on Computational Linguistics and the 44th annual meeting of the Association for Computational Linguistics*. Sydney, Australia. Association for Computational Linguistics. 81–88.
- [61] (Nguyen, 2009) Nguyen, H.-T. (2009). Des systèmes de TA homogènes aux systèmes de TAO hétérogènes. Thèse de doctorat, Université de Grenoble
- [62] (Och, 2003) Och, F. J. (2003). Minimum error rate training in statistical machine translation. *Proc. 41st Annual Meeting on Association for Computational Linguistics*. Sapporo, Japan. Association for Computational Linguistics. 160-167.
- [63] (Papineni, Roukos et al., 2002) Papineni, K., Roukos, S., Ward, T. and Zhu, W.-J. (2002). BLEU: a method for automatic evaluation of machine translation. *Proc. Association for Computational Linguistics*. Philadelphia, Pennsylvania, USA. 311–318.
- [64] (Parr and Quong, 1995) Parr, T. J. and Quong, R. W. (1995). ANTLR: A predicated - LL (k) parser generator. *Software: Practice and Experience* 25(7): 789-810.
- [65] (Pineau, 2004) Pineau, M. (2004). Comparaison de résultats de traduction (automatique ou non). Université Joseph Fourier, Rapport de TER.
- [66] (Pouliquen, Elizalde et al., 2013) Pouliquen, B., Elizalde, C., Junczys-Dowmunt, M., Mazenc, C. and García-Verdugo, J. (2013). Large-scale multiple language translation accelerator at the United Nations. *Proc. MT Summit*. Nice, France.
- [67] (Pouliquen and Mazenc, 2011) Pouliquen, B. and Mazenc, C. (2011). COPPA, CLIR and TAPTA: three tools to assist in overcoming the Pat-ent language barrier at WIPO. *Proceedings of the 13th Machine Translation Summit*: 24–30.
- [68] (Ramisch, 2008) Ramisch, C. (2008). Développement d'un site Web iMAG générique et instantiation sur des sites iMAG concrets. Rapport de PFE Ensimag, 2008. 29 p.
- [69] (Resnik, Olsen et al., 1999) Resnik, P., Olsen, M. B. and Diab, M. (1999). The Bible as a parallel corpus: Annotating the 'Book of 2000 Tongues'. *Computers and the Humanities* 33: 129–153.
- [70] (Resnik and Smith, 2003) Resnik, P. and Smith, N. A. (2003). The web as a parallel corpus. *Computational Linguistics* 29: 349–380.
- [71] (Salampasis and Hanbury, 2014) Salampasis, M. and Hanbury, A. (2014). PerFedPat: An integrated federated system for patent search. *World Patent Information* 38: 4–11.
- [72] (San Vicente and Manterola, 2012) San Vicente, I. and Manterola, I. (2012). PaCo2: A Fully Automated tool for gathering Parallel Corpora from the Web. *Proc. Eighth edition of the Language Resources and Evaluation Conference (LREC)* Istanbul, Turkey. 1–6.
- [73] (Schwenk, Barrault et al., 2015) Schwenk, H., Barrault, L., Blain, F., Bougares, F., Hazem, A. and Servan, C. (2015). MateCat - An Open Source CAT Tool with closely integrated User Specific Statistical Machine Translation. *Proc. TAO-CAT 2015*. Angers, France.
- [74] (Senez, 1998) Senez, D. (1998). Post-editing service for machine translation users at the European Commission. *Translating and the Computer* 20.
- [75] (Shi, Chen et al., 2013) Shi, X., Chen, Y. and Huang, X. (2013). Key Problems in Conversion from Simplified to Traditional Chinese Characters. *Proc. XIV Machine Translation Summit* Nice, France. p. 287-293.
- [76] (Silkroad, 2006) Silkroad (2006) 基于短语的统计机器翻译系统 "丝路" 1.0 版 (SilkRoad V1.0) 设计与使用说明.
- [77] (Smith, Quirk et al., 2010) Smith, J. R., Quirk, C. and Toutanova, K. (2010). Extracting parallel sentences from comparable corpora using document level alignment. *Proc. Conf of the Association for Computational Linguistics*. Uppsala, Sweden. 403–411.
- [78] (Snover, Madnani et al., 2009) Snover, M. G., Madnani, N., Dorr, B. and Schwartz, R. (2009). TER-Plus: paraphrase, semantic, and alignment enhancements to Translation Edit Rate. *Machine Translation* 23: 117–127.

- [79] (Steinberger, Pouliquen et al., 2006) Steinberger, R., Pouliquen, B., Widiger, A., Ignat, C., Erjavec, T., Tufis, D. and Varga, D. (2006). The JRC-Acquis: A multilingual aligned parallel corpus with 20+ languages. arXiv preprint cs/0609058.
- [80] (Tian, Wong et al., 2014) Tian, L., Wong, D., Chao, L., Quaresma, P., Oliveira, F., Li, S., Wang, Y. and Lu, Y. (2014). UM-Corpus: a large English-Chinese parallel corpus for statistical machine translation. Proc. *LREC 2014*. ELRA Reykjavik, Iceland.
- [81] (Tiedemann, 2012) Tiedemann, J. (2012). Parallel Data, Tools and Interfaces in OPUS. Proc. *eighth international conference on Language Resources and Evaluation (LREC)*. Istanbul, Turkey. 2214–2218.
- [82] (Uchino, Shirai et al., 2001) Uchino, H., Shirai, S., Yokoo, A., Oyama, Y. and Furuse, O. (2001). ALTFLASH: A Japanese-English Machine Translation System for Market Flash Reports. *EICE Transactions on Information and Systems, Pt. 2 (Japanese Edition)*: 1167-1174.
- [83] (Utiyama and Isahara, 2007) Utiyama, M. and Isahara, H. (2007). A Japanese-English patent parallel corpus. *Proceedings of MT summit XI*: 475–482.
- [84] (Varga, Halácsy et al., 2007) Varga, D., Halácsy, P., Kornai, A., Nagy, V., Németh, L. and Trón, V. (2007). Parallel corpora for medium density languages. *Amsterdam studies in the theory and history of linguistic science series 4*. 247.
- [85] (Vasconcellos and León, 1985) Vasconcellos, M. and León, M. (1985). SPANAM and ENGSPAN: machine translation at the Pan American Health Organization. *Computational Linguistics* 11(2-3): 122-136.
- [86] (Vo Trung, 2004) Vo Trung, H. (2004). Méthodes et outils pour utilisateurs, développeurs et traducteurs de logiciels en contexte multilingue. Thèse de doctorat, Université Stendhal.
- [87] (Vogel, Ney et al., 1996) Vogel, S., Ney, H. and Tillmann, C. (1996). HMM-based word alignment in statistical translation. *Association for Computational Linguistics*. 836–841.
- [88] (Wagner and Fischer, 1974) Wagner, R. A. and Fischer, M. J. (1974). The string-to-string correction problem. *Journal of the ACM (JACM)* 21(1): 168-173.
- [89] (Wang, 2015) Wang, H. (2015). Évaluation comparative de la qualité d'usage de plusieurs systèmes de TA français-chinois en fonction de la tâche de post-édition. *Mémoire de M2R*.
- [90] (Wang and Boitet, 2013) Wang, L. and Boitet, C. (2013). Online production of HQ parallel corpora and permanent task-based evaluation of multiple MT systems: both can be obtained through iMAGs with no added cost. Proc. *MT Summit XIV Workshop on Post-editing Technology and Practice*. Nice, France.
- [91] (Xiao, Zhu et al., 2012) Xiao, T., Zhu, J., Zhang, H. and Li, Q. (2012). NiuTrans: an open source toolkit for phrase-based and syntax-based machine translation. Proc. *ACL 2012 System Demonstrations*. Jeju Island, South Korea. 19–24.
- [92] (Zhao and Vogel, 2002) Zhao, B. and Vogel, S. (2002). Adaptive parallel sentences mining from web bilingual news collection. Proc. *IEEE International Conference on*. IEEE. pp. 745–748.

Table des définitions

- Définition 1.** Une *phrase* (*sentence* en anglais) est l'unité élémentaire d'un énoncé, formée de plusieurs mots ou groupes de mots, et qui présente un sens complet. (*TheFreeDictionary*).
- Définition 2*.** Un *syntagme* ou *groupe* (*phrase* en anglais) est un constituant possible d'une phrase. En général, un *titre* est un groupe nominal.
- Définition 3.** Un *segment* est l'unité de traduction de base des traducteurs humains. Il s'agit d'une phrase, d'un titre, ou d'un terme dans une nomenclature.
- Définition 4*.** Un *fragment* (*chunk* en anglais) est une partie d'un segment, qui peut être un *groupe syntaxique* (*phrase* en anglais) ou un simple n-gramme.
- Définition 5*.** Un segment peut contenir des *éléments non textuels*, ou *hors-texte*, comme des images, des formules, ou des balises, qui ont un rôle linguistique et une valeur non linguistique.
- Définition 6.** Un *segment monolingue* est un segment dont le contenu textuel est en une seule langue.
- Définition 7*.** Un *segment multilingue* est un segment dont le contenu textuel est dans plusieurs langues, chaque version étant considérée comme contenant exactement la même information, exprimée de façon correcte.
- Définition 8.** Un *segment monolingue multilingualisé (annoté)* est un objet contenant un segment « source » primaire, une ou plusieurs traductions (automatiques ou humaines ou automatiques post-éditées) pour une ou plusieurs langues, et des annotations, en général des objets, comme des arbres linguistiques, des graphes UNL, des résultats d'évaluation(s), et des références aux contributions ayant produit chaque objet non primaire.
- Définition 9*.** Un *segment multilingue multilingualisé (annoté)* est un objet contenant un segment multilingue, dans N langues « sources », et, dans M autres langues, une ou plusieurs traductions (automatiques ou humaines ou automatiques post-éditées), ainsi que des annotations, comme celles d'un segment monolingue multilingualisé et annoté.
- Définition 10*.** Le *chemin traductionnel* d'une annotation, en particulier d'une traduction ou d'une post-édition, est la suite des opérations l'ayant produite, ainsi que les intervenants humains impliqués, et les éventuels objets auxiliaires utilisés.
- Définition 11*.** Un *métasegment* est un segment comportant une ou plusieurs *variables*, éventuellement typées (nombre, date, balise faible...).
- Définition 12*.** Un *document* est un ensemble formé par un support et une information, celle-ci enregistrée de manière persistante. Nous nous intéressons aux documents textuels, qui contiennent des "segments" textuels.
- Définition 13*.** Un *métadocument* est un document pouvant contenir des métasegments.
- Définition 14*.** Un *pseudo-document* est défini par une référence (nom de fichier, url, uri) à un document qui peut varier au cours du temps.
- Définition 15.** *Contexte* : Le contexte *m-n* d'un segment source par rapport à un document, ou plus généralement à une instance d'un pseudo-document, est défini par :
- la liste des *m* segments (de même langue) qui le précèdent. - la liste des *n* segments (de même langue) qui le suivent. - l'instanciation des variables, s'il y en a, dans ces $m+n+1$ segments. Le contexte *m-n* d'un segment (cible) dans une version résultat de TA ou de PE, dans un segment monolingue, ou multilingue multilingualisé, est défini par :
- la liste des *m* segments (de même langue et de même version) qui le précèdent. - la liste des *n* segments (de même langue et de même version) qui le suivent. - l'instanciation des variables, s'il y en a, dans ces $m+n+1$ segments pour la même version.
- Définition 16*.** Un *corpus* est un ensemble usuellement fermé de documents homogènes du point de vue de leur structure, de leur(s) langue(s), de leur genre et de leur domaine.
- Définition 17*.** Un *corpus de traductions* est un corpus au sens précédent, contenant les traductions de tout ou partie de ses segments, dans une ou plusieurs langues.

Annexes

Annexe 1 : Corpus de la campagne d'évaluation de TA du projet TRANSAT

Corpus de tâche d'assistance (FT Assistance)

Composition de corpus :

1. Répartition des tours de parole dans les différentes sous-tâches
Santé : 50,85 %
Accidents, perte, vol : 24,51 %
Itinéraire : 16,28 %
Transports, et Location de voiture : 11,47 %
Non spécifique : 8,63 %
2. Tours de parole
Nombre de mots par tour de parole : moyenne = 6,2 ; min = 1 ; max = 30
Nombre de tours de parole à 1 phrase : 2083
Nombre de tours de parole à 2 phrases : 135
Nombre de tours de parole à 3 phrases : 5
3. Phrases
Nombre de mots par phrase : moyenne = 5,9 ; médiane = 5 ; min = 1 ; max = 28
Nombre de questions marquées : 717
Nombre d'affirmations (+ questions non marquées) : 1656
4. Vocabulaire
Nombre d'occurrences : 13833
Nombre d'occurrences différentes : 1319
Nombre des occurrences apparaissant 1 fois : 590 (45%)
Nombre des occurrences apparaissant 2 fois : 203 (15%)
Nombre des occurrences apparaissant 3 fois : 98 (7%)
Nombre des occurrences apparaissant 4 fois : 65 (5%)

Corpus de dialogues portant sur la restauration (FT Restaurant)

Composition de corpus :

1. Tours de parole
Nombre de mots par tour de parole : moyenne = 5.2 ; min = 1 ; max = 22
Nombre de tours de parole à 1 phrase : 1896

Nombre de tours de parole à 2 phrases : 97

Nombre de tours de parole à 3 phrases : 5

Nombre de tours de parole à 4 phrases : 2

2. Phrases

Nombre de mots par phrase : moyenne = 4,93 ; médiane = 5 ; min = 1 ; max = 19

Nombre de questions : 783

Nombre d'affirmations : 1330

3. Vocabulaire

Nombre d'occurrences : 10429

Nombre d'occurrences différentes : 855

Nombre des occurrences apparaissant 1 fois : 352 (41%)

Nombre des occurrences apparaissant 2 fois : 133 (16%)

Nombre des occurrences apparaissant 3 fois : 64 (7%)

Nombre des occurrences apparaissant 4 fois : 46 (5%)

Annexe 2 : Protocole d'évaluation pour le projet TRANSAT

Évaluation subjective à la NIST

Nous avons mis en œuvre un protocole légèrement différent du protocole proposé par le NIST dans les campagnes TIDES d'évaluation de systèmes de traduction automatique.

a) Fluidité

Pour l'évaluation de la fluidité, H. Blanchon a choisi de proposer une échelle de 3 notes, au lieu de 5 dans le protocole NIST standard :

- (F1) formulation parfaitement compréhensible sans effort, que le style soit écrit ou oral.
- (F2) formulation acceptable à l'oral, éventuellement compréhensible en faisant un effort.
- (F3) formulation non acceptable.

b) Adéquation

Pour l'évaluation de l'adéquation (transport de l'information pertinente de la source vers la cible), nous avons utilisé une échelle de cinq valeurs, comme dans le protocole NIST standard :

- (A1) Toute l'information est transportée.
- (A2) Presque toute l'information est transportée.
- (A3) La moitié de l'information est transportée.
- (A4) Peu d'information est transportée.
- (A5) Aucune information n'est transportée, ou il y a un contresens.

Évaluation objective

a) Évaluation objective fondée sur la distance d'édition

Nous avons choisi trois mesures: la distance d'édition en mots, la distance d'édition en caractères, et une distance d'édition pondérée.

Les distances d'édition [Damerau, 1964 ; Levenshtein, 1966 ; Wagner et Fischer, 1974] en mots et en caractères considèrent les opérations d'insertion, de suppression et de remplacement en attribuant à chacune de ces opérations un poids de 1.

Ces mesures permettent de se rendre compte du travail de correction nécessaire à l'obtention de traductions utiles pour la tâche à partir des traductions candidates produites par le système.

Ce genre de distance est aussi utilisé lors des évaluations du projet GALE dans le cadre de la méthode d'évaluation HTER [Przybocki et al., 2006 ; Snover et al., 2006].

La distance d'édition pondérée combine la distance d'édition en mots et la distance d'édition en caractères en donnant un poids de 0,2 à la première et un poids de 0,8 à la seconde. Il est en effet plus rapide de faire des manipulations sur les mots (un double clic pour la sélection) que sur des caractères individuels.

Cette méthode correspond à des propositions de [Blanchon et Boitet, 2007].

b) Évaluation objective n-gramme

Nous avons intégré dans SECTra_w les scripts fournis par NIST pour calculer BLEU et NIST. Cependant, ils n'ont pas été expérimentés pendant la réalisation du projet TRANSAT lui-même, à cause de la limite de temps d'implémentation de SECTra_w.1, mais seulement un peu plus tard (fin décembre 2007).

Annexe 3 : Un exemple du corpus B@bel

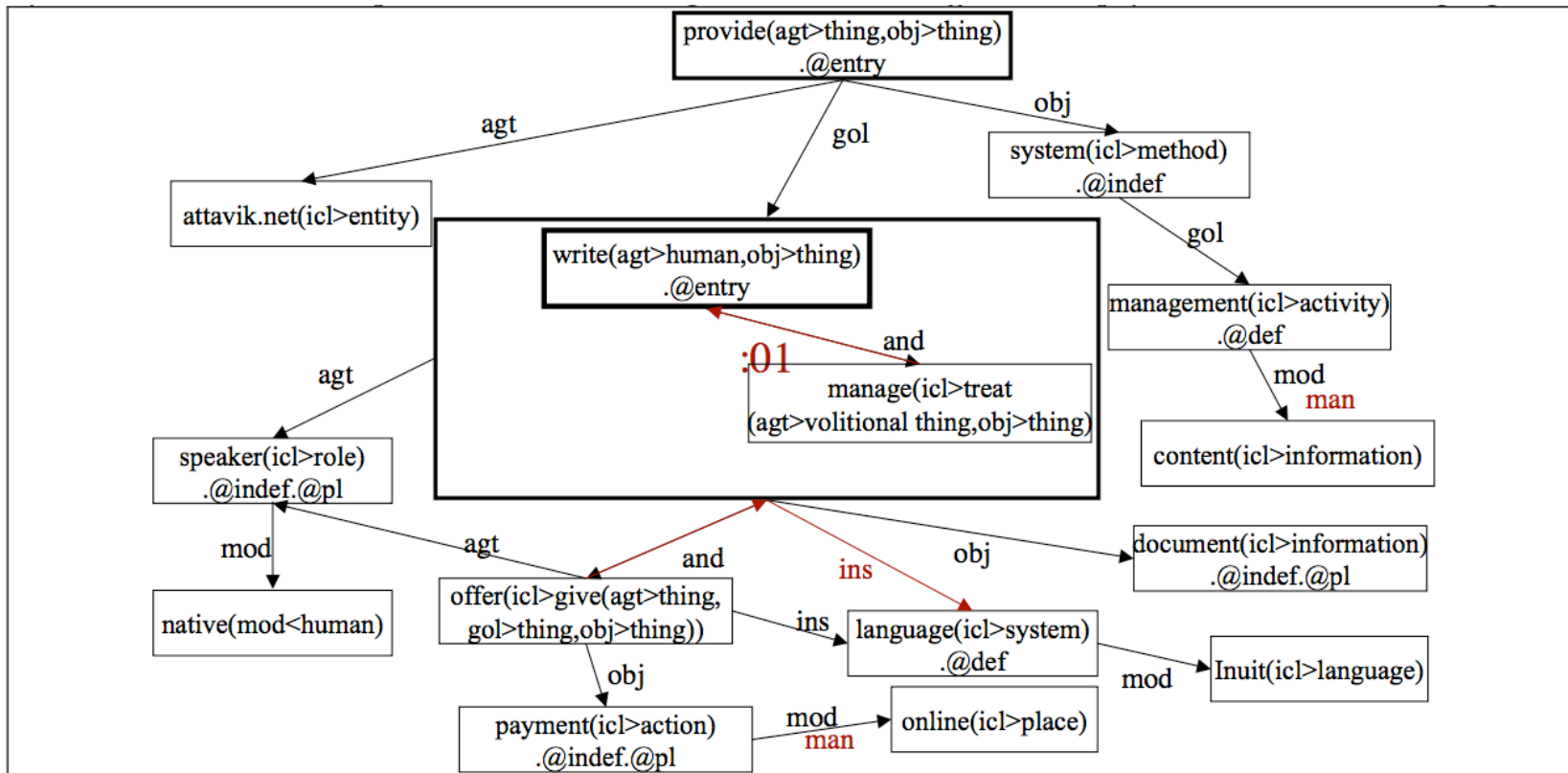
ID (n° en cours)	nb_car	nb_car_val	nb_mots	nb_tot_mots	nb_pages	h_debut	h_fin	duree_incr	duree_tot	mn_p_page	original (anglais)	traduction finale (français)	tr-SYSTRAN-5 (en_fr)
189	290	1196	35	8321	1,74				36	20,7	l. 1723 (11195/3779): 1196 car, 143 mots		
16601_na me_en_1	25	28	3	3	0,01	0	0	0	0	0,0	Multilingual Web Browser.	Navigateur Web multilingue.	Web browser multilingue.
16601_de scshort_e n_3	223	223	24	27	0,10	0	0	0	0	0,0	This project was carried out within Initiative B@bel in cooperation with SIL International to support efforts aimed at developing software/tools promoting multilingualism in cyberspace.	Ce projet a été mené au sein de l'initiative B@bel en coopération avec SIL international pour soutenir des efforts visant à développer les logiciels/outils favorisant le multilinguisme dans le Cyberspace.	Ce projet a été mis à exécution dans l'initiative B@bel en coopération avec SIL international pour soutenir des efforts visés développant le logiciel/outils favorisant le multilinguisme dans Cyberspace.
16601_de sclong_en _5	207	207	26	53	0,10	0	0	0	0	0,0	SIL International has developed Graphite engine which supports the display of complex and non-Roman scripts and is available for free download on the SIL International's website.	SIL international a développé le moteur Graphite qui supporte l'affichage des scripts complexes et non-romains et est disponible en téléchargement libre sur le site Web de SIL international.	SIL international a développé le moteur de graphite qui soutient l'affichage des manuscrits complexes et non-Romains et est disponible pour le téléchargement libre sur le site Web international de SIL.
16601_de sclong_en _6	239	239	28	81	0,11	0	0	0	0	0,0	The project will involve the incorporation of graphite's unique functionalities in other software applications, thereby contributing to the creation and dissemination of content in many currently lesser-used languages.	Le projet comportera l'incorporation des fonctionnalités uniques de Graphite dans d'autres applications logicielles, contribuant de ce fait à la création et à la diffusion du contenu dans beaucoup de langues actuellement moins utilisées.	Le projet comportera l'incorporation des fonctionnalités uniques du graphite dans d'autres applications de logiciel, contribuant de ce fait à la création et à la diffusion du contenu dans beaucoup de langues actuellement peu de-utilisées.
16601_de sclong_en _7	177	178	21	102	0,08	0	0	0	0	0,0	These products will also be freely disseminated with basic documentation facilitating the incorporation of Graphite by software developers in other products.	Ces produits seront également librement disséminés avec la documentation de base facilitant l'incorporation de Graphite par des réalisateurs de logiciel dans d'autres produits.	Ces produits également seront librement disséminés avec la documentation de base facilitant l'incorporation du graphite par des réalisateurs de logiciel dans d'autres produits.
16601_de sclong_en _8	25	28	3	105	0,01	0	0	0	0	0,0	Multilingual web browser.	Navigateur Web multilingue.	Web browser multilingue.
16601_de sclong_en _9	140	141	17	122	0,07	0	0	0	0	0,0	Web-page/site creation is one of the most common form of web publishing and information dissemination in cyberspace.	La création de page/sites Web est une des formes les plus communes d'édition sur le Web et de diffusion de l'information dans le cybersapce.	Le page Web/création d'emplacement est un de la forme la plus commune d'édition de Web et de diffusion de l'information dans le cybersapce.
16601_de sclong_en _10	197	197	30	152	0,12	0	0	0	0	0,0	By developing a beta version of a web browser that supports creation and viewing of web pages in Burmese the ability to create and disseminate multilingual information will be promoted.	En développant une version bêta d'un navigateur Web qui supporte la création et la visualisation des pages Web en Birman, la capacité de créer et diffuser l'information multilingue sera favorisée.	En développant une bêta version d'un web browser qui soutient la création et la visualisation des pages Web dans le Birman la capacité de créer et diffuser l'information multilingue sera favorisé.
16601_de sclong_en _11	79	79	10	162	0,04	0	0	0	0	0,0	The open-source Mozilla browser has been used for this development.	Le navigateur Mozilla à source ouvert a été employé pour ce développement.	Le navigateur de Mozilla d'ouvrir-source a été employé pour ce développement.

Annexe 4 : Corpus EOLSS

No	Titles of articles	UNL files	UW files	Delivery date d/m/y	Latest update d/m/y	snum	wnum	Remarks received L
D1	ETHICS AND SCIENCE	E1-37-05-14-TXT.unl	UWs	03/12/2007	26/06/2008	311	5633	ab, es, ru
D2	GROUND AND SOIL WATER CHARACTERISTICS	E2-03-05-TXT.unl	UWs	03/12/2007	01/07/2008	502	6911	ab, es, jp, ru
D3	HUMAN INTERACTION WITH LAND AND WATER	E2-24D-04-05-TXT.unl	UWs	03/12/2007	07/07/2008	512	7628	ab, es, ru
D4	THE DUBLIN PRINCIPLES	E2-24M-02-04-TXT.unl	UWs	03/12/2007	12/07/2008	591	9146	ab, es, ru
D5	TSUNAMIS	E4-06-01-06-TXT.unl	UWs	03/12/2007	17/07/2008	515	7940	ab, es, ru
D6	WATER TREATMENT : EQUIPMENT AND PROCESSES	E2-13-03-TXT.unl	UWs	08/01/2008	21/07/2008	397	7269	ab, es, ru
D7	CLIMATE CHANGE AND WATER RESOURCES	E2-24D-02-03-TXT.unl	UWs	08/01/2008	28/07/2008	351	6060	ab, es, ru
D8	ANALYSIS OF WATER QUALITY	E2-13-01-06-TXT.unl	UWs	08/01/2008	30/07/2008	454	6210	ab, es, ru
D9	ARTIFICIAL GROUNDWATER RECHARGE	E2-09-06-06-TXT.unl	UWs	07/01/2008	09/06/2008	512	6672	ab, es, ru
D10	BETWEEN THE GREAT RIVERS : WATER IN THE MIDDLE EAST AND NORTH AFRICA	E2-25-06-TXT.unl	UWs	07/01/2008	05/08/2008	466	6403	ab, es, ru
D11	BIOGRAPHIES OF EMINENT WATER RESOURCES PERSONALITIES	E2-23-01-01-TXT.unl	UWs	08/02/2008	11/08/2008	561	9690	es
D12	CHEMICAL PROPERTIES OF SOIL AND GROUND WATERS	E2-03-05-02-TXT.unl	UWs	11/02/2008	13/08/2008	590	9516	es, ru
D13	COMPOSITION AND STRUCTURE OF THE ATMOSPHERE	E4-02-01-01-TXT.unl	UWs	12/02/2008	20/08/2008	603	9824	es
D14	DRINKING WATER SUPPLY	E2-13-01-02-TXT.unl	UWs	25/02/2008	25/08/2008	581	9245	es, ru
D15	FOOD AND WATER DEMAND AND SUPPLY IN 2025	E2-24M-03-04-TXT.unl	UWs	07/02/2008	28/08/2008	430	6620	es
D16	GLOBAL AND REGIONAL FRESHWATER RESOURCES	E2-25-01-TXT.unl	UWs	10/03/2008	03/09/2008	747	13184	
D17	GLOBAL WATER NEEDS FOR THE FUTURE	E2-25-01-03-TXT.unl	UWs	17/03/2008	07/09/2008	463	7366	ru
D18	HEALTH PROBLEMS AND THEIR RESOLUTION, 02, 03, 04	E2-20B-04-TXT.unl, 02, 03, 04	UWs	11/04/2008	15/09/2008	796	12402	S:1-S:26 - ru-1
D19	INDUSTRIAL WATER	E2-19-02-04-TXT.unl	UWs	24/03/2008	16/09/2008	417	5960	better checked
D20	PERSPECTIVES OF GLOBAL WATER BALANCE AND REGIONAL WATER RESOURCES	E2-25-01-01-TXT.unl	UWs	10/04/2008	22/09/2008	523	8977	
D21	PROPERTIES OF RIVERS, STREAMS, LAKES AND WETLANDS	E2-03-04-TXT.unl	UWs	12/06/2008	08/07/2008	757	9800	
D22	RURAL WATER SUPPLY SYSTEMS	E2-14-03-03-TXT.unl	UWs	09/06/2008	08/07/2008	472	6472	
D23	WATER AND WASTEWATER TREATMENT	E2-13-TXT.unl	UWs	05/06/2008	08/07/2008	762	12622	
D24	WATER AS A FACTOR IN SOCIO-ECONOMIC DEVELOPMENT FUTURE TRENDS	E2-25-03-03-TXT.unl	UWs	02/06/2008	08/07/2008	611	9628	
D25	WATER SUPPLY FOR AGRICULTURE	E2-13-01-03-TXT.unl	UWs	19/06/2008	08/07/2008	495	8195	
						12,919		

Annexe 5 : Un exemple de graphe UNL avec correction

La phrase anglais : It [Attawik.net] provides a content management system that allows native speakers to write, manage documents and offer online payments in the Inuit language.



```
agt( provide(agt>thing,obj>thing).@entry,attavik.net(icl>entity))
obj( provide(agt>thing,obj>thing).@entry,system(icl>method).@indef)
gol( system(icl>method).@indef,management(icl>activity).@def)
mod( management(icl>activity).@def,content(icl>information))
gol( provide(agt>thing,obj>thing).@entry,:01)
and:01(write(agt>human,obj>thing).@entry,manage(icl>treat(agt>volitional thing,obj>thing)))
obj(:01,document(icl>information).@indef.@pl)
agt(:01,speaker(icl>role).@indef.@pl)
mod(speaker(icl>role).@indef.@pl,native(mod<human))
and(:01,offer(icl>give(agt>thing,gol>thing,obj>thing))
obj(offer(icl>give(agt>thing,gol>thing,obj>thing)),payment(icl>action).@indef.@pl)
mod(payment(icl>action).@indef.@pl,online(icl>place))
ins(offer(icl>give(agt>thing,gol>thing,obj>thing)),language(icl>system).@def)
mod(language(icl>system).@def,Inuit(icl>language))
agt(offer(icl>give(agt>thing,gol>thing,obj>thing)),speaker(icl>role).@indef.@pl)
```

Annexe 6 : Document de brevet du corpus CLEF-IP 2011 (EP-0000007-B2.xml)

1	<?xml version="1.0" encoding="UTF-8"?>
2	<!DOCTYPE patent-document
3	PUBLIC "-//MXW//DTD patent-document XML//EN" "http://www.ir-facility.org/dtds/patents/v1.4/patent-document.dtd">
4	<patent-document ucid="EP-0000007-B2" country="EP" doc-number="0000007" kind="B2" lang="FR" date="19841121" family-id="19728598" date-produced="20100220" status="new">
5	<bibliographic-data>
6	<publication-reference fvid="19066342" ucid="EP-0000007-B2" status="new">
7	<document-id status="new" format="original">
8	<country status="new">EP</country>
9	<doc-number>0000007</doc-number>
10	<kind>B2</kind>
11	<date>19841121</date>
12	<lang>FR</lang>
13	</document-id>
14	</publication-reference>
15	<application-reference mxw-id="PAPP16180943" ucid="EP-78200026-A" load-source="docdb" status="new" is-representative="N0">
16	<document-id format="epo" status="new">
17	<country status="new">EP</country>
18	<doc-number>78200026</doc-number>
19	<kind>A</kind>
20	<date>19780601</date>
21	<lang>FR</lang>
22	</document-id>
23	</application-reference>
24	<priority-claims status="new">
25	<priority-claim mxw-id="PPC19548354" ucid="LU-77489-A" status="new">
26	<document-id format="epo" status="new">
27	<country status="new">LU</country>
28	<doc-number>77489</doc-number>
29	<kind>A</kind>
30	<date>19770606</date>
31	</document-id>
32	</priority-claim>
33	</priority-claims>
34	<dates-of-public-availability status="new">
35	<intention-to-grant-date>
36	<date>19800715</date>
37	</intention-to-grant-date>
38	</dates-of-public-availability>
39	<technical-data status="new">
40	<classification-ipc status="new">
41	<edition>3</edition>
42	<main-classification status="new">C08F 10/00</main-classification>

43	<further-classification status="new">C08F 4/64</further-classification>
44	<further-classification status="new">C08F 4/02</further-classification>
45	</classification-ipc>
46	<classifications-ipc>
47	<classification-ipc mxw-id="PCL132952396" load-source="docdb" status="new">C08F 4/659 20060101ALI20060310RMJP </classification-ipc>
48	<classification-ipc mxw-id="PCL132963917" load-source="docdb" status="new">C08F 10/00 20060101A I20051008RMEP </classification-ipc>
49	<classification-ipc mxw-id="PCL132976798" load-source="docdb" status="new">C08F 10/00 20060101C I20051008RMEP </classification-ipc>
50	<classification-ipc mxw-id="PCL133045672" load-source="docdb" status="new">C08F 4/00 20060101C I20060521RMEP </classification-ipc>
51	<classification-ipc mxw-id="PCL133100514" load-source="docdb" status="new">C08F 4/00 20060101AFI20051220RMJP </classification-ipc>
52	<classification-ipc mxw-id="PCL133138324" load-source="docdb" status="new">C08F 4/64 20060101A I20060521RMEP </classification-ipc>
53	<classification-ipc mxw-id="PCL133141699" load-source="docdb" status="new">C08F 4/60 20060101ALI20060310RMJP </classification-ipc>
54	</classifications-ipc> <classification-ecla status="new">
55	<classification-symbol scheme="EC" mxw-id="PCL133183878">C08F 10/00+4/655D</classification-symbol>
56	</classification-ecla>
57	<invention-title mxw-id="PT16385698" lang="EN" load-source="patent-office" status="new">PROCESS FOR THE POLYMERISATION OF ALPHA-OLEFINS AND METHOD FOR PREPARING SOLID CATALYTIC COMPLEXES FOR USE IN THIS POLYMERISATION PROCESS</invention-title>
58	<invention-title mxw-id="PT72834363" lang="DE" load-source="patent-office" status="new">Verfahren zur Polymerisation von alpha-Olefinen und Verfahren zur Herstellung von in diesem Polymerisationsverfahren verwendbaren festen katalytischen Komplexen</invention-title>
59	<invention-title mxw-id="PT72834364" lang="FR" load-source="patent-office" status="new">Procédé pour la polymérisation des alpha-oléfines et procédé de préparation de complexes catalytiques solides utilisables pour cette polymérisation</invention-title> </technical-data> <parties>
60	<applicants>
61	<applicant mxw-id="PPAR77323081" sequence="1" format="intermediate" status="new">
62	<addressbook>
63	<name>SOLVAY & CIE (SOCIETE ANONYME)</name> </addressbook> </applicant>
64	</applicants>
65	<inventors>
66	<inventor mxw-id="PPAR77363925" sequence="1" format="intermediate" status="new">
67	<addressbook>
68	<name>BIENFAIT, CHARLES</name>
69	</addressbook>
70	</inventor>
71	<inventor mxw-id="PPAR269225690" sequence="1" format="original" status="new">
72	<addressbook>
73	<last-name>BIENFAIT, CHARLES</last-name>
74	<address>
75	<street>Mereldreef, 75</street>
76	<city>B-2850 Keerbergen</city>
77	<country status="new">BE</country>
78	</address>
79	</addressbook>
80	</inventor>
81	</inventors>
82	<assignees>
83	<assignee mxw-id="PPAR269225691" sequence="1" format="original" status="new">
84	<addressbook>

85	<last-name>SOLVAY & Cie (Société Anonyme)</last-name>
86	<address>
87	<street>Rue du Prince Albert, 33</street>
88	<city>B-1050 Bruxelles</city>
89	<country status="new">BE</country>
90	</address>
91	</addressbook>
92	</assignee>
93	</assignees>
94	</parties>
95	<international-convention-data>
96	<designated-states>
97	<ep-contracting-states>
98	<country mxw-id="DS101191" status="new">BE</country>
99	<country mxw-id="DS101192" status="new">CH</country>
100	<country mxw-id="DS101193" status="new">DE</country>
101	<country mxw-id="DS101194" status="new">FR</country>
102	<country mxw-id="DS101195" status="new">GB</country>
103	<country mxw-id="DS101196" status="new">LU</country>
104	<country mxw-id="DS101197" status="new">NL</country>
105	<country mxw-id="DS101198" status="new">SE</country>
106	</ep-contracting-states>
107	</designated-states>
108	</international-convention-data>
109	</bibliographic-data>
110	<copyright>User acknowledges that the Information Retrieval Facility (IRF) and its third party providers retain all right, title and interest in and to this xml under applicable copyright laws. User acquires no ownership rights to this xml including but not limited to its format. User hereby accepts the terms and conditions of the Licence Agreement set forth at http://www.ir-facility.org/legal/marec/data_licence </copyright>
111	</patent-document>

Annexe 7: Structure des données de minidictionnaires

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <!DOCTYPE SETRAMINIDICO "SETRAMINIDICO.dtd">
3 <minidico>
4   <description>
5     <format> minidico_sectra </format>
6     <segment> Hello World </segment>
7     <langueSource> eng </langueSource>
8     <langueCible> zho </langueCible>
9     <lemmatiseur> xip </lemmatiseur>
10    <serv_dico> Pivax </serv_dico>
11    <id> demo2_doc18_seg1059 </id>
12    <date> 20150705-10:59:54 </date>
13  </description>
14  <dictionnaires>
15    <dictionnaire>
16      <lemme> hello </lemme>
17      <traductions>
18        <traduction> (用于问候、接电话或引起注意) 你好, 喂, (用于询问所处的地方是否有人) 请问有人在吗 </traduction>
19        <traduction> (表示惊讶) 嘿 </traduction>
20      </traductions>
21    </dictionnaire>
22    <dictionnaire>
23      <lemme> world </lemme>
24      <traductions>
25        <traduction> 世界 </traduction>
26        <traduction> 地球和包括其所有的棲息生物, 及位於其上的所有東西(环境) </traduction>
27      </traductions>
28    </dictionnaire>
29  </dictionnaires>
30 </minidico>
```

Annexe 8 : 50 segments en « Vue SECTra/Post-édition »

Num	Segment source	Prétraduction	PE	D_{mix}	Tpe_1
Baidu					
01	Un nouveau plan de croissance pour l'Europe.	一个新的增长计划的欧洲。	欧洲的一个新的增长计划的欧洲。	5,8	07
02	Les dirigeants de l' Union européenne en sont maintenant convaincus : le redémarrage de la croissance en 2012 n'était-il pas l'item prioritaire du programme de la réunion du 30 janvier du Conseil européen.	欧盟领导人在现在相信：启动2012年的增长不是优先项目计划会议1月30日欧洲理事会。	欧盟领导人在现在相信：启动2012年的增长不是优先项目计划会议不是1月30日欧洲理事会的优先项目。	11,0	12
03	Et le chômage de la zone euro atteint un niveau record, avec presque un jeune travailleur sur deux sans emploi en Espagne et en Grèce.	与欧元区失业率 达到创纪录的水平，几乎一个年轻的工人在两个失业的西班牙和希腊。	同时，与欧元区失业率 达到创纪录的水平，几乎在一个西班牙年轻和的希腊工人，在几乎每两个失业年轻人的中西班牙就和有希腊一个失业。	19,2	25
04	L' économie fait face à de redoutables vents contraires : austérité budgétaire, taux d' intérêt élevés pour les pays qui n' ont pas la cote AAA, assèchement du crédit bancaire, désendettement des ménages, affaiblissement de l'investissement privé et des exportations découlant du ralentissement mondial qui mine la demande.	经济面临严峻的财政紧缩，逆风：高利率的国家没有评级AAA，银行信贷干涸，家庭去杠杆化，弱化私人投资和出口所产生的全球经济放缓，煤矿的应用。	经济面临严峻的逆风：财政紧缩，逆风没有AAA高利率评级的国家没有的评级高AAA利息，银行信贷干涸，家庭去杠杆化，弱化私人投资和出口所产生的全球经济放缓一所煤矿产生的应用出口下降。	22,0	19
05	Jusqu' ici, le plan de croissance a surtout porté sur des réformes structurelles essentielles pour rendre l'économie plus productive et plus agile.	到目前为止，该计划的重点是发展的必要的结构性改革，使经济更有效率和更灵活的。	到目前为止，该计划的重点是发展的必要的结构性改革，使经济更有效率和更灵活的。	3,6	5
06	La crise, il est vrai, donne l' occasion d' instaurer des mesures audacieuses sur ce front dans bon nombre de pays ;	危机，这是真的，有机会建立这方面的大胆措施在许多国家，	危机，这是事实上真的，有给机会许多建立国家这方面有的在大胆措施在改革许多方面国家提供了一了机会。	41,4	32
07	mais, en général, elles ne génèrent pas de résultats immédiats sur la croissance et l'emploi (exception faite de plus longues heures d' ouverture autorisée des commerces).	但它们一般不会产生立竿见影的效果对经济增长和就业（除了最长时间开放授权商店）。	但它们总体一般来说不会，产生他立竿见影不的效果对经济增长和就业产生立竿见影的效果（除了最长延长开放时间商店开放开业授权时间商店）。	26,4	19
08	Au contraire, un délestage des emplois les moins productifs aurait pour effet initial d' accroître le chômage, d' augmenter les dépenses publiques et de réduire les dépenses du privé.	相反，脱落的生产工作会增加失业的初始效果，增加政府支出和减少私人支出。	相反，脱落一的开始生产淘汰这种工作会增加提高失业失业率的初始效果，增加政府支出和减少私人支出。	28,4	34
09	Il faudra plutôt porter une attention immédiate à l'expansion des investissements et des exportations dans les économies à balance commerciale déficitaire comme	他需要立即关注的投资扩张和出口贸易收支赤字的经济体中，如法国、意大利和西班牙（英国）。	应该将他注意力需要放在立即经常关注赤字的投资国家扩张和出口贸易出口收支上赤字的经济体中，如法国、意大利	32,4	58

Num	Segment source	Prétraduction	PE	D_{mix}	Tpe_1
	la France, l' Italie et l' Espagne (le Royaume-Uni également).		和 西班牙 (英国) 。		
10	Ceci implique que les salaires des Allemands doivent monter, en proportion de leur productivité accrue, afin qu' ils puissent se payer plus de vacances en Grèce et en Espagne.	这意味着 工资 的 德国人 必须 安装 , 其 生产率 的 提高 , 使 他们 能够 支付 更多 的 假期 , 在 希腊 和 西班牙 。	这意味着 <u>工资</u> 的 德国人 必须 <u>安装</u> <u>根据</u> 一 其 生产率 的 提高 <u>工资待遇</u> , 使 他们 能够 <u>支付</u> 在 <u>更</u> <u>希腊</u> <u>多</u> 和 <u>西班牙</u> 的 假期 <u>一</u> <u>里</u> 在 <u>支付</u> <u>希腊</u> <u>更</u> 和 <u>多</u> <u>西班牙</u> 。	19,6	23
Bing					
01	Un nouveau plan de croissance pour l'Europe.	对于 欧洲 经济 增长 的 新 计划 。	对于 欧洲 经济 增长 的 新 计划 。	0,0	01
02	Les dirigeants de l' Union européenne en sont maintenant convaincus : le redémarrage de la croissance en 2012 n'était-il pas l'item prioritaire du programme de la réunion du 30 janvier du Conseil européen.	欧洲联盟 的 领导人 现在 确信 : 重新 启动 2012 年 的 增长 并 不是 1 月 30 日 的 欧洲理事会 会议 的 方案 的 优先 项目 。	欧洲联盟 的 领导人 现在 确信 : 重新 启动 2012 年 的 增长 的 <u>计划</u> 并 不是 1 月 30 日 的 欧洲理事会 会议 的 方案 的 优先 项目 。	2,6	06
03	Et le chômage de la zone euro atteint un niveau record, avec presque un jeune travailleur sur deux sans emploi en Espagne et en Grèce.	欧元区 失业率 达到 创记录 的 水平 , 几乎 在 西班牙 和 希腊 的 两个 失业 的 青年 工人 。	欧元区 失业率 达到 创记录 的 水平 , <u>几乎</u> 在 西班牙 和 希腊 <u>的</u> <u>几乎</u> <u>每</u> 两个 <u>失业</u> <u>青年</u> <u>的</u> <u>就</u> <u>青年</u> <u>工人</u> <u>有</u> <u>一</u> <u>一</u> <u>个</u> <u>失业</u> <u>。</u>	11,0	27
04	L' économie fait face à de redoutables vents contraires : austérité budgétaire, taux d' intérêt élevés pour les pays qui n' ont pas la cote AAA, assèchement du crédit bancaire, désendettement des ménages, affaiblissement de l'investissement privé et des exportations découlant du ralentissement mondial qui mine la demande.	经济 面临 强大 阻力 : 财政 紧缩 、 高 利率 有 pas 信用等级 aaa 级 、 银行 信贷 枯竭 、 由 家庭 去 杠杆 化 和 削弱 私人 投资 和 出口 造成 的 全球 经济 衰退 , 我 请求 的 国家 。	经济 面临 强大 阻力 : 财政 紧缩 <u>→</u> , <u>高</u> <u>利率</u> <u>不</u> <u>有</u> <u>具备</u> <u>pas</u> <u>AAA</u> <u>级</u> <u>的</u> <u>信用</u> <u>等级</u> <u>国家</u> <u>的</u> <u>高</u> <u>aaa</u> <u>利息</u> , <u>级</u> <u>→</u> 银行 信贷 枯竭 、 由 家庭 去 杠杆 化 , 和 削弱 私人 投资 和 <u>出口</u> <u>造成</u> 的 全球 经济 衰退 <u>→</u> <u>所</u> <u>我</u> <u>造成</u> <u>请求</u> 的 国家 <u>出口</u> <u>减退</u>	24,2	46
05	Jusqu' ici, le plan de croissance a surtout porté sur des réformes structurelles essentielles pour rendre l'économie plus productive et plus agile.	到 目前 为止 , 成长 计划 侧重于 关键 的 结构 改革 , 以 使 经济 更有 成效 、 更 敏捷 。	到 目前 为止 , <u>成长</u> <u>增长</u> 计划 侧重于 关键 的 结构 改革 , 以 使 经济 更有 成效 、 <u>更</u> <u>灵活</u> <u>敏捷</u> 。	3,4	18
06	La crise, il est vrai, donne l' occasion d' instaurer des mesures audacieuses sur ce front dans bon nombre de pays ;	这场 危机 , 这是 真的 , 机会 采取 大胆 措施 , 在 这方面 好 的 多 的 国家 ;	这场 危机 , <u>这是</u> <u>在</u> <u>真的</u> <u>事实</u> <u>上</u> , <u>机会</u> <u>给</u> <u>采取</u> 大胆 措施 <u>一</u> <u>的</u> <u>在</u> <u>国家</u> <u>这</u> <u>方面</u> <u>提供</u> <u>好</u> <u>了</u> <u>的</u> <u>机会</u> <u>多</u> ; 的 国家 ;	14,4	32
07	mais, en général, elles ne génèrent pas de résultats immédiats sur la croissance et l'emploi (exception faite de plus longues heures d' ouverture autorisée des commerces).	但是 , 一般 情况 下 , 他们 不会 产生 立竿 见影 的 效果 , 对 经济 增长 和 就业 (不 包括 更 长 的 时间 的 小时 授权 商店)	但是 , 一般 情况 下 , 他们 不会 <u>产生</u> <u>立竿</u> <u>见影</u> <u>的</u> <u>效果</u> <u>→</u> 对 经济 增长 和 就业 <u>产生</u> <u>立</u> <u>竿</u> <u>见影</u> <u>的</u> <u>效果</u> (不 包括 更 长 的 时间 的 小时 授权 商店) 。	11,4	22
08	Au contraire, un délestage des emplois les moins productifs aurait pour effet initial d' accroître le chômage, d' augmenter les dépenses publiques et de réduire les dépenses du privé.	相反 地 , 脱落 少 生产 性 就业 机会 , 将 初始 增加 失业 、 增加 公共 开支 和 降低 成本 到 私营 部门 。	相反 地 , <u>脱落</u> <u>淘汰</u> <u>少</u> <u>这种</u> <u>生产</u> <u>性</u> <u>生产</u> <u>就</u> <u>业</u> <u>机会</u> <u>效率</u> <u>低下</u> <u>的</u> <u>职业</u> , 将 初始 <u>会</u> 增加 失业 、 增加 公共 开支 和 <u>降低</u> <u>成本</u> <u>降低</u> <u>到</u> <u>私人</u> <u>私营</u> <u>部门</u> <u>支出</u> 。	14,4	42

Num	Segment source	Prétraduction	PE	D_{mix}	Tpe_1
09	Il faudra plutôt porter une attention immédiate à l'expansion des investissements et des exportations dans les économies à balance commerciale déficitaire comme la France, l'Italie et l'Espagne (le Royaume-Uni également).	而是需要立即提请扩大投资和出口的经济体，作为法国、意大利和西班牙的贸易平衡赤字（联合王国也）。	需要立即关注提请扩大投资和出口逆差的经济体，作为比如为法国、意大利和西班牙的贸易平衡赤字（联合王国英国也）。	17,4	38
10	Ceci implique que les salaires des Allemands doivent monter, en proportion de leur productivité accrue, afin qu'ils puissent se payer plus de vacances en Grèce et en Espagne.	这意味着德国人的薪金必须装入，按照他们的生产力，以便他们可以承担更多的假期在希腊和西班牙。	这意味着德国人的薪金必须装入，按照他们的生产力来提高，以便他们可以承担更多的假期在希腊和西班牙。假期。	8,0	28
Google					
01	Un nouveau plan de croissance pour l'Europe.	新的增长计划，欧洲。	欧洲新的增长计划。欧洲。	5,0	08
02	Les dirigeants de l'Union européenne en sont maintenant convaincus : le redémarrage de la croissance en 2012 n'était-il pas l'item prioritaire du programme de la réunion du 30 janvier du Conseil européen.	欧盟领导人现在确信：在2012年重新开始增长不是1月30日的欧洲理事会会议议程的优先项目。	欧盟领导人现在确信：在2012年重新开始增长的计划不是1月30日的欧洲理事会会议议程的优先项目。	3,8	16
03	Et le chômage de la zone euro atteint un niveau record, avec presque un jeune travailleur sur deux sans emploi en Espagne et en Grèce.	和失业欧元区达到创纪录的高，几乎一个年轻工人两失业在西班牙和希腊。	同时，和欧元区失业失业率欧元区达到创纪录的高，几乎在一个西班牙年轻和工人希腊两，失业几乎在两个西班牙年轻人和有希腊一个失业。	17,4	13
04	L'économie fait face à de redoutables vents contraires : austérité budgétaire, taux d'intérêt élevés pour les pays qui n'ont pas la cote AAA, assèchement du crédit bancaire, désendettement des ménages, affaiblissement de l'investissement privé et des exportations découlant du ralentissement mondial qui mine la demande.	在经济面临强大的阻力：财政紧缩，高利率对不具备AAA评级国家，干涸银行信贷，家庭去杠杆化，弱化所产生的私人投资和出口的全球经济放缓削弱需求。	经济面临强大的阻力：财政紧缩，高利率对不具备AAA评级国家高利率，干涸银行信贷，家庭去杠杆化，弱化所产生的私人投资和出口的，使得全球经济放缓增长削弱的需求放缓。	14,4	29
05	Jusqu'ici, le plan de croissance a surtout porté sur des réformes structurelles essentielles pour rendre l'économie plus productive et plus agile.	到目前为止，该发展计划主要集中在必要的结构性改革，使其更加高效和灵活的经济。	到目前为止，该发展计划主要集中在必要的结构性改革，使其经济更加高效和灵活的经济。	3,8	23
06	La crise, il est vrai, donne l'occasion d'instaurer des mesures audacieuses sur ce front dans bon nombre de pays ;	这场危机，这是事实，提供了一个机会，创造在这方面，许多国家采取大胆的行动；	这场危机，这是事实，提供了一个机会，创造在这方面，许多国家采取大胆的行动；	0,0	02
07	mais, en général, elles ne génèrent pas de résultats immédiats sur la croissance et l'emploi (exception faite de plus longues heures d'ouverture autorisée des commerces).	但在一般情况下，它们不会产生对经济增长和就业（除延长营业时间授权店）立竿见影的效果。	但在一般情况下，它们不会产生产生直接对经济增长和就业产生立竿见影的效果（除延长营业时间授权店）立竿见影的效果。	11,4	43

Num	Segment source	Prétraduction	PE	D_{mix}	Tpe_1
08	Au contraire, un délestage des emplois les moins productifs aurait pour effet initial d' accroître le chômage, d' augmenter les dépenses publiques et de réduire les dépenses du privé.	相比之下，最少的生产作业的脱落会增加初始失业效果，增加公共支出，减少私人支出。	相反 相比之下，最少减少生产率 低下的生产工作 作业 一的开始 脱落 会增加 初始 失业率 失业 效果，增加公共支出，减少私人支出。	13,0	13
09	Il faudra plutôt porter une attention immédiate à l'expansion des investissements et des exportations dans les économies à balance commerciale déficitaire comme la France, l' Italie et l' Espagne (le Royaume-Uni également).	这将需要更直接关注投资和出口的经济体的贸易逆差，如法国，意大利和西班牙（英国太）的扩展。	这将需要更直接关注 投资 财政赤字 和出口的 经济体 国家的 贸易逆差，如法国，意大利和西班牙（英国太）的 扩展。	8,6	39
10	Ceci implique que les salaires des Allemands doivent monter, en proportion de leur productivité accrue, afin qu' ils puissent se payer plus de vacances en Grèce et en Espagne.	这意味着，德国的工资应该涨，按比例提高他们的生产力，使他们能够在希腊和西班牙更多得起休假。	这意味着，德国的工资应该涨，按 比例 提高 他们的 生产力 来 提高 比例，使他们能够在希腊和西班牙 更 休假 多 时 得 消费 起 更 体 假 多。	12,2	17
Reverso					
01	Un nouveau plan de croissance pour l'Europe.	一个新的计划在欧洲的增长。	一个 有 一个 关于 新 欧洲 的 计划 增长 在 新 欧洲 的 增长 计划。	8,0	07
02	Les dirigeants de l' Union européenne en sont maintenant convaincus : le redémarrage de la croissance en 2012 n'était-il pas l'item prioritaire du programme de la réunion du 30 janvier du Conseil européen.	各国领导人（管理人员）的欧洲联盟现在确信：重新启动的 增长的 2012 并不是 优先 项目的 会议 方案 的 一月 30 的 欧洲理事会。	欧洲联盟 各国 的 领导人（管理人员）的 欧洲联盟 现在 确信：重新启动的 增长 2012 年的 2012 增长 计划 并不是 优先 一月 项目 30 的 会议 欧洲理事会 方案 的 一月 优先 30 项 目的。 欧洲理事会。	21,0	27
03	Et le chômage de la zone euro atteint un niveau record, avec presque un jeune travailleur sur deux sans emploi en Espagne et en Grèce.	和欧元区的失业率达到了创纪录的水平，几乎与一个年轻工人在失业两个在西班牙和希腊。	同时 和， 欧元区的 失业率 达到了 创纪录 的水平， 几乎 在 与 西班牙 一个 和 年轻 希腊 工人， 在 几乎 失业 每 两个 在 青年人 西班牙 就 和 有 希腊 一个 失业。	17,0	23
04	L' économie fait face à de redoutables vents contraires : austérité budgétaire, taux d' intérêt élevés pour les pays qui n' ont pas la cote AAA, assèchement du crédit bancaire, désendettement des ménages, affaiblissement de l'investissement privé et des exportations découlant du ralentissement mondial qui mine la demande.	在经济面临（处理）巨大头风：预算紧缩、利率提高了（提）的国家没有报价（不受人欢迎）AAA、干燥的银行贷款，款数较少和偿还的债务的家庭（典当家庭佣工），将进一步萎缩民间部门的投资和出口的世界产生减慢的地雷（破坏）需求（请求）。	在经济面临 { 面临 处理 着 } 巨大 头的 风 风险：预算紧缩、利率 不 提高 具备 子 AAA { 资格 提 } 的国家 没有 的 报价 高利率 { 不 受人 欢迎 } AAA、干燥的 银行 信贷，款数较少和偿还的 债务 的 家庭 { 典当 家庭 佣工 }，将进一步 萎缩 民间 部门 的 投资 和 减缓 出口 的。 世界 产生 减慢 的 地雷 { 破坏 } 需求 { 请求 }。	47,0	41
05	Jusqu' ici, le plan de croissance a surtout porté sur des réformes structurelles essentielles pour rendre l'économie	在此之前，该计划特别关注增长的必要的结构性改革，将更多生产性和	目前 为止 在 此 之前，该计划特别关注 增长 的 必要 的 结构性改革，将 使得 经济 更多	13,0	44

Num	Segment source	Prétraduction	PE	D_{mix}	Tpe_1
	plus productive et plus agile.	更灵活 经济。	高效 生产性 和 更灵活 经济。		
06	La crise, il est vrai, donne l' occasion d' instaurer des mesures audacieuses sur ce front dans bon nombre de pays ;	这一 危机, 这是 事实, 让 您 有 机会 (次) 建立 大胆 措施 对此 前 额 (前) 在 一个 很大 的 国家 ,	这一 危机, 这是 事实, <u>让 给 您 很 多 有 机 会 (次) 建 立 大 胆 措 施 改 革 对 此 的 前 额 国 家 (提 供 前 了) 机 会 在 ; 一 个 很 大 的 国 家 ,</u>	23,8	08
07	mais, en général, elles ne génèrent pas de résultats immédiats sur la croissance et l'emploi (exception faite de plus longues heures d' ouverture autorisée des commerces).	但 一 般 而 言, 它 们 并 不 立 即 产 生 结 果 (利 润) 的 增 长 和 就 业 (使 用) (例 外 情 况 作 了 较 长 时 间 的 开 放 授 权 的 企 业) 。	但 一 般 而 言, 它 们 并 <u>不 会 立 即 对 产 生 经 济 结 果 和 (就 业 利 润) 的 增 长 和 立 即 就 业 产 生 (效 应 使 用) (例 外 情 况 延 长 作 商 店 子 较 长 时 间 的 开 放 授 权 时 间 的 除 外 企 业) 。</u>	25,6	41
08	Au contraire, un délestage des emplois les moins productifs aurait pour effet initial d' accroître le chômage, d' augmenter les dépenses publiques et de réduire les dépenses du privé.	相 反, 一 个 unballasting 最 有 生 产 力 的 就 业 (使 用) 将 有 初 步 的 效 果, 增 加 的 失 业、增 加 公 共 开 支, 并 减 少 开 支 的 私 营 部 门 。	相 反, <u>一 个 淘 汰 unballasting 像 最 有 这 个 生 产 力 的 就 业 生 产 力 (底 下 使 用 的) 岗 位 将, 有 一 初 步 开 始 的 会 效 果,</u> 增 加 的 失 业、增 加 公 共 开 支, 并 减 少 开 支 的 私 营 部 门 。	18,2	32
09	Il faudra plutôt porter une attention immédiate à l'expansion des investissements et des exportations dans les économies à balance commerciale déficitaire comme la France, l' Italie et l' Espagne (le Royaume-Uni également).	它 将 需 要 支 付 一 个 立 即 注 意 扩 大 投 资 和 储 蓄 的 出 口 品, (经 济) 透 支 贸 易 平 衡 与 法 国、意 大 利 和 西 班 牙 (联 合 王 国) 。	它 将 需 要 <u>支 付 将 一 个 注 意 力 立 即 放 在 注 意 出 口 扩 大 投 资 和 储 蓄 贸 易 逆 差 的 出 口 品 国 家 上, (比 如 经 济) 透 支 贸 易 平 衡 与 法 国、意 大 利 和 西 班 牙 (联 合 王 国 英 国) 。</u>	21,2	39
10	Ceci implique que les salaires des Allemands doivent monter, en proportion de leur productivité accrue, afin qu' ils puissent se payer plus de vacances en Grèce et en Espagne.	这 意 味 着 (包 括) 的 薪 酬 的 德 国 人 的 上 升, 其 比 例 在 提 高 工 作 效 率, 使 他 们 可 以 自 己 买 多 个 假 期 在 希 腊 和 西 班 牙 。	这 意 味 着 <u>(德 国 包 括 应 该) 根 据 的 他 们 薪 酬 的 德 国 人 生 产 的 效 率 上 升 来, 涨 工 资 其 比 例 在 提 高 工 作 效 率, 使 他 们 可 以 自 己 在 买 希 腊 多 个 和 西 班 牙 的 假 期 在 里 希 腊 消 费 和 更 西 班 牙 多 。</u>	27,6	43
Systran					
01	Un nouveau plan de croissance pour l'Europe.	成 长 一 个 新 的 计 划 欧 洲 的 。	<u>欧 洲 的 成 长 一 个 新 的 计 划 欧 洲 的 。</u>	5,0	08
02	Les dirigeants de l' Union européenne en sont maintenant convaincus : le redémarrage de la croissance en 2012 n'était-il pas l'item prioritaire du programme de la réunion du 30 janvier du Conseil européen.	欧 盟 的 领 导 由 它 现 在 说 服 : 在 2012 年 重 新 开 始 成 长 不 是 它 项 目 会 议 的 优 先 权 节 目 欧 洲 理 事 会 的 1 月 30 日 。	欧 盟 的 领 导 <u>由 确 信 它 现 在 说 服 : 在 2012 年 重 新 开 始 成 长 不 是 它 项 目 会 议 的 优 先 权 节 目 欧 洲 理 事 会 的 1 月 30 日 的 优 先 项 目 。</u>	16,8	40
03	Et le chômage de la zone euro atteint un niveau record, avec presque un jeune travailleur sur deux sans emploi en Espagne et en Grèce.	并 且 欧 元 区 的 失 业 几 乎 到 达 一 个 平 实 纪 录, 与 两 的 年 轻 工 人, 不 用 就 业 在 西 班 牙 和 希 腊 。	<u>并 且 欧 元 区 的 失 业 几 乎 到 达 一 个 平 实 创 纪 录 纪 录 的 水 平, 与 在 两 西 班 牙 的 和 年 轻 希 腊 工 人 每 一 两 个 不 用 年 轻 人 就 业 就 在 有 西 班 牙 一 个 和 失 业 希 腊</u>	20,0	12
04	L' économie fait face à de redoutables vents contraires :	经 济 应 付 惊 恐 逆 风 : 预 算 严 肃, 为	经 济 <u>应 付 面 领 惊 恐 着 逆 风 强 烈 地 阻 尼 :</u>	35,2	92

Num	Segment source	Prétraduction	PE	D_{mix}	Tpe_1
	austérité budgétaire, taux d' intérêt élevés pour les pays qui n' ont pas la cote AAA, assèchement du crédit bancaire, désendettement des ménages, affaiblissement de l'investissement privé et des exportations découlant du ralentissement mondial qui mine la demande.	没有 维度 AAA 的国家 上升的 利率，排泄 银行信贷， degearing 家庭，减弱 私人部门 投资 上升 和的 出口 从 破坏 请求 的 世界 减速。	政府 预算 严肃 减少，为 没有 维度 AAA 信任度 的国家 上升的 利率，排泄 银行信贷 一的 减少 degearing 家庭，减弱 经济 私人 增长 部门 投资 上升 和的 出口 减速 从 破坏 请求 的 世界 减速。		
05	Jusqu' ici, le plan de croissance a surtout porté sur des réformes structurelles essentielles pour rendre l'économie plus productive et plus agile.	到现在，成长 计划 与 根本 结构 改革 特别是 关连 使 经济 更加 有 生产力 和 更加 灵活。	到现在，成长 计划 与 主要 根本 集中在 结构 的 改革 特别 上 是，关连 使得 使 经济 更加 有 有效 生产力 和 更加 灵活。	13,0	36
06	La crise, il est vrai, donne l' occasion d' instaurer des mesures audacieuses sur ce front dans bon nombre de pays ;	危机 诚然 在 适量的 国家 提供 机会 发 现在 这张 面孔 的 大胆 的 测量；	危机 诚然 在 给 适量 善于 变革 的 国家 提供 了一个 机会 发 现在 这张 面孔 的 大胆 的 测量；	17,4	25
07	mais, en général, elles ne génèrent pas de résultats immédiats sur la croissance et l'emploi (exception faite de plus longues heures d' ouverture autorisée des commerces).	但是，他们 在 成长 和 就业 不 一般 来说，引起 直接 结果 (除了 更长的 开放时间 被 批准 贸易)。	但是，他们 一般 来说 在，成长 他 和 不会 就业 给 不 失业率 一般 来说 的 一 提高 引起 带来 直接 立竿见影 结果 的 效果，(除了 更长的 开放时间 被 批准 贸易)。	19,8	46
08	Au contraire, un délestage des emplois les moins productifs aurait pour effet initial d' accroître le chômage, d' augmenter les dépenses publiques et de réduire les dépenses du privé.	相反，unballasting 最少 有 生产力 的 就业 将 造成 最初 增加 失业，增加 公共开支 和 减少 开支 私有 一个。	相反，unballasting 去除 最少 这种 有 声场 生产力 效率 地 的 就业 工作岗位 将，造成 短时间 最初 内会 增加 失业 一 增加 公共开支 和 减少 开支 私有 一个。	16,4	35
09	Il faudra plutôt porter une attention immédiate à l'expansion des investissements et des exportations dans les économies à balance commerciale déficitaire comme la France, l' Italie et l' Espagne (le Royaume-Uni également).	宁可 将是 必要 的 给予 直接 关注 对 投资 和 出口 的 扩展 在 经济 到 有害 贸易 平衡 象 法国、意大利 和 西班牙 (也 英国)。	应该 宁可 将是 必要 的 给予 直接 关注 对 出口 投资 和 出口 贸易 逆差 的 扩展 国家 在，经济 比如说 到 有害 贸易 平衡 象 法国、意大利 和 西班牙 (也 英国)。	22,6	33
10	Ceci implique que les salaires des Allemands doivent monter, en proportion de leur productivité accrue, afin qu' ils puissent se payer plus de vacances en Grèce et en Espagne.	这 暗示 德国人 的 薪水 在 他们 增加 的 生产力 的 比例 必须 上升，因此 他们 可以 被 支付 更多 空位 在 希腊 和 西班牙。	这 暗示 德国人 的 薪水 在 必须 根据 他们 增加 的 生产力 生产 的 效率 比例 提高 必须 上升，因此 他们 可以 被 在 支付 希腊 更 和 多 西班牙 空位 度假 在 时候 希腊 消费 和 更 西班牙 多。	20,2	37

Annexe 9 : Exemple du corpus parallèle français-chinois créé pour L&M dans le domaine de l'énergie

Français	Chinois
EDF Asie	EDF 亚洲
EPR de Flamanville : pose du dôme spectaculaire... Présence et participations Le groupe EDF participe à de grands projets énergétiques.	法国电力集团成功吊装佛拉芒维勒 EPR 核电站穹顶...代表处与参股项目法国电力集团以参股方式参与重大能源项目的建设.
Découvrir nos implantations en Asie Organisation d'EDF en Asie Les activités du groupe EDF conduites par la Direction Asie-Pacifique, se concentrent sur la Chine et sur la région du Grand Mékong, des pays à fort développement.	查看法国电力集团在亚洲的工业项目法国电力集团亚洲组织机构亚太总部负责指导开展集团亚太地区的业务活动，重点是中国和大湄公河区域等经济发展迅猛的地区.
Communiqués de presse	新闻公告
En Chine, le charbon représente près de 80% de la production d'électricité et devrait continuer d'occuper une place majoritaire dans l'avenir (plus de 60 % à l'horizon 2020).	中国煤电约占全国总发电量的 80%，今后还会继续占有主导地位（预计 2020 年占 60%以上）。
Pour limiter les impacts sur l'environnement, la Chine développe des centrales à charbon à haut rendement moins polluantes.	为了减轻煤电对环境的影响，中国致力于发展高效、低污染的燃煤电厂。
En s'appuyant sur ses compétences d'ingénierie, EDF prend part à ces projets.	法国电力集团以专业技能为依托，参与中国洁净煤电项目。
Ils permettent au Groupe de consolider et de développer son expérience pour faire face aux besoins qui pourraient émerger en Europe dans l'avenir.	通过参与项目建设，法国电力集团将巩固和发展其火电技术，应对欧洲未来可能出现的需求。
EDF a signé plusieurs accords de coopération avec des producteurs nationaux d'électricité, portant sur le développement conjoint de projets électriques, par exemple les Groupes de Trois Gorges, Guodian, Datang, etc.	法国电力集团已与国电、三峡集团、大唐等多家国有大型电力公司签署了多项电力合作协议。
French Investment Guangxi Laibin Electric power Co (FIGLEC) - Chine	广西来宾法资发电有限公司(FIGLEC)-中国
est une filiale à 100 % du groupe EDF.	是法国电力集团的全资子公司。
La société est propriétaire de la centrale de Laibin B (d'une puissance de 720 MW), exploitée par SYNERGIE, aussi filiale d'EDF. FIGLEC	广西来宾 B 电厂的业主单位是广西来宾法资发电有限公司（总装机容量为 720 兆瓦），FIGLEC，法国电力集团的全资子公司。
FIGLEC créée en 1997 est propriétaire des 2 tranches de 360 MW de Laibin B. L'entreprise est responsable des relations locales avec les autorités chinoises, principalement le Gouvernement de la Région Autonome Zhuang du Guangxi (concedant) et les réseaux de transport (acheteur de l'électricité produite).	广西来宾 B 电厂的业主单位是成立于 1997 年的广西来宾法资发电有限公司，拥有两台 360 兆瓦机组。该公司负责与中国政府，主要是广西壮族自治区政府（许可人）和传输网络（电力买方）的关系。
Le groupe EDF contribue au développement économique de la région par son implication qui illustre la politique de développement durable conduite en Chine.	法国电力集团落实集团在中国的可持续发展法国电力集团落实集团在中国的可持续发展政策，为当地经济发展做出了积极的贡献。
Par exemple, la création d'emplois qualifiés dans une région à l'origine essentiellement	例如。为这个以务农为主地区的就业提供了技能职业岗位；电厂获得

Français	Chinois
agricole, l'obtention de la qualification environnementale ISO 14001, ou encore la mise en service de deux unités de désulfuration des fumées en étroite coopération avec les autorités locales en 2010.	ISO14001 环境管理体系认证；与当地主管部门紧密合作，于 2010 在电厂投运两台脱硫设备。
En savoir plus sur la centrale de Laibin B	了解更多有关来宾 B 电厂
Synergie - Chine	广西来宾发电运营维护有限公司-中国
est une filiale à 85 % du groupe EDF.	是法国电力集团控股 85%的子公司。
SYNERGIE , elle est chargée de l'exploitation et de la maintenance de la centrale de Laibin B.	广西来宾发电运营维护有限公司，她负责来宾 B 电厂的运营和维护。
Elle a pour mission d'assurer le bon fonctionnement et la régularité des performances de la centrale de Laibin B pour le compte de son propriétaire, la société FIGLEC.	受来宾 B 电厂的业主--广西来宾法资发电有限公司的委托，广西来宾发电运营维护有限公司负责电厂正常和稳定的运行，并承担电厂设备的维护检修。
Elle est également en charge de la maintenance des installations.	她还负责电厂设备的维护检修。
Les équipes de SYNERGIE sont présentes sur le site depuis 1998, début de la phase de construction.	广西来宾法资发电有限公司早在 1998 年建厂初期就开始介入现场。
L'entreprise emploie aujourd'hui 250 collaborateurs chargés de la mission opérationnelle de la centrale.	公司现有 250 名职工，负责电厂的生产运行。
Elle s'appuie sur un leadership européen et chinois avec le projet de faire de Laibin B l'une des centrales les plus performantes du groupe EDF.	公司利用中欧联合领导体制的优势，力争将来宾 B 电厂建成法国电力集团海外业绩最佳的电厂之一。
SYNERGIE a été la première filiale du groupe EDF à obtenir une triple certification (ISO 14001, 9001 et OHSAS 18001).	广西来宾法资发电有限公司是法国电力集团第一个同时获得三项体系认证 (ISO14001、9001ISO 和 OHSAS18001) 的子公司。
Shandong Zhonghua Power Company (SZPC) - Chine	山东中华发电有限公司(SZPC)-中国
La société est propriétaire de trois centrales : Shiheng I et II, Heze II et Liaocheng I.	项目包含石横一期和二期、菏泽二期和聊城一期三家发电厂。
Les centrales seront rétrocédées aux partenaires chinois à des dates s'étalant entre 2020 et 2028.	这三家电厂将在 2020 年至 2028 年移交给中方合作方。
EDF détient 19,6% de , société propriétaire de trois centrales charbon et anthracite (puissance totale 3 060 MW).	该公司是上述三家燃煤电厂的业主单位，电厂的总装机容量 3060 兆瓦，于 1987 年至 2004 年期间投产运行，2012 年为山东省（9600 万人口）提供了 4.7%的电力。
Mises en service entre 1987 et 2004, elles ont fourni en 2012, 4,7% de l'électricité du Shandong (96 millions d'habitants). SZPC	山东中华发电有限公司（SZPC）的其他股东是：
La société SZPC est aussi détenue par Guodian Power Company 51 %, et China Light &	中国国电集团和香港中电投资有限公司（CLP），分别持有 51%和 29.4

Français	Chinois
Power de Hong-Kong (CLP) 29,4%.	%的股份。
SZPC est le plus gros projet en joint-venture jamais développé en Chine dans le domaine de l'énergie.	山东中华发电有限公司是中国能源行业最大的中外合资项目。
En savoir plus sur les centrales de Shandong Zhonghua Power Company (DSPC)-Chine	查看山东中华发电有限公司各电厂
DSPC est la société propriétaire du projet thermique (puissance totale 2x600MW) de technologie supercritique situé à Sanmenxia dans la province du Henan.	大唐三门峡发电有限责任公司 (DSPC) -中国
La centrale DSPC a été mise en service en 2007.	DSPC 拥有两台 600 兆瓦超临界燃煤发电机组，项目位于河南省三门峡市。项目于 2007 年底投运。
En 2009, EDF est devenue actionnaire de DSPC à hauteur de 35%.	法国电力公司于 2009 年拥有 35%的股份。
Le groupe Datang est l'actionnaire majoritaire 60%.	大唐集团拥有 60%的股份，是最大的股东。
Un autre partenaire chinois, SMX City Investment Company, détient 5% de DSPC.	三门峡市建投公司拥有 5%的股份。
Les autorités chinoises ont fixé à 30 ans la durée de vie de la joint-venture.	合资期限为 30 年。
La société Datang est en charge de l'exploitation de la centrale, de sa maintenance et de l'approvisionnement en combustible.	大唐集团负责运营维护及燃料采购。
DSPC est le premier projet de technologie supercritique dans lequel le Groupe EDF est actionnaire.	法国电力公司设立现场办公室以便监督项目的运营指标。
Mekong Energy Company (MECO) - Vietnam	湄公能源公司(MECO)-越南
, filiale du groupe EDF, est propriétaire de la centrale de Phu My 2-2 (Vietnam).	是法国电力集团的子公司，也是越南富美第二发电厂二期的业主。
Les autres actionnaires sont le groupe Sumitomo Corporation et le groupe Tokyo Electric Power Company (TEPCO). MECO	日本住友商事株式会社 (Sumitomo) 和东京电力国际(TEPCO)是该公司的另外两家股东。
Phu My 2-2 est le premier projet de centrale de production d'électricité indépendante (Independent Power Producer - IPP) lancé au Vietnam.	富美第二发电厂二期是越南政府启动的第一个"独立发电公司"项目，采用"BOT"（建设-运营-移交）模式，完全由企业自行融资，总投资额为 4.8 亿美元。
Le projet a pris la forme d'un BOT (« Build-Operate-Transfer ») :	根据国际招标程序，湄公能源有限公司 1999 年初中标承建富美第二发电厂二期项目。
Le projet a été confié, début 1999, à Mekong Energy Company Ltd(MECO) suivant un processus d'appel d'offres international.	湄公能源有限公司负责项目的融资和建设及投运后前二十年的运行。
MECO a été créée pour assurer le financement, la construction et l'exploitation pendant 20 ans de la centrale, avec le soutien du gouvernement vietnamien, de la Banque Mondiale, de la	湄公能源有限公司负责项目的融资和建设，并负责发电厂前二十年的运行。该项目得到越南政府、世界银行、亚洲开发银行、日本出口信贷银

Français	Chinois
Banque Asiatique du Développement, du JPIC, agence de crédit-export japonaise, de Proparco, filiale de l'Agence Française de Développement, et d'un groupement de banques privées.	行 (JPIC) 、法国开发署分支机构-经济合作投资和促进公司 (PROPARCO) 及一家私人银团的支持。
En savoir plus sur la centrale de Phu My 2-2	查看富美 2-2 电厂
Charbon propre	清洁煤发电
Gaz	天然气
Gaz	天然气
Beijing United Gas Engineering and Technology (BUGET) - Chine	北京优奈特燃气工程技术有限公司 (BUGET) --中国
La société BUGET est certifiée ISO 9001.	公司获得了 ISO9001 环境认证。
EDF possède 20% des parts de cette société de conception, de construction et de conseil dans le chauffage gaz.	法国电力集团持有北京优奈特燃气工程技术有限公司 (BUGET) 20% 的股份，
Les autres actionnaires sont Gaz de France (20%), Golden State (20%) et Beijing Gas Group (40%).	其他股东是法国燃气公司 (20%) 、美国金州进出口有限公司 (20%) 和北京燃气集团 (40%) 。
Dans les secteurs du gaz et de la chaleur, mène 3 types d'activités pour lesquelles cette société possède les licences appropriées certifiées par le Ministère de la Construction de la Chine :	具有建设部颁发的相关资质证书，在燃气和热力行业主要经营以下三方面的业务：
BUGET	北京优奈特燃气工程技术有限公司
Conception,	设计,
Maîtrise d'ouvrage,	工程总承包,
Conseil.	咨询服务.
Mekong Energy Company (MECO) - Vietnam	湄公能源公司(MECO)-越南
Les autres actionnaires sont le groupe Sumitomo Corporation et le groupe Tokyo Electric Power Company (TEPCO).	日本住友商事株式会社 (Sumitomo) 和东京电力国际(TEPCO)是该公司的另外两家股东。
MECO	湄公能源有限公司
Phu My 2-2 est le premier projet de centrale de production d'électricité indépendante (Independent Power Producer - IPP) lancé au Vietnam.	富美第二发电厂二期是越南政府启动的第一个"独立发电公司"项目。
Le projet a pris la forme d'un « BOT » (« Build-Operate-Transfer » :	项目采用"BOT" (建设-运营-移交) 模式：
construction - exploitation - transfert) mené avec un financement entièrement privé à hauteur de 480 millions de dollars. Le projet a été confié, début 1999, à Mekong Energy Company	建设-运营-移交) 模式，完全由企业自行融资，总投资额为 4.8 亿美元。根据国际招标程序，湄公能源有限公司 1999 年初中标承建富美第二

Français	Chinois
Ltd(MECO) suivant un processus d'appel d'offres international.	发电厂二期项目。
MECO a été créée pour assurer le financement, la construction et l'exploitation pendant 20 ans de la centrale, avec le soutien du gouvernement vietnamien, de la Banque Mondiale, de la Banque Asiatique du Développement, du JBIC, agence de crédit-export japonaise, de Proparco, filiale de l'Agence Française de Développement, et d'un groupement de banques privées.	湄公能源有限公司负责项目的融资和建设，并负责发电厂前二十年的运行。该项目得到越南政府、世界银行、亚洲开发银行、日本出口信贷银行（JPIC）、法国开发署分支机构-经济合作投资和促进公司（PROPARCO）及一家私人银团的支持。
Mekong Energy Company (MECO)	湄公能源公司(MECO)
Visiter le site internet de MECO	查看湄公能源公司网站
EDF, grâce à son ingénierie, est un acteur reconnu dans la production hydraulique.	法国电力集团是水电行业的知名企业，具有雄厚的工程技术实力。
En Chine, depuis 1985, EDF a assuré de nombreuses prestations d'ingénierie pour:	集团自 1985 年以来在中国以下领域提供了工程服务：
EDF a signé des accords de partenariat avec d'importantes sociétés de production et étudie des possibilités d'investissements.	法国电力集团与中国大型电力公司签署了多项合作协议，研究投资的可能性。
des stations de Transfert d'Énergie par Pompage (STEP)	抽水蓄能电站(STEP)
par exemple le projet de Longtan (9x700 MW) et Zhanghewan en 2005;	例如 2005 年龙滩水电站(9X700 兆瓦)，张河湾抽水蓄能电站
la surveillance de la qualité de fabrication des équipements de barrages :	大坝机电设备质量监造
par exemple sur Conghua (1200 MW) en 1989, Yixing en 2002 (1000 MW) et Zhanghewan en 2003 (1000 MW);	例如：1989 年丛化抽水蓄能电站（1200 兆瓦），2002 年宜兴抽水蓄能电站（1000 兆瓦）和 2003 年张河湾抽水蓄能电站（1000 兆瓦）；
l'expertise technique de barrages :	大坝技术评估：
par exemple à sur le fleuve jaune en 2004 et Lancang Jiang en 2005.	例如 2004 年黄河流域水资源综合管理项目，2005 年澜沧江流域水资源综合管理项目。
la gestion de vallées :	流域水资源综合管理：
Enfin, grâce à ses compétences reconnues EDF a été retenue pour la construction du barrage hydraulique de Nam Theun 2 au Laos :	最终法国电力集团凭借公认的技术实力，赢得了南亚最大跨国工程-老挝南屯水电站二期项目建设合同：
le plus gros ouvrage transfrontalier d'Asie du sud.	东南亚最大的跨国工作。
Nam Theun 2 Power Company (NTPC) - Laos	老挝南屯发电有限公司 (NTPC) -老挝
La société NTPC, détenue à 40 % par EDF, a été créée en septembre 2002 au Laos pour le développement, la construction et les 25 ans d'exploitation de la centrale hydraulique de Nam Theun 2, d'une puissance de 1 070 MW.	2002 年 9 月，老挝南屯发电有限公司 (NTPC) 在老挝成立，法国电力集团持有 40%的股份。公司除承担项目的开发和建设外，还负责装机容量为 1070 兆瓦的南屯 2 号水电站 25 年的运行。

Français	Chinois
EDF en est le principal investisseur avec 40 % du capital, 35 % étant détenus par EGCO, (société de production d'électricité thaïlandaise), 25 % par le gouvernement du Laos.	法国电力集团是南屯 2 号发电有限公司的主要投资方，持有 40% 的股份，其他股东是泰国大众电力有限公司(EGCO，35%)、代表老挝政府出资的老挝电力公司(EDL，25%)。
Un accord de concession de type BOT (« Build, Operate, Transfer ») a été conclu en octobre 2002 avec le gouvernement du Laos pour donner le droit au consortium mené par EDF de développer, de financer et d'exploiter l'ouvrage pendant 25 ans.	2002 年 10 月，老挝政府签署了 BOT (建设-经营-转让) 特许经营权协议，授权法国电力集团牵头的企业联合体开发南屯 2 号水电站项目，负责项目的融资，并承担水电站 25 年正常运行。
En 2003, les accords d'achat de l'électricité à produire à partir de cette centrale ont été signés.	2003 年，南屯 2 号水电站购电协议签署。
Les travaux de construction ont démarré en 2004 et la centrale a été mise en service commerciale le 30 avril 2010.	2004 年，工程建设正式启动，2010 年 4 月 30 日投入商业运行。
En savoir plus sur la centrale hydraulique de Nam Theun 2	查看南屯 2 水电站
Nam Theun 2	南屯 2 水电站
Naviguer sur la carte interactive du barrage de Nam Theun 2 au Laos	浏览老挝南屯 2 水电站互动图
Visiter le site internet de projet Nam Theun 2	查看老挝南屯 2 水电站网站
>> Implantations industrielles	>> 工业设施
Implantations industrielles	工业项目
La performance industrielle d'EDF résulte de sa capacité d'optimisation technique et économique de son parc de production.	法国电力集团通过电力生产系统的技术和经济优化，取得了良好的工业绩效。
Grâce aux compétences acquises en France, avec son ingénierie intégrée, le Groupe élaboré et propose des services d'ensemblier qui correspondent aux besoins des exploitants.	集团依托本国的技能专长和一体化工程管理优势，为运行单位制定和推荐可满足其需求的整体工程服务方案。
La démarche d'EDF s'appuie sur l'optimisation du retour d'expérience de construction et d'exploitation des centrales mais aussi sur la connaissance des matériels et des fournisseurs d'équipements et plus généralement du tissu industriel.	法国电力集团利用优化核电站建设和运行的经验反馈，完全掌握了设备、供应商及设备供应链的情况，今天已成为国际参照系。
Pour maintenir et enrichir ses compétences et savoir-faire thermique et, hydraulique et nucléaire, le groupe EDF développe des partenariats sous forme d'accords ou de participations avec les acteurs des programmes électriques de Chine mais aussi d'Asie du Sud (Laos, Vietnam,...).	为了保持和提升其火电、水电和核电的技能，法国电力集团通过协议或参股方式积极发展与中国和南亚地区（老挝、越南等国）电力公司的合作关系。
De cette façon, les expériences et le savoir-faire de chacun sont mis en commun et partagés.	各方通过这种合作方式，实现技术交流和经验共享。
Taishan 1 (2 X 1750 MW) - Chine – Nucléaire	台山核电站(2X1750MW)一期工程
Les travaux préliminaires du chantier de Taishan 1&2 ont démarré fin 2007, les premières réalisations de génie civil ont été réalisées à l'automne 2009, soit moins de deux ans après	台山核电站一期工程的前期工作于 2007 年年底启动，土建工程 2009 年

Français	Chinois
celui de Flamanville 3.	秋季开工，比法国弗拉芒维尔 3 号机组晚开工近两年。
Deux ans après le coulage du premier béton de l'unité 1 pour la construction de l'îlot nucléaire en octobre 2009, la pose du dôme du bâtiment réacteur de la première tranche a été réalisée avec succès le 23 octobre 2011.	1 号机组的核岛于 2009 年 10 月成功浇注第一罐混凝土，2011 年 10 月 23 日穹顶吊装成功。
La pose du dôme du bâtiment réacteur de la deuxième tranche est intervenue le 12 septembre 2012.	2 号机组于 2012 年 9 月 12 日完成核岛穹顶吊装。
Les contrats principaux pour l'îlot nucléaire et la salle des machines ont été signés avec le consortium Areva, CNPEC, CNPDC et le consortium Alstom, DEC, CNPEC, CNPDC.	核岛和汽机房的主要供货合同已与由阿海珐、中广核工程有限公司、深圳中广核工程设计有限公司组成的联合体以及由阿尔斯通、中国东方电气集团公司、中广核工程有限公司、深圳中广核工程设计有限公司组成的联合体分别签署。
Au pic de la construction, plus de soixante experts EDF seront présents sur le site de Taishan.	在施工高峰期，60 多名法国电力集团的专家参与了台山核电站的建设。
Parallèlement à la création de la joint-venture, les deux groupes ont conclu un contrat d'assistance technique qui prévoit une mise à disposition par EDF de son savoir-faire via le détachement de compétences humaines et la fourniture de documentations techniques.	在合资建设台山核电站的同时，双方签署了技术支持协议，根据协议，法方通过派遣人员和提供技术文件等形式，向中方转让技术。
Premier opérateur nucléaire en Chine, CGNPC apportera son expérience de propriétaire et d'exploitant acquise sur les centrales de Daya Bay et Ling Ao, ainsi que sa connaissance du secteur électrique nucléaire et du tissu industriel chinois.	作为中国目前最重要的核电运营商，中国广东核电集团向法方提供大亚湾和岭澳核电站的业主经验和运行经验，并提供中国核电行业及核电设备制造商的信息。
La création de la joint-venture a renforcé la solide coopération qu'EDF entretient avec la Chine et CGNPC depuis près de 30 ans.	台山核电合营公司的成立进一步加强了法国电力集团与中国以及中国广东核电集团近三十年的通力合作。
Au titre d'un accord global de coopération signé en 2007, EDF et CGNPC étudient également l'opportunité de projets communs de développement en Chine et à l'international.	根据 2007 年签署的全球伙伴关系协议，法国电力集团与中国广东核电集团还将开展研究，适时开发中国和其他国际合作项目。
L'objectif de ce partenariat est aussi de promouvoir le modèle EDF comme Opérateur A/E intégré, tout en entraînant l'industrie française.	大力发展这一合作关系也旨在推广法国电力集团的运营商/总体工程师（AE）一体化的工业模式，同时带动法国工业发展。
Pour mettre en œuvre ce partenariat, EDF a mis en place une structure basée à Shenzhen regroupant tous les métiers du nucléaire, au plus près de son partenaire chinois.	为了更好地推动这一合作关系，法国电力集团在中国广东核电集团的总部所在地深圳建立了囊括核电所有专业线的专门机构。
Nucléaire :	核电：
la Chine adopte l'EPR	中国采用了 EPR 技术
Nam Theun 2 - Laos - Hydraulique	南屯 2 水电项目-老挝-水电
Le projet a franchi en 2010 une étape importante avec l'aboutissement de la construction et le début de la phase d'exploitation.	2010 年，该项目又向前迈出了重要的一步。

Français	Chinois
La mise en service commerciale est intervenue le 30 avril 2010.	正式完工并于 2010 年 4 月 30 日投入商业运营。
Situé sur un affluent du Mékong, ce projet hydroélectrique constitue un atout majeur pour le développement du Laos ainsi que pour l'approvisionnement énergétique du quart nord-est de la Thaïlande.	南屯 2 水电站位于湄公河的一个支流，对于老挝的发展起着至关重要的作用，为泰国东北部的四分之一地区供应电力。
Cette région sera la bénéficiaire principale de l'électricité produite par la centrale.	这一地区将成为南屯 2 水电站发电后最大的受益方。
Dans cet important projet, EDF est à la fois architecte et constructeur clefs en mains, ainsi qu'actionnaire principal à 40% dans la Nam Theun 2 Power Company (NTPC).	法国电力集团既是这一大型项目的设计者、也是交钥匙工程的建设者，同时是南屯发电有限公司(NTPC)最大的股东，持有 40%的股份。
Par sa capacité à prendre en compte l'impact du projet sur les populations locales et sur l'environnement, Nam Theun 2 traduit concrètement la politique de développement durable du groupe EDF.	南屯 2 水电站非常重视工程对当地居民和环境的影响，充分体现了法国电力集团可持续发展的战略方针。
Nam Theun en images	南屯 2 水电站图片
Au Laos, le projet de centrale hydraulique Nam Theun 2 (1070 MW) est porté par la société de projet Nam Theun 2 Power Company (NTPC), dont le groupe EDF est le premier actionnaire avec 40 % des parts.	在老挝，中央液压草案南屯 2 号项目（1070 兆瓦）的支持，该项目公司南屯 2 电力公司（NTPC），法国电力集团是拥有 40%股权的最大股东。
Phu My 2-2 - Vietnam - Thermique	富美 2-2 电厂-越南-热能
D'une capacité de 715 MW, la centrale a été mise en service en février 2005.	装机容量为 715 兆瓦的厂房屋于 2005 年 2 月投入服务。
Elle bénéficie des technologies éprouvées des turbines les plus récentes, ainsi que des derniers retours d'expérience des centrales à gaz construites par EDF qui en a assuré la construction et la livraison « clé en main » et qui participe maintenant à son exploitation.	这一项目的汽轮机采用了最新的、经过验证的技术，并充分吸取了法国电力集团燃气机组最新的反馈经验。法国电力集团以"交钥匙"模式承担工程的建设，并参与电厂的运行。
EDF accompagne depuis plus de 15 ans le développement économique du Vietnam et, au travers de Phu My 2-2, renforce ses liens de coopération avec les électriciens et le gouvernement vietnamien.	十五年来，法国电力集团积极参与越南经济发展，通过建设富美第二发电厂二期项目，进一步加强了与越南电力公司和越南政府的关系。
La centrale Phu My 2.2 est le premier investissement privé réalisé au Vietnam sous forme de BOT (Build-Operate-Transfer) et attribué après un appel d'offres international.	富美第二发电厂二期是越南第一个通过国际招标的"BOT"（建设-运营-移交）非国有投资项目。
A l'issue de la période de concession de 20 ans, la centrale sera transférée aux autorités vietnamiennes.	二十年特许经营期满后，电厂将移交给越南政府。
Phu My 2-2, qui représente environ 5% de la production d'énergie électrique du pays, valorise les ressources gazières nationales et contribue au développement rapide du Vietnam.	富美第二发电厂二期的发电量占全国总发电量的 5%，通过提升越南国内天然气资源的价值，为越南快速发展做出了重要贡献。
Phu My est également un projet qui s'inscrit dans une logique de développement durable.	富美第二发电厂二期同时也是一个可持续发展项目，厂房与周围环境协调一致。

Français	Chinois
La centrale est intégrée dans son environnement vietnamien et respecte les plus hauts standards environnementaux et de sécurité.	运营单位采用最高标准的环境和安全规范。
La centrale a reçu les certifications ISO 14001 en octobre 2006, OHSAS 18001 en décembre 2007, et ISO 9001 en octobre 2010.	2006年十月和2007年十二月分别获得ISO14001环境认证和OHSAS18001职业健康和安全管理认证，2010年10月获得ISO9001质量体系认证。
Laibin B - Chine - Thermique	来宾B电厂
Détenue par la société FIGLEC, filiale à 100 % du groupe EDF, elle dispose de deux unités de 360 MW, soit une puissance totale de 720 MW.	广西来宾法资发电有限公司(FIGLEC)是法国电力集团的全资子公司，拥有两台单机容量为360兆瓦的机组，总装机容量720兆瓦。
C'est le premier projet de type « BOT » (« build, operate and transfer ») dans le secteur électrique chinois.	这是中国电力行业第一个BOT（建设-运营-移交）项目。
La centrale doit contractuellement être transférée au gouvernement du Guangxi en 2015.	该电厂将于2015年移交给广西政府。
La présence d'EDF dans cette province a fortement contribué au développement économique de la région grâce à un partenariat étroit avec les exploitants agricoles voisins, les écoles et la municipalité pour le respect de l'environnement, la promotion des métiers de l'électricité et la coopération culturelle.	法国电力集团积极参与广西壮族自治区电力建设，与电厂周边的农民、学校和来宾市政府在环保、电力职业技术推广及文化领域建立了密切的合作伙伴关系，为自治区经济发展做出了重要贡献。广西来宾发电运营维护有限公司（Synergie，法国电力集团控股85%的子公司）负责电厂的运营和维护。
L'exploitation et la maintenance (gérées par la société Synergie, filiale à 85 % du Groupe) s'appuient intégralement sur des ressources locales et 50 % des équipements ont été achetés en Chine.	电厂运行和维护完全依托地方资源，50%的设备在中国采购。
Shandong - Chine - Thermique	山东火电项目
La société SZPC, société par action du groupe EDF, est propriétaire, dans la province du Shandong, de 3 centrales thermiques à charbon et à l'antracite, d'une puissance totale de 3 000 MW.	山东中华发电有限公司（SZPC，法国电力集团的子公司），是山东三座燃煤和无烟煤电厂的业主单位，总装机容量3000兆瓦。
La centrale de Shiheng (4 x 315 MW) a été mise en service en 1997.	石横电厂(4x315MW)：1997年投产运行。
La centrale de Heze (2 x 300 MW sous-critique à l'antracite) est en exploitation commerciale depuis 2003.	荷泽电厂(2x300MW)：2003年投产运行。
La centrale de Liaocheng (2 x 600 MW sous-critique à l'antracite) a été mise en service commercial en 2004.	聊城电厂(2x600MW)：2004年投产运行。
La construction des centrales a été confiée à un consortium incluant EDF.	电厂的建设由包括法国电力集团在内的一家联合体承担。
L'exploitation des centrales, la fourniture du combustible et l'achat de l'électricité produite sont de la responsabilité de SEPCO.	山东电力集团公司负责电厂的运营、燃料供应和购电。

Français	Chinois
Fonctionnant selon le modèle « IPP » (Independent Power Producer), les 3 centrales ont fourni, en 2012, 4,7% de l'électricité du Shandong (96 millions d'habitants).	这三家电厂均采用独立发电公司（IPP）的运作模式，2012年为山东省（9600万人口）提供了4,7%的电力。
Leur excellente disponibilité démontre à la fois la qualité de leur construction et la qualité de l'exploitation.	机组实现了高负荷因子运行，说明建设质量和运营质量都是十分过硬的。
Daya Bay - Guangdong, Chine - Nucléaire	广东大亚湾-广州，中国-核电站
De 1984 à 1994, EDF a dirigé, pour le compte de la société propriétaire, la construction et la mise en service de deux réacteurs PWR 2 x 1 000 MW, grâce au savoir-faire de son ingénierie.	1984年至1994年，法国电力集团受业主公司委托，利用其工程技术实力，负责组织2x1000兆瓦压水堆的建设和运行。
En pointe, près de 100 salariés d'EDF ont travaillé sur le site.	在建设高峰期，近一百名法国电力集团工程技术人员在现场工作。
Les excellentes performances de cette centrale constituent une des principales références du Groupe en Chine.	大亚湾核电站出色的运行绩效是法国电力集团在中国的主要参照系之一。
Aujourd'hui, EDF participe à son exploitation par l'intermédiaire de contrats d'assistance et de transfert de technologie.	目前，法国电力集团与业主以技术支持和技术转让合同方式，参与大亚湾核电站的运行管理。
Ling Ao 1 et 2- Guangdong, Chine - Nucléaire	广东岭澳1-2期核电站，中国-核电站
Identique à celle de Daya Bay, cette centrale en est la suite logique :	岭澳核电站是大亚湾核电站的后续项目，采用了翻版技术：
quelques dizaines d'ingénieurs d'EDF ont assisté et conseillent encore les responsables locaux sur la construction et l'exploitation des deux premières tranches mises en service en mai 2002 et janvier 2003 et des deux secondes tranches de Ling Ao phase II mises en service en septembre 2010 et en août 2011.	法国电力集团几十名工程师为岭澳核电站2002年5月和2003年1月运行投产的一期两台机组以及2010年9月和2011年8月运行投产的二期两台机组的建设和运行都提供了技术支持和咨询服务。
Dans ce cadre, EDF a signé en avril 2005 un contrat d'assistance avec CNPEC (China Nuclear Power Energy Co.).	2005年4月，法国电力集团与中国广东核电工程公司签署了一份新的技术支持合同，同时还长期参与中方核电运营公司的单项作业。
D'autres opérations ponctuelles associent régulièrement EDF et les opérateurs chinois du nucléaire.	同时还长期参与中方核电运营公司的单项作业。
En présence des deux Présidents de République française et chinoise, le Président d'EDF a signé le 26 octobre 2006 un accord de partenariat industriel avec CGNPC (China Guangdong National Power Co.) qui a prolongé la coopération entre les partenaires dans le nucléaire où EDF souhaite investir, en associant les compétences de toute la filière française.	在中法两国国家元首的见证下，法国电力集团董事长与中国广东核电集团董事长2006年10月26日签署了工业合作伙伴协议。双方将利用法国核电的技术和经验，继续在法国电力集团希望投资的核电领域开展合作。
EDF en Asie	法国电力集团业务遍布亚太地区
Découvrir la carte des implantations en Asie	查看亚太地区代表机构和子公司分布图
Nucléaire	核电

Français	Chinois
Le développement de l'énergie nucléaire est un enjeu majeur pour la Chine et le reste du monde dans le cadre de la préservation de l'environnement et de la réduction de l'effet de serre.	发展核能对中国和全球都具有非常重要的意义，是保护环境和减缓温室效应十分重要的途径。
C'est pourquoi l'industrie nucléaire chinoise connaît une expansion croissante depuis les 20 dernières années en s'appuyant sur des compétences de haut niveau en matière de recherche et développement.	因此，二十年来，中国加大了研发力度，核电建设快速发展。
Dans ce cadre, EDF a mis en œuvre une stratégie d'intégration dans des partenariats de long terme, avec un accompagnement de l'industrie nucléaire chinoise vers son autonomie à travers un programme nucléaire réussi.	在这种背景下，法国电力集团实施长期合作伙伴关系一体化战略，支持中国核电建设，实现核电自主化目标。
De nombreux projets sont menés avec une coopération franco-chinoise très active, notamment dans le cadre de la mise en œuvre opérationnelle des centrales nucléaires de Daya Bay, Ling Ao (phases 1 et 2), mais aussi Qinshan 2 et Tianwan.	中法双方积极推动开展了多个合作项目，如大亚湾核电站和岭澳核电站（1 期和 2 期）的运行、秦山 2 期和田湾核电站的相关项目。法国电力集团自 1984 年开始负责大亚湾项目的各项技术工作：
A Daya Bay, EDF a assuré dès 1984 la responsabilité technique de la conception d'ensemble, la surveillance des fabrications, la maîtrise d'ouvrage de la construction et la mise en exploitation de la centrale.	工程总体设计、设备监造、建设总承包和核电站运行。2009 年法国电力集团和中国广东核电集团合资建立了广东台山核电合营有限公司（TNPJVC），共同建设和运营台山核电站两台 EPR 机组。
Daya Bay et Ling Ao - Chine	大亚湾核电站和岭澳核电站
Après avoir conduit la conception, la construction et la mise en service en 1994 de Daya Bay (2 réacteurs nucléaires de 1000 MW chacun), puis assisté la société propriétaire China Guangdong Nuclear Power Co.	法国电力集团承担了大亚湾核电站（两台 1000 兆瓦机组）的设计和建設，1994 年投产运行。
(CGNPC) pour la construction des deux tranches de la centrale de Ling Ao phase I (2 × 1 000 MW), mises en service respectivement en 2002 et 2003, EDF apporte aujourd'hui une assistance à la société Daya Bay Nuclear Operation and Management Co. Ltd dans le domaine de l'exploitation.	此后，集团协助业主单位中国广东核电集团建设岭澳核电站一期的两台机组（2x1000 兆瓦），分别于 2002 年和 2003 年投产运行。
Les performances enregistrées depuis leur mise en service constituent une des principales références du Groupe en Chine.	目前集团仍为大亚湾核电运营管理有限公司提供运行技术支持。
EDF est également intervenue en assistance à la filiale de CGNPC, China Nuclear Power Energy Corporation (CNPEC) sur le projet Ling Ao phase II qui consistait à construire deux nouvelles tranches de 1 000 MW sur ce site.	大亚湾和岭澳核电的运行绩效已成为法国电力集团在中国的重要参照项目。
Les deux unités de Ling Ao phase II ont été mises en service respectivement en septembre 2010 et en août 2011.	另外，法国电力公司为岭澳核电站二期（两台 1000 兆瓦机组）工程管单位——中广核工程公司（中广核集团子公司）也提供技术支持，两台机组已分别于 2010 年九月和 2011 年八月投产运行。
En savoir plus sur la centrale de Daya Bay	查看大亚湾核电站

Français	Chinois
En savoir plus sur la centrale de Ling Ao	台山核电站一期工程
Taishan 1&2 - Chine	这是中国签署的第一个中外合资核电建设项目协议。
Premier accord signé pour un investisseur étranger en Chine dans la production d'électricité nucléaire.	2006年10月，法国电力集团和中国广东核电集团宣布建立工业合作伙伴关系，2008年8月10日在北京签署了关于合资建立广东台山核电合营有限公司（TNPJVC）协议。
Faisant suite au partenariat industriel annoncé en octobre 2006, EDF et l'électricien chinois China Guangdong Nuclear Power Holding Company (CGNPC) ont concrétisé le 10 août 2008, à Pékin, les accords finaux de création d'une joint-venture dénommée «Taishan Nuclear Power Joint Venture Company Limited » (TNPJVC), dont l'objet est de construire et d'exploiter deux centrales nucléaires de technologie EPR à Taishan, dans la province du Guangdong, sur le modèle du réacteur EPR actuellement construit par EDF à Flamanville en Normandie.	双方将参照法国电力集团弗拉芒维尔核电站（位于诺曼底地区）建设EPR机组的模式，共同建设和运营台山核电站两台EPR机组。法国电力集团持有台山核电合营有限公司30%的股份，经营期50年。
La participation d'EDF au sein de TNPJVC s'élève à 30 % pour 50 ans, soit la durée maximale autorisée pour une joint-venture en Chine.	这是目前中国允许的中外合资经营企业的最高期限。
Le Groupe devient ainsi pour la première fois investisseur dans la production nucléaire dans ce pays.	法国电力集团因此成为第一个在中国核电领域投资的外国企业。
En savoir plus sur la centrale de Taishan 1	查看台山核电站一期工程
>> Activités d'EDF en Asie	>>法国电力集团亚洲地区业务概览
Présentation des activités d'EDF en Asie	法国电力集团亚洲地区业务发展概况
Activités du Groupe en Chine	法国电力集团在华业务发展概况
Présent depuis près de 30 ans en Chine au travers de prestations de conseil dans les domaines et ,le groupe EDF est aujourd'hui l'un des plus importants investisseurs étrangers dans la production d'électricité par ses participations dans des d'une puissance totale installée de 4920 MW.	近三十年来，法国电力集团为中方提供了大量和咨询服务，并投资建设了总装机容量为4920兆瓦的火电项目，现已成为中国最大的外国电力投资商之一。
EDF développe des partenariats lui ouvrant de nouvelles perspectives d'investissement, dans le nucléaire, le thermique charbon technologiquement le plus avancé, l'efficacité énergétique, les énergies renouvelables, le transport, la distribution d'électricité.	法国电力集团与中方建立的合作关系使它在核电、最先进的、能源效率、可再生能源、以及输配电领域都获得了新的投资发展前景。
nucléaire	核电
hydraulique	水电
centrales thermiques au charbon	煤电技术
Activités du Groupe en Asie du Sud-Est	法国电力集团在东南亚地区业务发展概况

Français	Chinois
L'activité du groupe EDF en Asie du Sud-Est est centrée sur le développement du secteur électrique de la zone du Grand Mékong dont la Thaïlande et le Vietnam sont les moteurs.	法国电力集团东南亚地区的主要业务是开发大湄公河地区的电力项目。泰国和越南是该地区经济增长的发动机。
La région du Mékong offre des opportunités de type « Independent Power Plants » (IPP), comme et (Vietnam).	湄公河地区为集团提供了以"独立发电厂 (IPP)"模式建设老挝和越南燃气电厂的机遇。
Dans cette perspective, EDF étudie l'intérêt de sa participation, au travers de partenariats, à la conception, à la construction et à l'exploitation de nouvelles centrales de production thermique, hydraulique et, à plus long terme, nucléaire.	鉴于这一前景，法国电力集团正在开展研究，拟通过合作方式，参与火电和水电新项目以及今后核电项目的设计、建设和运行。
Nam Theun 2	南屯 2 号
Phu My2.2	富美 2.2
L'enjeu asiatique	来自亚洲的挑战
Partenariats industriels	工业合作伙伴关系
La performance industrielle d'EDF résulte de sa capacité d'optimisation technique et économique de son parc de production.	法国电力集团通过对电力生产设施的技术和经济优化，取得了良好的工业绩效。
Grâce aux compétences acquises en France, avec son ingénierie intégrée, le Groupe élabore et propose des services d'ensemblier qui correspondent aux besoins des exploitants.	集团依托本国的技能专长和一体化工程管理优势，为运行单位制定和推荐满足其需求的整体工程服务方案。
La démarche d'EDF s'appuie sur l'optimisation du retour d'expérience de construction et d'exploitation des centrales mais aussi sur la connaissance des matériels et des fournisseurs d'équipements et plus généralement du tissu industriel.	法国电力集团利用优化核电站建设和运行的经验反馈，完全掌握了设备、供应商及工业制造体系的情况，目前已成为国际参照样板。
Pour maintenir et enrichir ses compétences et son savoir-faire thermique, hydraulique et nucléaire, le groupe EDF développe des partenariats sous forme d'accords ou de participations avec les acteurs des programmes électriques de Chine mais aussi d'Asie du Sud (Laos, Vietnam,...).	为了保持和提升其火电、水电和核电的技能和实力，法国电力集团通过协议或参股方式积极发展与中国和南亚地区（如老挝、越南等国）电力公司的合作伙伴关系。
De cette façon, les expériences et le savoir-faire de chacun sont mis en commun et partagés.	各方通过这种合作方式实现技术交流与经验共享。
Transport et distribution :	输配电：
des prestations de service	技术支持服务
Partenaire de State Grid Corporation of China, de China Southern Grid et d'acteurs locaux de gestion de réseaux, EDF fournit de nombreuses prestations d'assistance technique et de formation dans le domaine de la gestion des réseaux du transport et de la distribution.	法国电力集团为其合作伙伴——国家电网公司、中国南方电网公司以及多家地方电网公司提供了多项关于输配电网管理的技术服务和培训。
La présence du groupe EDF en Chine évolue vers un partenariat industriel de long terme reposant sur des projets d'investissements conjoints avec des sociétés chinoises dans le domaine de la conception, de la construction et de l'exploitation dans le secteur électrique.	法国电力集团在中国的发展方向是依托与中国企业共同投资项目而结成长期工业合作伙伴关系。目前主要业务集中在电力项目的工程设计、建

Français	Chinois
	设和运行，适时还将在输配电领域开展合作。
Depuis septembre 2011, ERDF-I est présent en Chine et promeut, dans la continuité des contacts et visites préalablement réalisés avec EDF Chine, les coopérations et projets potentiels avec cinq partenaires principaux.	法国配电公司国际部自 2011 年九月起在中国设立人员，在法国电力集团在中国输配电方面做的一些前期工作的基础上继续同五大主要合作伙伴发展潜在合作项目。
La cible à long terme est de participer au management des distributeurs pour améliorer leur performance en utilisant du savoir-faire et des technologies.	长期目标是运用自身的技能和技术帮助提高配电管理质量。
Mentions Légales	法律申明
La médiathèque EDF	法国电力公司文库
Organisation en Asie > La délégation Japon-Corée	亚洲地区组织机构>日韩代表团
La délégation Japon-Corée	法国电力集团日韩代表团
La Délégation EDF Japon - Corée a été créée en 2005 pour coordonner les échanges et projets entre le groupe EDF et ses partenaires Japonais et Coréens en cohérence avec sa stratégie globale. Le bureau de représentation est localisé à Tokyo.	法国电力集团驻日韩代表团成立于 2005 年，总部设在日本东京，在集团全球战略的框架下，负责与日韩合作伙伴开展沟通，做好项目的组织协调工作。
La mission de la délégation Japon-Corée Elle a pour mission de faciliter les coopérations bilatérales EDF - Japon / Corée sur tous les thèmes liés aux activités du groupe ainsi que de conduire une veille technologique et économique.	日韩代表团的职责代表团的职责是开展经济和技术跟踪，满足日韩相关单位的要求，协助引导与法国电力集团开展业务合作。在高附加值领域，可通过建立伙伴关系、技术交流等形式开展合作。
Grâce à cette présence locale, la délégation Japon - Corée permet de mettre en œuvre des échanges et partenariats sur des thèmes qui représentent de fortes valeurs ajoutées pour les métiers du Groupe.	在核能、能源和研究领域，法国电力集团与日韩两国电力公司具有共同的战略利害关系。
Depuis l'accident nucléaire de Fukushima, l'amélioration globale des critères de sûreté nucléaire conduit à un renforcement des efforts de coopération entre EDF et les organisations nucléaires Japonaises.	由于福岛核事故，核安全标准全面改进导致了加强法国电力集团和日本核组织之间的合作努力的。
Riches d'industries technologiques de pointe, le Japon et la Corée représentent de multiples enjeux pour le Groupe EDF.	丰富的高科技产业，日本和韩国都为法国电力集团多重挑战。
Les défis énergétiques du Japon post-Fukushima et de la Corée en plein essor économique pourraient permettre des progrès technologiques importants dans les nouvelles technologies de l'énergie.	日本福岛事件和韩国经济蓬勃发展后的能源挑战可能导致新能源技术显著的技术进步。
Découvrez la délégation Japon-Corée	查看日韩代表团
Contact	联系方式
Urban Toranomom, BLDG, 5F, 1-16-4 Toranomom, Minato-ku, Tokyo 105-0001 Japan	虎门市，，虎门港区大厦 5 楼 1-16-4，日本东京 105-0001

Français	Chinois
EDF - délégation Japon-Corée	法国电力集团日韩代表团
Nous contacter	联系我们
Organisation en Asie La direction Asie-Pacifique	亚洲地区组织机构亚太区总部
La direction Asie-Pacifique	法国电力集团亚太区总部
Les activités d'EDF se concentrent sur la Chine et sur la région du Grand Mékong, zones à fort développement, dans lesquelles EDF souhaite consolider et poursuivre ses investissements.	法国电力集团亚太区致力于中国和大湄公河区域等经济发展迅猛地区的业务发展，巩固和继续开发投资项目。
Au Japon et en Corée du Sud, EDF assure une veille active et procède à des échanges technologiques (ingénierie, R&D).	在日本和韩国主要开展工程和研发信息跟踪和技术交流。法国电力集团参与中国等亚洲国家的发展，这对保证其工业技术长期处于国际先进水平具有十分重要的意义。
Le principal enjeu de la présence du Groupe en Asie et particulièrement en Chine est C'est en effet dans cette zone que se déplace actuellement le centre de gravité de l'industrie électrique mondiale, soutenu par d'importants programmes de R&D et de construction de moyens de production et de transport.	在重要的研发规划和发输电项目的支撑下，世界电力工业重心目前已转移到亚太地区。
Grâce à sa présence depuis plus de 20 ans dans la zone et les relations mises en place au niveau institutionnel et industriel, EDF a aujourd'hui le maintien au niveau mondial de sa compétence industrielle sur le long terme.	法国电力集团二十多年来始终参与亚洲地区的发展，与政府部门和企业建立了密切的关系。
l'opportunité de prendre une part active à son développement, dans une optique de valorisation de son savoir-faire dans le nucléaire comme dans les autres technologies de production.	今天，适逢亚洲发展机遇，法国电力集团将积极参与，充分利用其核电和其它电力行业的技术资源和优势。
Une équipe au service du projet industriel du Groupe en Asie	服务于工业项目集团的亚洲团队
Directeur de la direction Asie-Pacifique	亚太区总裁
Directeur adjoint Ressources Humaines et Services Corporate	人力资源和企业沟通副总裁
Martin Leys	马乐思
Directeur de la division Chine	中国首席执行官
Directeur financier	财务总经理
Jing Zeng	曾静
Directeur de la division Asie du Sud	南亚首席执行官
Jean-Christophe Philbé	Jean-ChristophePhilbe

Français	Chinois
Directeur de la délégation Japon-Corée	驻日韩总代表
Thierry Knockaert	ThierryKnockaert
EDF en Asie	法国电力集团业务遍及亚太地区
Découvrir les bureaux de représentations de EDF en Asie	查看法国电力集团亚太地区代表处和子公司
Organisation en Asie La division Asie du Sud	亚洲地区组织机构南亚分部
La dynamique économique de l'Asie du Sud incite EDF à y développer des projets avec ses partenaires industriels, afin de poursuivre ses activités d'ingénierie des réseaux électriques et de production hydraulique et thermique.	亚洲地区经济增长强劲，激励法国电力集团与其工业合作伙伴携手开发该地区的电网、水电和火电项目。
Des partenariats stratégiques	战略合作伙伴关系
Au sein de l'Association of South East Asian Nations (ASEAN), les pays de la zone du « Grand Mékong » (Thaïlande, Vietnam, Laos, Birmanie et Cambodge) poursuivent l'objectif d'une intégration progressive du bassin du Mékong au niveau électrique, économique et politique, dont la Thaïlande et le Vietnam sont les moteurs.	东南亚国家联盟（ASEAN）大湄公河地区各国（泰国、越南、老挝、缅甸和柬埔寨）正在逐步推动湄公河流域电力、经济和政治一体化目标。
Le groupe EDF est le partenaire stratégique de nombreux pays membres de l'ASEAN.	泰国和越南是该地区的经济引擎。法国电力集团与许多东盟成员国结成了战略合作伙伴关系。
La mission de la Division Asie du Sud est de développer des infrastructures électriques avec ses partenaires locaux, et recouvrent deux types d'enjeux :	与当地合作伙伴共同开发电力基础设施是法国电力集团南亚分部的职责，主要包括两个层面：
technique et économique d'une part et acceptation par les populations locales d'autre part.	技术经济和公众接受。
La Division est basée à Bangkok. Elle pilote des bureaux de représentation d'Hanoi (Vietnam) et de Delhi (Inde).	南亚分部设在曼谷，负责指导河内（越南）、万象（老挝）和德里（印度）代表处的工作。
La Division assure également le pilotage des deux filiales de la région :	南亚分部还负责指导法国电力集团该地区两家子公司的工作：
et .	和.
MECO	湄公能源公司
NTPC	南屯 2 发电公司
Les bureaux de représentation en Asie du Sud	南亚分部各代表处
Les bureaux de représentation assurent l'interface avec les différents acteurs régionaux du secteur de l'énergie en Asie du Sud.	各代表处负责与南亚地区能源机构和企业开展接口与沟通工作。
Découvrez le bureau de représentation en Thaïlande	查看泰国代表处

Français	Chinois
Découvrez le bureau de représentation au Vietnam	查看越南代表处
Découvrez le bureau de liaison en Inde	查看印度联络处
Division Asie du Sud en Thaïlande	法国电力集团泰国代表处
Hanoi Central Office Building Suite 1705 -17th floor 44 B Ly Thuong Kiet Hanoi, Vietnam	河内中央办公楼 1705 室，17 楼 44B 阮文黎商信武文杰河内，越南
Division Asie du Sud au Vietnam	法国电力集团越南代表处
Division Asie du Sud en Inde	法国电力集团印度联络处
Organisation en Asie La division Chine	亚洲地区组织机构中国分部
Présente en Chine depuis près de 20 ans, la Division Chine d'EDF est installée à Pékin avec l'objectif d'aider le Groupe à devenir un des maîtres d'ouvrage du plus grand programme électrique au monde.	二十年前，法国电力集团在北京设立了在中国分部，其职责是协助集团进入中国这一全球最大的电力市场，使其成为项目业主单位之一。
La stratégie du Groupe EDF en Asie	法国电力集团亚洲发展战略
Le groupe EDF a mis en place une stratégie d'intégration grâce à des partenariats de long terme. Il peut ainsi prendre part aux nouveaux développements technologiques dans les domaines , et .Les activités d'EDF en Chine sont portées par :	法国电力集团依托长期合作伙伴关系，制定了一体化发展战略，将参与和领域的新技术开发项目，法国电力集团中国分部的主要业务部门是：
thermique	火电
Le département « Ingénierie industrielle » :	工业工程部：
dispose d'un réseau de compétences et expertise en conception et conduit les activités techniques d'EDF en Chine.	拥有设计方面的技术与评估网络，负责组织管理法国电力集团在华各项技术活动；
Il contrôle les performances des filiales d'EDF en Chine et nourrit en retour la maîtrise industrielle d'EDF en relation étroite avec les énergéticiens avec qui le Groupe développe ses projets.	监管集团在华子公司的绩效，加强与能源项目协作单位的合作，向子公司提供集团的工程管理经验；
Le département « Investissement et développement » :	投资与开发部：
développe les projets d'investissement d'EDF en Chine.	开发法国电力集团在中国的投资项目；
Le département « Energie nucléaire » :	核能部：
porte les activités d'EDF dans le nucléaire sous toutes ses formes: services et investissements.	通过咨询服务、投资等各种形式，开展集团核能领域的各项业务；
Le département « Support » :	支持部：
gère les ressources humaines et financières, et assure la logistique nécessaire aux activités de la Division Chine.	负责管理中国分部的人力资源、财务和后勤服务。

Français	Chinois
Les bureaux de représentation d'EDF en Chine	法国电力集团中国分部各代表处
Les bureaux de représentations en Chine assurent l'interface avec les différents acteurs régionaux du secteur de l'énergie.	法国电力集团中国分部各代表处负责与当地能源机构和企业开展接口与沟通工作。
Découvrez les bureaux de représentation en Chine	查看中国分部各代表处
58, Mao Ming Nan Road Jin Tai Office Building, Room 407 Shanghai 200020 China	中国 200020 上海市卢湾区茂名南路 5 8 号锦泰办公楼 407 室
EDF à Shanghai	法国电力集团中国分部上海联络处
LA6215, Technical Centre, Daya Bay Site Shenzhen 518124 China	中国 518124 深圳大亚湾核电基地技术中心 LA6215
EDF à Daya Bay	法国电力集团中国分部大亚湾代表处
Nos filiales en Chine	在中国的子公司
Filiales d'EDF en Asie	在中国的子公司
Sites de production	工业设施
Implantations	代表机构和子公司
Retour	返回
Sites de production	生产基地
Présentation de EDF en Asie	法国电力集团亚洲地区业务概况
Les activités du groupe EDF conduites par la Direction Asie-Pacifique, se concentrent sur la Chine et sur la région du Grand Mékong, des pays à fort développement.	法国电力集团亚太区总部致力于中国和大湄公河区域等经济发展迅猛地区的业务发展。
L'investissement dans le secteur de la production électrique en Asie et particulièrement en Chine, constitue un enjeu industriel pour le groupe EDF	投资亚洲国家，尤其是中国的电力行业对法国电力集团具有重要产业意义。
En complément des projets comme l'EPR, les nouveaux projets dans cette zone donneront au Groupe l'accès aux innovations technologiques et lui permettront dans le même temps de valoriser son savoir-faire industriel, en particulier nucléaire.	除 ERP 项目外，亚洲地区的新项目将为集团提供技术创新机会，并有利于充分利用其核电等行业的技术资源。
EDF maintiendra, ainsi, ses atouts concurrentiels et technologiques dans un contexte de compétition internationale pour la relance du programme nucléaire mondial, pour l'équipement de pays émergents et dans la perspective du renouvellement du parc français.	面对当前全球核电复苏的国际竞争局面，法国电力集团将继续保持技术竞争优势，积极参与新兴市场国家的建设，同时着手本国核电站的更新换代。
Les activités d'EDF en Asie	法国电力集团亚洲地区业务概况
Les activités d'EDF en Asie (Chine, Thaïlande, Vietnam, Laos) s'inscrivent dans une stratégie de développement.	推动亚洲地区（中国、泰国、越南、老挝）的业务发展是实施法国电力集团发展战略的重要举措。

Français	Chinois
Elles se traduisent par la participation du Groupe dans des projets de longue durée, marqués par un esprit de coopération et de partenariat, dans le respect des valeurs sociales et environnementales du Groupe.	法国电力集团积极参与长期发展项目，致力于建立合作和合作伙伴关系，信守企业的社会和环境价值理念。
Elles sont pilotées par la direction Asie-Pacifique du Groupe, installée à Pékin.	法国电力集团亚太区总部设在北京，负责指导该地区的业务发展。
Découvrez les installations industrielles et les filiales associées en Chine, au Laos et au Vietnam.	查看中国、老挝、越南的相关工业项目和子公司。
Découvrez les partenariats industriels autour du nucléaire, du gaz naturel, du charbon propre et de l'hydraulique.	查看核电、天然气、洁净煤和水电项目的合作伙伴。
L'organisation d'EDF en Asie	亚洲地区的组织机构
Le groupe EDF en Asie est représenté par la direction Asie-Pacifique. Le Directeur est Monsieur Hervé Machenaud.	法国电力集团驻亚洲的代表机构是亚太区总部，马识路是亚太区总裁。
La direction est pleinement responsable de la conduite des activités du groupe EDF pour l'ensemble de la région.	亚太区总部全面负责指导法国电力集团亚太地区的业务发展。
L'organisation EDF en Asie est composée:	亚太区各机构包括：
D'une équipe au service du projet industriel du Groupe en Asie.	一个为集团亚洲工业项目提供服务的团队：
De deux organisations au service des activités du Groupe et de ses filiales locales.	两个为集团及地区子公司拓展业务服务的团队：
D'un bureau de représentation EDF.	一个法国电力集团代表处：
La délégation Japon-Corée	驻日韩代表团
Henderson Centre, Tower 2, Floor 12 18 Jianguomennei Avenue Beijing 100005 China	中国 100005 北京建国门内大街 18 号恒基中心 2 座 12 层
EDF - direction Asie-Pacifique	法国电力集团亚太区总部
Actualités	新闻
Résultats semestriels 2013 en hausse...	法国电力集团 2013 年上半年业绩...
EPR de Flamanville :	法国 EPR 核电站建设的重要里程碑：
pose du dôme spectaculaire...	法国电力集团成功吊装佛拉芒维勒 EPR 核电站穹顶...
Présence et participations	代表处与参股项目
Le groupe EDF participe à de grands projets énergétiques.	法国电力集团以参股方式参与重大能源项目的建设
Découvrir nos implantations en Asie	查看法国电力集团在亚洲的工业项目
Organisation d'EDF en Asie	法国电力集团亚洲组织机构

Français	Chinois
Les activités du groupe EDF conduites par la Direction Asie-Pacifique, se concentrent sur la Chine et sur la région du Grand Mékong, des pays à fort développement.	亚太总部负责指导开展集团亚太地区的业务活动，重点是中国和大湄公河区域等经济发展迅猛的地区。
Découvrir l'organisation d'EDF en Asie	查看法国电力集团亚洲组织机构
Toutes nos activités en 3D	法国电力集团业务三维展示
Le contenu de cette page nécessite une version plus récente d'Adobe Flash Player.	本页的内容需要最新的 AdobeFlashPlayer 版本。
Les dates clés du Groupe	法国电力集团重大事件
création d'EDF (statut d'EPIC) mise en service sur une période de 20 ans de 54 tranches nucléaires représentant une puissance installée de 56 600 MW.	法国电力公司成立（工商企事业单位）54 台核机组在 20 年内陆续建成并投入运行，总装机容量 56600 兆瓦。
EDF devient une Société Anonyme à conseil d'administration introduction en Bourse Démarrage de la construction de l'îlot nucléaire de l'EPR à Flamanville.	法国电力公司改制，成为董事会制股份有限公司。公司股票正式上市。
Aujourd'hui, le parc nucléaire d'EDF SA dispose de 58 tranches nucléaires et d'une capacité installée de 63 100 MW	目前，法国电力集团共有 58 台核机组，总装机容量 63100 兆瓦。
à partir des années 70 :	从上世纪 70 年代起：
Identité d'entreprise	企业特性
Valeurs d'entreprise	企业价值观
>> Valeurs > Ethique et gouvernance	>>企业价值观>行为准则与管理
L'ambition d'EDF est de contribuer à bâtir un avenir meilleur fondé sur le développement durable.	“遵循可持续发展理念，为创建更美好的未来做贡献”是法国电力集团的远大抱负。
Nos solutions sûres et fiables nous permettent de répondre à la croissance des besoins énergétiques dans le monde, tout en luttant contre le changement climatique et en préservant les ressources naturelles.	我们掌握各种安全可靠的解决方案，可满足全球不断增长的能源需求，并能应对气候变化和保护自然资源的挑战。
Nous construisons un avenir meilleur en produisant des énergies propres et en réduisant constamment leur impact environnemental.	我们生产清洁能源，不断减少对环境的影响，努力创建更美好的未来。
Les valeurs d'EDF:	法国电力集团的企业价值观：
Dès les années 2000, il nous est apparu nécessaire de formuler l'éthique et les valeurs qui guident notre action et nos relations avec nos parties prenantes.	自 2000 年以来，集团确定了必要的行为准则和企业价值观，指导我们的行为和合作伙伴的关系。
Les 5 valeurs qui animent notre démarche :	激励我们行动的五大价值观是：
le respect de la personne,	尊重个人,

Français	Chinois
la responsabilité environnementale,	环境责任,
la recherche de la performance,	争取绩效,
l'engagement de solidarité,	团结互助,
l'exigence d'intégrité.	廉洁自律.
Les sociétés du Groupe travaillent actuellement à l'élaboration d'un référentiel qui sera traduit par une politique éthique commune.	目前，集团各公司正在制定诠释共同行为准则方针的基准参照。
Pour en savoir plus	了解更多的相关信息
Global Compact EDF est signataire du depuis 2001	法国电力集团 2001 年签署了联合国全球契约 (GlobalCompact)
EDF participe au Global Reporting Initiative (GRI)	法国电力集团参与全球报告倡议组织 (GRI) 的相关工作
L'alerte éthique nous permet de recueillir les manquements au respect de nos valeurs	举报机制有利于我们掌握违背价值观的行为
L'accord sur la responsabilité sociale du groupe EDF	法国电力公司的社会责任的一致
La médiation est chargée de rechercher un règlement amiable aux litiges entre le groupe EDF et ses clients, fournisseurs et partenaires.	调解是负责寻求和解，在法国电力集团及其客户之间，以及供应商及合作伙伴之间的纠纷。
Impartialité et transparence sont les conditions de réussite de cette démarche	公正性和透明度是这种方法的成功的条件
Le groupe EDF	法国电力集团
Les engagements du groupe EDF	法国电力集团的承诺
Valeurs Objectifs et engagements	企业价值观目标与承诺
Répondre à la demande croissante d'énergie tout en parant aux risques climatiques et à la raréfaction des ressources : c'est le défi lancé aux énergéticiens.	既能满足日益增长的能源需求，又要应对气候变化风险和资源枯竭问题，这就是能源企业面临的挑战。
C'est pourquoi le développement durable est au cœur de notre stratégie.	为此，可持续发展成为法国电力集团的战略核心。
Notre politique traduit la volonté du Groupe de « changer l'énergie ensemble » en apportant des solutions réalistes.	“共同改变能源结构”体现了法国电力集团的意愿，为此，我们将提供切实可行的技术方案。
La lutte contre le changement climatique commence par la maîtrise de nos impacts environnementaux.	应对气候变化首先应从控制环境影响做起。
Nous nous sommes fixé 9 engagements pour répondre à 3 enjeux prioritaires	我们制定了九项承诺，优先应对以下三方面的挑战：
La lutte contre le changement climatique et la préservation de la biodiversité :	应对气候变化，保护生物的多样性：
Un enjeu environnemental	环境的挑战

Français	Chinois
rester, en tant que Groupe, le moins émetteur de CO des grands énergéticiens européens,	法国电力集团将继续保持欧洲大型能源公司一氧化碳排放最少企业的地位；
adapter notre parc de production et nos offres au changement climatique,	电力生产设施和服务项目适应气候变化的要求；
réduire notre impact environnemental, notamment sur la biodiversité.	减少对环境，特别是对生物多样性的影响。
Faciliter l'accès à l'énergie et développer des liens de proximité avec les territoires :	方便获取能源和发展与地区保持紧密的联系：
favoriser l'accès à l'énergie et l'éco-efficacité énergétique,	为普及能源使用提供便利，发挥能源项目的生态效益；
développer dans la durée la proximité avec les territoires où nous opérons,	与地方合作，在集团经营区域内长期实施就近供电政策；
contribuer à l'effort éducatif sur les questions liées à l'énergie.	支持能源教育的各种活动。
Contribuer au débat sur le développement durable par le dialogue, l'information et la communication :	支持关于可持续发展的辩论，可采用对话、信息与交流等各种形式：
Un enjeu de gouvernance	管理方面的挑战
poursuivre le développement des politiques et le partage des valeurs au sein du Groupe, en relation avec les parties prenantes,	与有关方面合作，在集团内部继续制定相关政策，共享企业价值观；
communiquer et rendre compte des activités et résultats du Groupe en matière de développement durable,	交流和通报集团在可持续发展方面开展的工作和取得的成果；
participer au débat sur le développement durable au niveau national et international.	参与法国和国际有关可持续发展的讨论。
Pour aller plus loin	进一步了解其他相关信息
La politique de développement durable d'EDF Energy	EDF Energy 公司（法国电力集团在英国的子公司）可持续发展政策
La politique de développement durable d'Edison	意大利 Edison 能源公司（法国电力集团拥有该公司的股份）可持续发展政策
Valeurs Partenariats	企业价值观合作伙伴关系
Les partenariats prolongent l'engagement économique, social et environnemental du Groupe dans la société.	合作伙伴关系是集团参与社会各个方面（经济、社会、环境等）活动的延伸。
Portant ses valeurs, ils constituent un moyen de consolider des relations de proximité avec ses parties prenantes, les associations et les collectivités territoriales notamment.	合作伙伴宣传集团的企业价值观，是集团加强和密切与协会、地方行政单位和有关方面联系的手段。
Nous construisons un avenir meilleur en produisant des énergies propres et en réduisant constamment leur impact environnemental.	我们生产清洁能源，不断减少对环境的影响，努力建设更美好的未来。
Le Groupe a choisi de privilégier deux champs d'intervention :	集团择先支持以下两方面的活动：

Français	Chinois
l'environnement et l'énergie,	环境和能源；
le lien social et la solidarité.	社会关系和互助。
les partenariats contribuent à la protection de la nature (protection d'espaces et d'espèces) et à la lutte contre le changement climatique (réduction de l'empreinte carbone), ainsi qu'au soutien aux sciences, à l'innovation et aux technologies.	合作伙伴为自然保护（空间环境和物种）作贡献，应对气候变化（减少碳足迹），支持科学、创新和技术活动。
Dans le domaine de l'environnement et de l'énergie,	在环境和能源方面：
les partenariats concernent la diversité et la proximité (accès à l'énergie et l'éco-efficacité, accès à l'emploi et à l'insertion, lutte contre l'exclusion...) ainsi que la santé et la recherche médicale (études médicales et neurosciences...).	合作伙伴开展的活动涉及能源的多样性和就近供电（普及能源使用和提高生态效益、就业与融入、反对社会排斥...）以及卫生和医学研究（医学研究、神经科学...）。
Dans le domaine du lien social et de la solidarité,	在社会关系和互助方面：
Nos partenariats sont mis œuvre par les sociétés et les fondations qui composent le Groupe :	合作伙伴的活动由相关公司和基金会具体实施，法国电力集团是这些公司和基金会的成员：
La Fondation européenne pour les énergies de demain	欧洲未来能源基金会
La Fondation EDF Diversiterre	法国电力集团全球能源多元化基金会
La Fondation RTE	输电网基金会
La Fondation Edison	爱迪生基金会
La Fondation agir pour l'emploi	促进就业行动基金会
Le principe de solidarité s'inscrit au cœur de la Fondation EDF Diversiterre.	互助原则是法国电力集团全球能源多元化基金会的核心理念。
Créée en 2007, elle a pris le relais des vingt années de mécénat du groupe EDF sous l'égide de la Fondation de France.	基金会成立于 2007 年，是法国电力集团赞助活动的执行机构，取代了过去二十多年来一直由法兰西基金会（Fondation de France）行使的职能。
EDF Diversiterre soutient des actions dans trois domaines :	基金会支持三方面的活动：
la nature et la biodiversité, la solidarité et la santé, la culture et le patrimoine.	自然和生物多样性，互助和卫生，文化和遗产。
Parmi ses partenaires figurent notamment :	集团的合作伙伴主要包括：
la Fondation Nicolas Hulot,	法国 Nicolas Hulot 基金会,
les Réserves naturelles de France,	法国自然保护基金会,
le Conservatoire du littoral,	海岸保护研究院.

Français	Chinois
les Restos du cœur,	赈济协会,
l'Association française contre les myopathies (AFM),	法国肌病防治协会 (AFM),
le Téléthon.	马拉松电视剧基金会.
Dans le monde	在其他国家
Dans le cadre de la construction du barrage de Nam Theun, au Laos, EDF s'est associé à l'Institut Pasteur pour réaliser l'accompagnement médical et le suivi sanitaire des populations déplacées.	在老挝南屯水电站项目框架下，法国电力集团与巴斯德研究所合作，对大坝移民进行医疗诊治和跟踪。
EDF a conclu un partenariat avec (UICN) pour soutenir la mise à jour et l'édition de la liste rouge des espèces menacées.	法国电力集团与(UICN)签订了合作伙伴协议，支持濒危物种红色名录的更新和出版工作。
l'Union Internationale pour la conservation de la nature	国际自然保护联盟
À Liverpool (Royaume-Uni), des agents volontaires d'EDF et d'EDF Energy ont contribué, dans le cadre d'un partenariat avec l'association Habitat pour l'humanité, à la construction de logements à efficacité énergétique renforcée, destinés à des personnes vulnérables.	在利物浦（英国），法国电力集团和 EDFEnergy 公司的志愿者在人居住房协会合作框架下，积极推动弱势群体高能效住房的建设。
EDF pilote des projets en collaboration avec des ONG, des organisations internationales et des associations locales.	法国电力集团与非政府组织、国际组织和地方协会合作，参与一些项目的指导管理。
Les partenariats conclus notamment avec ou Droit à l'énergie, visent à fournir un accès à l'énergie à des populations vulnérables.	农村电气化机构开展合作，试验推广“南北经验转让”的新形式。
Avec l' EDF teste et promeut de nouvelles formes de transfert d'expériences Nord-Sud, notamment entre les agences d'électrification rurale d'Afrique francophone.	法国电力集团、法国环境与能源署和非洲法语国家,
Avec la , EDF mène des opérations d'électrification au Laos, au Burkina-Faso, au Sénégal et à Madagascar.	法国电力集团与,在老挝、布基纳法索、塞内加尔和马达加斯加实施电气化建设项目。
Fondation Énergies pour le monde	世界能源基金会
Sponsoring sportif	体育赞助
Les partenariats noués ont pour objectif d'ancrer EDF dans le mouvement olympique et handisport et de soutenir les sports qui ont un lien direct avec les métiers de l'entreprise.	法国电力集团建立了相关的合作伙伴关系，以巩固在奥林匹克运动的和残疾人体育事业方面的地位，同时，还支持与集团有直接业务关系的体育运动项目。
EDF est sponsor officiel et partenaire développement durable des jeux Olympiques et Paralympiques de Le Groupe alimentera les Jeux en gaz et en électricité, notamment par de l'énergie renouvelable. Pour la flamme olympique, EDF fournira un combustible à faible émission de carbone.	法国电力集团是的赞助商和可持续合作伙伴，为这两次活动提供天然气和可再生能源电力，并为奥运火炬提供低碳燃料。

Français	Chinois
Londres 2012.	2012 年伦敦奥运会和残奥会
Partager les savoirs pour innover	知识共享，推动创新
L'union fait la force, en matière d'innovation comme ailleurs.	在企业技术创新方面，集团采用“团结就是力量”的方式。
Notre politique R&D s'inscrit dans une logique de codéveloppements (avec des partenaires industriels et des organismes de recherche).	与工业伙伴和研究机构联合开展研发是我们的研发方针。
Des laboratoires communs permettent de mutualiser les compétences de haut niveau et de partager expériences et bonnes pratiques.	联合实验室有利于高层次人才交流，是分享经验和探索的平台。
Un tiers des dépenses de notre R&D est lié à des études sur la protection de l'environnement.	集团三分之一的研发经费用于环境保护研究。
>> Identité	>>企业特性
Premier producteur nucléaire mondial, bien implanté dans les grands pays d'Europe, le groupe EDF investit pour une croissance industrielle durable, en portant ses priorités sur trois axes :	法国电力集团是全球最大的核电生产商，已成功进入欧洲主要国家的电力市场，为实现企业长期发展，集团投资将优先考虑以下三个层面的因素：
Etre un leader du renouveau du nucléaire dans le monde,	成为全球核电复兴的引领者；
Développer les énergies renouvelables et l'éco-efficacité énergétique,	发展可再生能源，提高能源生态效率；
Renforcer ses positions dans le Monde.	加强集团在欧洲的地位。
159 740 collaborateurs dans le monde	世界各地员工人数 159740
EDF s'appuie sur des équipes motivées et compétentes.	法国电力公司依靠积极和有能力的团队。
Le Groupe mène une politique active de recrutement, offrant à des milliers de jeunes la possibilité de participer au changement énergétique.	法国电力集团依托团队激励与业务技能优势，实施积极的招聘政策，为众多青年提供了参与改变能源结构的机会。
Les équipes d'EDF se rassemblent autour de trois valeurs clés :	法国电力集团团队的五个核心价值观是：
respect	尊重个人,
solidarité	争取绩效,
responsabilité	环境责任,
Des engagements pour un développement durable	可持续发展承诺
La politique de développement durable du Groupe engage toutes ses sociétés sur trois enjeux :	法国电力集团可持续发展方针指引所属公司围绕以下三大目标开展工作：
changement climatique et biodiversité ;	缓解气候变化，注重生物多样性；

Français	Chinois
accès à l'énergie et proximité territoriale ;	普及能源使用，实施经营区域内就近供电政策；
contribution au débat sur le développement durable.	促进关于可持续发展的辩论。
R&D :	研发：
523 millions d'euros de budget, 2 000 personnes	5.23 亿欧元预算，2000 名员工
En France, 2 000 personnes - dont 30 % de femmes - travaillent pour la recherche d'EDF selon trois grands axes :	在法国，2000 名员工，其中 30%女职工在法国电力集团研发机构工作，研发范围涉及三大领域：
Limiter les émissions de CO2 avec des alternatives aux énergies fossiles	取代化石燃料，以便限制二氧化碳的排放；
Faire bénéficier les clients des nouvelles technologies	让客户享用新技术；
Contribuer à la sécurité des réseaux électriques	为电网的安全作贡献。
EDF dans le monde	法国电力集团业务遍及世界各地
Découvrez les implantations du groupe EDF dans le monde	查看法国电力集团世界各地的代表机构和子公司
Press release	出版物
Résultats semestriels 2013 en hausse	法国电力集团 2013 年上半年业绩
EDF Group	EDFGroup
Résultats annuels 2012	2012 年全年业绩公告
Lire le communiqué	阅读新闻稿
Download the press release (331Kb)	下载新闻稿 (331KB)
Chine :	中国：
la construction du réacteur EPR numéro 1 de Taishan franchit une nouvelle étape majeure avec la mise en place de la cuve.	压力容器的安装就位标志着台山 EPR™反应堆 1 号机组的建设进入新的重要阶段。
Download the press release (52Kb)	下载新闻稿 (52KB)
Chine : succès de la pose du dôme du réacteur EPR de Taishan 1	中国广东台山核电站一号机组核岛反应堆厂房穹顶吊装成功
Paris, le 24 octobre 2011 - La pose du dôme sur le bâtiment réacteur de la tranche 1 de l'EPR de Taishan en Chine a été réalisée hier avec succès.	巴黎，2011 年 10 月 24 日-中国广东台山核电站一号机组反应堆厂房穹顶于昨日吊装成功。
Cette opération, coordonnée par le maître d'ouvrage Taishan Nuclear Power Joint Venture Company (TNPJVC) - Joint Venture détenue à 70% par CGNPC et à 30% par EDF - intervient un peu plus de 2 ans après le coulage du béton du radier du bâtiment réacteur.	吊装作业是在（由中广核集团持股 70%、法国电力集团持股 30%共同组建）项目业主-台山核电合营有限公司的协调下，在一号机组反应堆厂房筏基的第一罐混凝土浇筑两年多之后完成的。

Français	Chinois
Premières propositions d'EDF à l'ASN post-Fukushima	法国电力集团 2011 年第一季度财务信息
Pacte d'actionnaire d'Edison :	爱迪生股东公约 :
position d'EDF	法国电力集团的立场
Résultats Annuels 2010	法国电力集团 2010 年年度公报
Amélioration de la performance industrielle Des provisions exceptionnelles Une flexibilité financière retrouvée	改善工业效率出色的收益重现财务灵活性
L'énergie est notre avenir, économisons-la !	能源是我们的未来，节约能源！

Annexe 10 : Exemple d'extraction de bisgments à partir du corpus MultiUN

Français	Chinois
Rappelle également que le Sommet mondial de 2005 a décidé de renforcer la contribution des organisations non gouvernementales, de la société civile, du secteur privé et d'autres parties prenantes aux efforts de développement national ainsi qu'à la promotion du partenariat mondial en faveur du développement, et qu'il a encouragé les partenariats publics-privés dans les domaines suivants : réalisation de nouveaux investissements et création d'emplois; financement du développement; recherche de solutions aux problèmes de santé par le traitement et la recherche; et promotion de la science et de la technique au service du développement dans les secteurs de la santé, de l'agriculture, de la protection de l'environnement, de l'utilisation durable des ressources naturelles et de la gestion de l'environnement, de l'énergie, des forêts et de l'incidence des changements climatiques;	“5.又回顾 2005 年世界首脑会议决心加强非政府组织、民间社会、私营部门和其他利益有关者在国家发展努力及促进全球发展伙伴关系中的贡献，并鼓励在下列领域建立公私伙伴关系：创造新投资和就业；发展筹资；通过治疗和研究解决健康方面的挑战；以及在卫生、农业、养护、可持续利用资源和环境管理、能源、林业和气候变化的影响等领域推动科学和技术促进发展；
Encourage le système des Nations Unies à continuer d'adhérer à une conception commune et systématique des partenariats qui mette davantage l'accent sur l'impact, la responsabilité et la durabilité, sans imposer une quelconque rigidité aux accords de partenariat et en tenant dûment compte des principes régissant les partenariats énoncés dans la résolution 58/129;	“6. 鼓励联合国系统继续就伙伴关系制定一种共同的、系统的做法，其中更加侧重影响、问责制和可持续性，但不要在伙伴关系协定中施加不当的僵硬限制，并适当考虑到第 58/129 号决议规定的各项伙伴关系原则；
Encourage également des pratiques commerciales responsables telles que celles énoncées dans le Pacte mondial;	“7. 又鼓励负责任的商业做法，例如全球契约促进的那种商业做法；
Souligne l'importance d'une bonne gestion des entreprises et de la responsabilité sociale de ces dernières et encourage le Bureau du Pacte mondial à continuer de favoriser des pratiques commerciales responsables, à promouvoir l'échange de pratiques optimales et à favoriser une action positive par l'apprentissage, le dialogue et les partenariats;	“8. 着重指出良好的公司治理和公司社会责任的重要性，鼓励全球契约办事处通过学习、对话和伙伴关系，继续推动负责任的商业做法、促进分享最佳做法和促进积极行动；
Encourage le Bureau du Pacte mondial à veiller à ce que les enseignements pertinents tirés des partenariats, notamment avec les milieux d'affaires, contribuent au processus de réforme de l'Organisation des Nations Unies en cours;	“9. 鼓励全球契约办事处确保从伙伴关系包括从商业界学习到相关经验教训用于支持正在进行的联合国改革工作；
Se félicite de la nomination par le Secrétaire général d'un conseiller spécial pour le Pacte mondial;	“10. 欢迎秘书长为全球契约任命了一名特别顾问；
Prie le Secrétaire général de prendre toutes autres dispositions pertinentes pour consolider les partenariats en renforçant leurs études d'impact, leur champ d'action stratégique et leur prise en main au niveau local et en améliorant la gestion des partenariats grâce à une formation appropriée à tous les niveaux, à la mise en commun des meilleures pratiques, à la rationalisation des directives des Nations Unies pour les partenariats et à l'amélioration des procédures de sélection des partenaires;	“11. 请秘书长进一步采取适当行动来加强伙伴关系的影响评估、战略重点和当地所有权，以便巩固伙伴关系，并在所有有关各级进行充分的培训、分享最佳做法、精简联合国伙伴关系准则和改善伙伴选择程序，以便加强伙伴关系的管理；
Se félicite des méthodes novatrices adoptées par les organes et organismes compétents des Nations Unies, ainsi que les organismes issus des Accords de Bretton Woods et l'Organisation mondiale du commerce, pour tirer le meilleur parti des partenariats afin de mieux mettre en œuvre leurs objectifs et programmes, en particulier pour ce qui est du développement et de l'élimination de la pauvreté, et les encourage à continuer	“12. 欢迎联合国的相关组织和机构以及布雷顿森林机构和世界贸易组织采取创新方法，利用伙伴关系来更好地落实它们的目标和方案，特别是谋求发展和消灭贫穷，并鼓励它们

d'étudier ces possibilités, compte tenu des différents mandats, modes de fonctionnement et buts des organes et organismes ainsi que des rôles spécifiques des partenaires non étatiques concernés;	进一步探讨这类可能性，同时铭记各个机构的不同任务、运作方式和目标，以及相关的非国营部门伙伴的特别作用；
Recommande, dans ce contexte, que les partenariats visent également à éliminer la discrimination, notamment à caractère sexiste, en matière d'emploi et de profession;	“13.在这方面建议伙伴关系还应促进在就业和职业方面消灭歧视，包括基于性别的歧视；
Lance à nouveau un appel :	“14.再次呼吁：

Annexe 11 : Script de filtrage de corpus

1	#!/usr/bin/perl
2	use locale;
3	use Encode;
4	use utf8;
5	use IO::Handle;
6	
7	
8	my \$fileSource = \$ARGV[0];
9	my \$fileTarget = \$ARGV[1];
10	my \$outSource = \$fileSource.".out";
11	my \$outTarget = \$fileTarget.".out";
12	
13	open(INS,"<:encoding(UTF-8)", \$fileSource);
14	open(OUTS,">:encoding(UTF-8)", \$outSource);
15	open(INT,"<:encoding(UTF-8)", \$fileTarget);
16	open(OUTT,">:encoding(UTF-8)", \$outTarget);
17	
18	
19	
20	my @tabSource;
21	my @tabTarget;
22	
23	
24	while(<INS>){
25	my \$ls = \$ _;
26	chomp(\$ls);
27	push @tabSource, \$ls;
28	}
29	while(<INT>){
30	my \$lt = \$ _;
31	chomp(\$lt);
32	push @tabTarget, \$lt;
33	}
34	for (my \$i = 0; \$i < @tabSource; \$i++){
35	\$tabSource[\$i] =~ s/http+:\V[^\s\()]*\/url\/sg;
36	\$tabSource[\$i] =~ s/www\.[^\s\()]*\/url\/sg;
37	\$tabSource[\$i] =~ s/[^\s\()]*@[^\s\()]*\/mail\/sg;
38	
39	\$tabTarget[\$i] =~ s/http+:\V[^\s\()]*\/ url /sg;

40	\$tabTarget[\$i] =~ s/www\.[^\s\(\)]*/url /sg;
41	\$tabTarget[\$i] =~ s/^[^\s\(\)]*@[^\s\(\)]*/mail /sg;
42	
43	my @tokenSource = split(/ /, \$tabSource[\$i]);
44	my @tokenTarget = split(/ /, \$tabTarget[\$i]);
45	my \$nbTokenS = @tokenSource;
46	my \$nbTokenT = @tokenTarget;
47	
48	# if (\$nbTokenS lt 81 && \$nbTokenT lt 81 && nbTokenS gt 5 && \$nbTokenT gt 5 && \$nbTokenS < \$nbTokenT*1.5 && \$nbTokenS > \$nbTokenT*0.66){
49	if (\$nbTokenS lt 81 && \$nbTokenT lt 81) {
50	my \$nbNum = 0;
51	my \$nbPunc = 0;
52	my \$nbBlock = 0;
53	my \$nbMaj = 0;
54	
55	foreach my \$keyS(keys @tokenSource){
56	if (\$tokenSource[\$keyS] =~ /[0-9a-zA-Z]/){
57	\$nbNum++;
58	}
59	if (\$tokenSource[\$keyS] =~ /\.,;:"'«»'&*=+V\ #\√≤τΣ>Δ<@\{\}\ \[\(\)\δ\~} < /){
60	\$nbPunc++;
61	}
62	if (\$tokenSource[\$keyS] =~ /[!•n■□◆√≤τΣ>Δ<(-)□(□)四(五)六(七)八(九)(+¼)/){
63	\$nbBlock++;
64	}
65	}
66	
67	foreach my \$keyT(keys @tokenTarget){
68	if (\$tokenTarget[\$keyT] =~ /[0-9]/){
69	\$nbNum++;
70	}
71	if (\$tokenTarget[\$keyT] =~ /[A-Z]{2,50}/){
72	\$nbMaj++;
73	}
74	if (\$tokenTarget[\$keyT] =~ /\.,;:"'«»'&*=+V\ #\√≤τΣ>Δ<@\{\}\ \[\(\)\δ\~} < /){
75	\$nbPunc++;
76	}
77	if (\$tokenTarget[\$keyT] =~ /[!•n■□◆√≤τΣ>Δ<(-)□(□)四(五)六(七)八(九)(+¼)/){
78	\$nbBlock++;
79	}
80	}

81	
82	my \$n1 = \$nbPunc + \$nbNum;
83	my \$n2 = \$nbTokenS + \$nbTokenT;
84	my \$score = \$n1/\$n2;
85	
86	if (\$score < 0.6 and \$nbBlock ==0 and \$nbMaj/\$nbTokenT < 0.2){
87	print OUTS "\$tabSource[\$i]\n";
88	print OUTT "\$tabTarget[\$i]\n";
89	}
90	}
91	}
92	
93	
94	close (INS);
95	close (INT);
96	close (OUTS);
97	close (OUTT);

Annexe 12 : Source du programme pour calculer D_{mix}

1	#!/usr/bin/perl
2	
3	use locale;
4	use Encode;
5	use utf8;
6	
7	
8	
9	class editDist():
10	
11	# initialization
12	def __init__(self):
13	self.coefStr = 0.2
14	self.coefWord = 0.8
15	self.coefDelete = 1
16	self.coefInsert = 1
17	self.coefSubstitute = 1
18	
19	# set the coef of word distance and string distance
20	def setStrWordCoef(self, coefStr, coefWord):
21	self.coefStr = coefStr
22	self.coefWord = coefWord
23	
24	# set the 3 operations coefs: delete, insert, substitute
25	def setOpeCoef(self, coefDelete, coefInsert, coefSubstitute):
26	self.coefDelete = coefDelete
27	self.coefInsert = coefInsert
28	self.coefSubstitute = coefSubstitute
29	
30	# word distance function: 1) preparing by white space, wipe all the white spaces 2) count the editing distance in word level
31	def wordDist(self, input1, input2):
32	list1 = input1.strip().split() # tokenize
33	list2 = input2.strip().split() # tokenize

34	lenth1 = len(list1)
35	lenth2 = len(list2)
36	if lenth1 == 0:
37	return lenth2*self.coefInsert
38	if lenth2 == 0:
39	return lenth1*self.coefDelete
40	distanceMatrix = [[x * self.coefDelete + y * self.coefInsert for y in range(lenth2 + 1)] for x in range(lenth1 + 1)] # distance matrix
41	for i in range(lenth1):
42	for j in range(lenth2):
43	deletion = distanceMatrix[i][j + 1] + self.coefDelete # if delete
44	insertion = distanceMatrix[i + 1][j] + self.coefInsert # if insert
45	if list1[i] != list2[j]:
46	substitution = distanceMatrix[i][j] + self.coefSubstitute # if substitute
47	else:
48	substitution = distanceMatrix[i][j]
49	compt = [deletion, insertion, substitution]
50	distanceMatrix[i + 1][j + 1] = min(compt) # mini distance
51	return distanceMatrix[lenth1][lenth2]
52	
53	# string distance function: count the editing distance in string level, include white spaces
54	def strDist(self, input1, input2):
55	lenth1 = len(input1)
56	lenth2 = len(input2)
57	if lenth1 == 0:
58	return lenth2 * self.coefInsert
59	if lenth2 == 0:
60	return lenth1 * self.coefDelete
61	distanceMatrix = [[x * self.coefDelete + y * self.coefInsert for y in range(lenth2 + 1)] for x in range(lenth1 + 1)] # distance matrix
62	for i in range(lenth1):
63	for j in range(lenth2):
64	deletion = distanceMatrix[i][j + 1] + self.coefDelete # if delete
65	insertion = distanceMatrix[i + 1][j] + self.coefInsert # if insert
66	if input1[i] != input2[j]:
67	substitution = distanceMatrix[i][j] + self.coefSubstitute # if substitute
68	else:

69	substitution = distanceMatrix[i][j]
70	compt = [deletion, insertion, substitution]
71	distanceMatrix[i + 1][j + 1] = min(compt)
72	return distanceMatrix[lenth1][lenth2] # mini distance
73	
74	# combo distance function: combine string distance and word distance with 2 coefs
75	def comboDist(self, input1, input2):
76	return self.strDist(input1, input2) * self.coefStr + self.wordDist(input1, input2) * self.coefWord

Annexe 13 : 100 bisegments anglais-français extraits de CLEF-IP 2011

No	Anglais	Français
1	An additive for a lubricating oil or hydrocarbon fuel obtainable by a process which comprises reacting a polyamino alkenyl or alkyl succinimide with a compound of the general formula: wherein W is oxygen or sulfur; X is oxygen or sulfur; R ₄ is an alkylene group of 2 or 3 carbon atoms optionally substituted by from 1 to 3 alkyl groups of 1 or 2 carbon atoms each; and R ₅ is hydrogen or alkyl of from 1 to 20 carbon atoms.	Additif pour une huile lubrifiante ou un combustible hydrocarboné, pouvant être obtenu par un procédé qui consiste à faire réagir un polyamino-alcényl- ou alkyl-succinimide avec un composé de formule générale dans laquelle W représente l'oxygène ou le soufre ; X représente l'oxygène ou le soufre ; R ₄ représente un groupe alkylène ayant 2 ou 3 atomes de carbone, facultativement substitué avec 1 à 3 groupes alkyle ayant chacun 1 ou 2 atomes de carbone ; et R ₅ représente l'hydrogène ou un groupe alkyle ayant 1 à 20 atomes de carbone.
2	An additive as claimed in Claim 1, wherein R ₄ is alkylene of 2 or 3 carbon atoms.	Additif suivant la revendication 1, dans lequel R ₄ représente un groupe alkylène ayant 2 ou 3 atomes de carbone.
3	An additive as claimed in Claim 1 or 2, wherein R ₅ is hydrogen or alkyl of from 1 to 10 carbon atoms.	Additif suivant la revendication 1 ou 2, dans lequel R ₅ représente l'hydrogène ou un groupe alkyle ayant 1 à 10 atomes de carbone.
4	An additive as claimed in Claim 1, 2 or 3, wherein W and X are both oxygen.	Additif suivant la revendication 1, 2 ou 3, dans lequel W et X représentent l'un et l'autre l'oxygène.
5	An additive as claimed in Claim 1, 2 or 3, wherein W is sulfur and X is oxygen.	Additif suivant la revendication 1, 2 ou 3, dans lequel W représente le soufre et X représente l'oxygène.
6	An additive as claimed in Claim 1, 2 or 3, wherein W and X are both sulfur.	Additif suivant la revendication 1, 2 ou 3, dans lequel W et X représentent l'un et l'autre le soufre.
7	An additive as claimed in any preceding claim, wherein the reaction is conducted at from 0°C to 250°C.	Additif suivant l'une quelconque des revendications précédentes, dans lequel la réaction est conduite à une température de 0°C à 250°C.
8	An additive as claimed in any preceding claim, wherein the molar charge of the compound of Formula I to the basic nitrogen of the polyamino moiety of the polyaminoalkenyl or alkyl succinimide is in the range from 0.2:1 to 5:1.	Additif suivant l'une quelconque des revendications précédentes, dans lequel le rapport molaire de la charge de composé de formule I à l'azote basique du groupement polyamino du polyamino-alcényl- ou alkyl-succinimide est compris dans l'intervalle de 0,2:1 à 5:1.
9	A lubricating oil composition comprising an oil of lubricating viscosity and 0.2 to 10 percent by weight of an additive as claimed in any preceding claim.	Composition d'huile lubrifiante, comprenant une huile de viscosité propre à la lubrification et 0,2 à 10 % en poids d'un additif suivant l'une quelconque des revendications précédentes.
10	A fuel composition comprising a hydrocarbon boiling in the gasoline or diesel range and an additive as claimed in any one of Claims 1 to 8.	Composition de combustible comprenant un hydrocarbure bouillant dans la plage de l'essence ou du combustible diesel et un additif suivant l'une quelconque des revendications 1 à 8.
11	A fuel composition as claimed in Claim 10, wherein the additive is present in an amount of from 10 to 10,000 weight parts per million.	Composition de combustible suivant la revendication 10 dans laquelle l'additif est présent en une quantité de 10 à 10 000 parties par million, en poids.
12	An output apparatus for outputting information to a recording medium of the kind which comprises an ink sheet (9) and output means for outputting information to a recording medium (2) through the ink sheet (9), a first feed mechanism (33,24,1) for feeding said recording medium (2), a second feed mechanism (14, 25, 30, 31, 53, 75) for feeding the ink sheet (9), a third feed mechanism (14, 25, 30, 31, 54, 58, 113-115, 67) for feeding an ink correction sheet (11), means for driving the first, second and third feed mechanisms, characterised in that the said driving means comprises a single motor there being first transfer means (20, 24, 25, 27) for transferring the driving force generated by said motor selectively to (a) said first feed mechanism (33, 24, 1) and (b) to one of said second feed mechanism (14, 25, 30, 31, 53, 75) and said third feed mechanism (14, 25, 30, 31, 54, 58,	Appareil périphérique de sortie destiné à délivrer en sortie une information vers un support d'enregistrement du type qui comprend une feuille encreuse (9) et des moyens de sortie destinés à délivrer en sortie une information vers un support d'enregistrement (2) à travers la feuille encreuse (9), un premier mécanisme d'alimentation (33, 24, 1) destiné à l'alimentation dudit support d'enregistrement (2), un deuxième mécanisme d'alimentation (14, 25, 30, 31, 53, 75) destiné à l'alimentation de la feuille encreuse (9), un troisième mécanisme d'alimentation (14, 25, 30, 31, 54, 58, 113-115, 67) destiné à l'alimentation d'une feuille encreuse de correction (11), des moyens destinés à entraîner les premier, deuxième et troisième mécanismes d'alimentation, caractérisé en ce que lesdits moyens d'entraînement comprennent un moteur unique, des premiers moyens de transfert (20, 24, 25, 27) étant destinés à transférer

No	Anglais	Français
	113-115, 67); and second transfer means (73, 57b) for switching the driving force generated by said motor (19) between said second feed mechanism and said third feed mechanism in response to a change in direction of rotation of said motor (19).	la force d'entraînement générée par ledit moteur sélectivement vers (a) ledit premier mécanisme d'alimentation (33, 24, 1) et (b) l'un dudit deuxième mécanisme d'alimentation (14, 25, 30, 31, 53, 75) et dudit troisième mécanisme d'alimentation (14, 25, 30, 31, 54, 58, 113-115, 67) ; et des seconds moyens de transfert (73, 57b) destinés à commuter la force d'entraînement générée par ledit moteur (19) entre ledit deuxième mécanisme d'alimentation et ledit troisième mécanisme d'alimentation en réponse à un changement de sens de la rotation dudit moteur (19).
13	An output apparatus according to claim 1, wherein said first transfer means comprises a clutch mechanism (20, 64, 25, 27) with a serrated engagement member (25b).	Appareil périphérique de sortie selon la revendication 1, dans lequel lesdits premiers moyens de transfert comprennent un mécanisme d'embrayage (20, 64, 25, 27) comportant un élément strié de prise (25b).
14	An output apparatus according to claim 1, wherein said third feed mechanism also is operable to shift the correction ink sheet (11) vertically relative to the recording medium (2).	Appareil périphérique de sortie selon la revendication 1, dans lequel ledit troisième mécanisme d'alimentation peut également être mis en oeuvre pour décaler la feuille encreuse de correction (11) verticalement par rapport au support d'enregistrement (2).
15	An apparatus according to claim 1, further comprising key means (5) for inputting data and control means (201) for controlling said apparatus in accordance with the data input by said key means.	Appareil selon la revendication 1, comportant en outre des moyens à touche (5) destinés à introduire des données et des moyens de commande (201) destinés à commander ledit appareil conformément aux données introduites par lesdits moyens à touche.
16	An apparatus according to claim 4, wherein said control means (201) controls said motor (19) to make an additional rotation by a predetermined amount of angle of rotation immediately after the change in the direction of rotation of said motor.	Appareil selon la revendication 4, dans lequel lesdits moyens de commande (201) commandent ledit moteur (19) pour produire une rotation supplémentaire d'un angle de rotation d'une valeur prédéterminée immédiatement après le changement du sens de rotation dudit moteur.
17	An apparatus according to claim 4, wherein said control means (201) is divided into a keyboard controller (202), an apparatus controller and a print controller (204).	Appareil selon la revendication 4, dans lequel lesdits moyens de commande (201) sont divisés en un dispositif de commande (202) de clavier, en un dispositif de commande d'appareil et en un dispositif (204) de commande d'impression.
18	An apparatus according to claim 1, wherein said first transfer means comprises a clutch mechanism (24, 25, 27) and a solenoid (20) and wherein when said solenoid is in its off state, said clutch mechanism is connected to said second and third feed mechanism but not to said first feed mechanism.	Appareil selon la revendication 1, dans lequel lesdits premiers moyens de transfert comprennent un mécanisme d'embrayage (24, 25, 27) et une bobine (20) et dans lequel, lorsque ladite bobine est dans son état hors circuit, ledit mécanisme d'embrayage est relié auxdits deuxième et troisième mécanismes d'alimentation, mais non audit premier mécanisme d'alimentation.
19	An apparatus according to claim 5, wherein the predetermined amount of angle of rotation is about 15°.	Appareil selon la revendication 5, dans lequel la valeur prédéterminée de l'angle de rotation est d'environ 15°.
20	An apparatus according to claim 5, wherein the direction of rotation of said motor (19) is a direction in which the ink sheet (9) is fed in response to said second transfer means (73, 57b) and said second feed mechanism (14, 25, 30, 31, 53, 75).	Appareil selon la revendication 5, dans lequel le sens de la rotation dudit moteur (19) est un sens dans lequel la feuille encreuse (9) est avancée en réponse auxdits seconds moyens de transfert (73, 57b) et audit deuxième mécanisme d'alimentation (14, 25, 30, 31, 53, 75).
21	An apparatus according to claim 1, wherein said apparatus can erase one character by hitting it by a hammer (10) through the ink correction sheet (11) a plurality of times and said third feed mechanism (14, 25, 30, 31, 54, 58, 113-115, 67) moves the ink correction sheet (11) upward or downward during the hitting operation.	Appareil selon la revendication 1, dans lequel ledit appareil peut effacer un caractère en le frappant plusieurs fois à l'aide d'un marteau (10) à travers la feuille encreuse (11) de correction et ledit troisième mécanisme d'alimentation (14, 25, 30, 31, 54, 58, 113-115, 67) déplace la feuille encreuse de correction (11) vers le haut ou vers le bas durant l'opération de frappe.
22	An apparatus according to claim 1, wherein said first feed mechanism (1, 24, 33) feeds said recording medium (2) in a forward or reverse direction in response to the direction of	Appareil selon la revendication 1, dans lequel ledit premier mécanisme d'alimentation (1, 24, 33) fait avancer ledit support d'enregistrement (2) vers l'avant ou en sens inverse en réponse

No	Anglais	Français
	rotation of said motor (19).	au sens de la rotation dudit moteur (19).
23	An apparatus according to any previous claim wherein said motor (19) is a stepper motor.	Appareil selon l'une quelconque des revendications précédentes, dans lequel ledit moteur (19) est un moteur pas à pas.
24	Apparatus for generating from an input analog signal having first and second signal elements that repeat at a first frequency a repetitive signal having a predetermined phase and frequency relationship to the first signal element, the input analog signal being digitized by an analog-to-digital converter (6) under control of the repetitive signal to generate a succession of digital words, the apparatus also including: means for detecting (4) the second signal element within the input analog signal to initially produce a write window; means for storing (14) the succession of digital words from the analog-to-digital converter during the write window; means for processing (16) the succession of digital words from the storing means to generate a control word, and a variable digitally controlled oscillator (8) responsive to said control word to generate said repetitive signal with said predetermined phase and frequency relationship to said first signal element, characterised by address control means (12) arranged to control said storing means (14) to cause successive digital words to be written from said analog-to-digital converter (6) to said storing means during the write window, said write window encompassing said first and second signal elements, and to cause the written digital words to be read from said storing means to said processing means (16) when said write window is closed, said processing means being arranged to adjust the position of said write window in response to the digital words from the storing means and to adjust said control word to maintain said predetermined phase and frequency relationship between said repetitive signal and said first signal element.	Appareil destiné à produire, à partir d'un signal analogique d'entrée possédant des premier et deuxième éléments de signal qui se répètent à une première fréquence, un signal répétitif ayant une relation prédéterminée de phase et de fréquence avec le premier élément de signal, le signal analogique d'entrée étant mis sous forme numérique par un convertisseur analogique-numérique (6) sous commande du signal répétitif de façon à produire une succession de mots numériques, l'appareil comportant également : un moyen servant à détecter (4) le deuxième élément de signal à l'intérieur du signal analogique d'entrée afin de produire initialement une fenêtre d'écriture ; un moyen servant à emmagasiner (14) la succession de mots numériques venant du convertisseur analogique-numérique pendant la fenêtre d'écriture ; un moyen servant à traiter (16) la succession de mots numériques venant du moyen d'emmagasinage afin de produire un mot de commande, et un oscillateur à commande numérique variable (8) qui répond audit mot de commande en produisant ledit signal répétitif suivant ladite relation prédéterminée de phase et de fréquence avec ledit premier élément de signal, caractérisé par : un moyen de commande d'adresse (12) destiné à commander ledit moyen d'emmagasinage (14) afin de faire écrire, dudit convertisseur analogique-numérique (6) audit moyen d'emmagasinage, des mots numériques successifs pendant la fenêtre d'écriture, ladite fenêtre d'écriture renfermant lesdits premier et deuxième éléments de signal, et à faire lire, les mots numériques écrits dudit moyen d'emmagasinage audit moyen de traitement (16), lorsque ladite fenêtre d'écriture est fermée, ledit moyen de traitement étant destiné à ajuster la position de ladite fenêtre d'écriture en réponse aux mots numériques venant du moyen d'emmagasinage et à ajuster ledit mot de commande afin de maintenir ladite relation prédéterminée de phase et de fréquence entre ledit signal répétitif et ledit premier élément de signal.
25	Apparatus as claimed in claim 1 wherein said storing means (14) is a random access memory and said processing means (16) is a microprocessor.	Appareil selon la revendication 1, où ledit moyen d'emmagasinage (14) est une mémoire vive et ledit moyen de traitement (16) est un microprocesseur.
26	A semipermeable microcompartment which is artificially prepared by reassembly of proteinaceous macromolecules and which is defined by a peripheral membrane consisting substantially of a layer of said macromolecules, each of which comprises a relatively hydrophilic moiety and a relatively hydrophobic moiety and wherein the majority of such macromolecules forming the membrane are disposed with their relatively hydrophilic moieties orientated outwardly from the microcompartment and their relatively hydrophobic moieties orientated inwardly towards the interior of the microcompartment.	Microcompartment semiperméable, qui est préparé artificiellement par réassemblage de macromolécules protéiques et qui est défini par une membrane périphérique composée substantiellement d'une couche desdites macromolécules, dont chacune comprend un motif relativement hydrophile et un motif relativement hydrophobe, caractérisé en ce que la majorité de ces molécules, formant la membrane, sont disposées avec leur motif relativement hydrophile orienté vers l'extérieur du microcompartment et leur motif relativement hydrophobe orienté vers l'intérieur du microcompartment.
27	A microcompartment according to claim 1 wherein substantially all of said macromolecules are disposed as aforesaid.	Microcompartment selon la revendication 1 ou 2, caractérisé en ce qu'il présente la forme générale d'une sphère.
28	A microcompartment according to claim 1 or claim 2 which has the shape generally of a sphere. A microcompartment according to claim 1 or claim 2 which has the shape generally of a partially collapsed sphere.	Microcompartment selon la revendication 1 ou 2, caractérisé on ce qu'il présente la forme générale d'une sphère en partie affaissée.
29	A microcompartment according to claim 1 or claim 2 which has the shape generally of a	Microcompartment selon la revendication 1 ou 2, caractérisé en ce qu'il présente la forme

No	Anglais	Français
	sphere inserted to fit into the central space of an annular disc.	générale d'une sphère insérée de façon à être adaptée à l'espace central d'un disque annulaire.
30	A microcompartment according to any of claims 3 to 5 which has an overall diameter in the range of from about 0.1 to about 100 microns.	Microcompartment selon l'une quelconque des revendications 3 à 5, caractérisé en ce qu'il présente un diamètre total compris entre environ 0,1 et environ 100 µm.
31	A microcompartment according to any of claims 3 to 6 which has a wall thickness in the range of from about 0.01µm (100Å) to about 0,1µm (1000Å).	Microcompartment selon l'une quelconque des revendications 3 à 6, caractérisé en ce qu'il a une épaisseur de paroi comprise entre environ 0,01 µm (100 Å) et environ 0,1 µm (1000 Å).
32	A microcompartment according to any of the preceding claims wherein the said macromolecules comprise protein molecules.	Microcompartment selon l'une quelconque des revendications précédentes, caractérisé en ce que lesdites macromolécules comprennent des molécules de protéines.
33	A microcompartment according to any of claims 1 to 7 wherein the said macromolecules comprise glycoprotein molecules.	Microcompartment selon l'une quelconque des revendications 1 à 7, caractérisé en ce que lesdites macromolécules comprennent des molécules de glycoprotéines.
34	A microcompartment according to any of claims 1 to 7 wherein the said macromolecules comprise at least two of protein, glycolipid and glycoprotein molecules.	Microcompartment selon l'une quelconque des revendications 1 à 7, caractérisé en ce que lesdites macromolécules comprennent au moins deux molécules de protéines, de glycolipides et de glycoprotéines.
35	A microcompartment according to any of claims 8 to 10 wherein the membrane is derived from materials comprised in a naturally occurring cell membrane.	Microcompartment selon l'une quelconque des revendications 8 à 10, caractérisé en ce que la membrane est dérivée de substances contenues dans une membrane cellulaire d'origine naturelle.
36	A microcompartment according to any of claims 8 to 10 wherein the membrane is derived from eukaryotic cells or cells of prokaryotes.	Microcompartment selon l'une quelconque des revendications 8 à 10, caractérisé en ce que la membrane est dérivée de cellules eucaryotes ou de cellules procaryotes.
37	A microcompartment according to any of claims 8 to 10 wherein the membrane is derived from materials comprised in whole blood, red blood cell membranes, casein or egg white constituents.	Microcompartment selon l'une quelconque des revendications 8 à 10, caractérisé en ce que la membrane est dérivée de substances contenues dans le sang entier, dans les membranes d'hématies, dans la caséine ou dans des constituants du blanc de l'oeuf.
38	A microcompartment according to any of claims 11 to 13 wherein the membrane contains about 90% of the proteinaceous materials in a naturally occurring cell membrane.	Microcompartment selon l'une quelconque des revendications 11 à 13, caractérisé en ce que la membrane contient environ 90 % de substances protéiques d'une membrane cellulaire d'origine naturelle.
39	A microcompartment according to any of the preceding claims wherein the membrane encloses phospholipids.	Microcompartment selon l'une quelconque des revendications précédentes, caractérisé en ce que la membrane renferme des phospholipides.
40	A microcompartment according to any of the preceding claims wherein the membrane encloses one or more magnetic particles.	Microcompartment selon l'une quelconque des revendications précédentes, caractérisé en ce que la membrane renferme une ou plusieurs particules magnétiques.
41	A microcompartment according to any of the preceding claims wherein the membrane encloses an inert bead.	Microcompartment selon l'une quelconque des revendications précédentes, caractérisé en ce que la membrane renferme un bourrelet inerte.
42	A microcompartment according to any of claims 1 to 17 wherein the membrane encloses an antibody.	Microcompartment selon l'une quelconque des revendications 1 à 17, caractérisé en ce que la membrane renferme un anticorps.
43	A microcompartment according to any of claims 1 to 17 wherein the membrane encloses a naturally-occurring or artificially produced proteinaceous substance.	Microcompartment selon l'une quelconque des revendications 1 à 17, caractérisé en ce que la membrane renferme une substance protéique d'origine naturelle ou produite artificiellement.
44	A microcompartment according to any of claims 1 to 17 wherein the membrane encloses an enzyme or an apo-enzyme.	Microcompartment selon l'une quelconque des revendications 1 à 17, caractérisé en ce que la membrane renferme une enzyme ou une apo-enzyme.
45	A microcompartment according to claim 20 wherein the membrane encloses an apo-enzyme together with a corresponding coenzyme.	Microcompartment selon la revendication 20, caractérisé en ce que la membrane renferme une apo-enzyme en même temps qu'une co-enzyme correspondante.
46	A microcompartment according to any of the preceding claims wherein the membrane surface presents a specific ligand.	Microcompartment selon l'une quelconque des revendications précédentes, caractérisé en ce que la surface de la membrane présente un ligand spécifique.

No	Anglais	Français
47	A microcompartment according to claim 22 wherein the membrane surface presents one or more substances selected from the group consisting of a hormone receptor, an antigen, an antibody and an enzyme.	Microcompartment selon la revendication 22, caractérisé en ce que la surface de la membrane présente une ou plusieurs substances choisies dans le groupe comprenant un récepteur d'hormones, un antigène, un anticorps et une enzyme.
48	A microcompartment according to any of the preceding claims wherein the membrane is comprised of pharmaceutically compatible materials and encloses a pharmacologically active substance.	Microcompartment selon l'une quelconque des revendications précédentes, caractérisé en ce que la membrane est formée de substances pharmaceutiquement compatibles et renferme une substance pharmacologiquement active.
49	A microcompartment according to claim 24 wherein the pharmacologically active substance has anti-inflammatory properties.	Microcompartment selon la revendication 24, caractérisé en ce que la substance pharmacologiquement active a des propriétés anti-inflammatoires.
50	A microcompartment according to claim 25 wherein the pharmacologically active substance is a steroid.	Microcompartment selon la revendication 25, caractérisé en ce que la substance pharmacologiquement active est un stéroïde.
51	A microcompartment according to claim 24 wherein the pharmacologically active substance has anticancer properties.	Microcompartment selon la revendication 24, caractérisé en ce que la substance pharmacologiquement active a des propriétés anticancéreuses.
52	A microcompartment according to claim 27 wherein the pharmacologically active substance is cisplatin.	Microcompartment selon la revendication 27, caractérisé en ce que la substance pharmacologiquement active est le cisplatine.
53	A microcompartment according to claim 27 wherein the pharmacologically active substance is doxorubicin.	Microcompartment selon la revendication 27, caractérisé en ce que la substance pharmacologiquement active est la doxorubicine.
54	A microcompartment according to claim 24 wherein the pharmacologically active substance has central nervous system activity.	Microcompartment selon la revendication 24, caractérisé en ce que la substance pharmacologiquement active a une activité sur le système nerveux central.
55	A microcompartment according to claim 24 wherein the pharmacologically active substance has peripheral nervous system activity.	Microcompartment selon la revendication 24, caractérisé en ce que la substance pharmacologiquement active a une activité sur le système nerveux périphérique.
56	A microcompartment according to claim 24 wherein the pharmacologically active substance has analgetic activity.	Microcompartment selon la revendication 24, caractérisé en ce que la substance pharmacologiquement active a une activité analgésique.
57	A microcompartment according to claim 24 wherein the pharmacologically active substance has local anesthetic activity.	Microcompartment selon la revendication 24, caractérisé en ce que la substance pharmacologiquement active a une activité anesthésique locale.
58	A microcompartment according to claim 24 wherein the pharmacologically active substance has narcotic activity.	Microcompartment selon la revendication 24, caractérisé en ce que la substance pharmacologiquement active a une activité narcotique.
59	A microcompartment according to claim 24 wherein the pharmacologically active substance has antidepressant activity.	Microcompartment selon la revendication 24, caractérisé en ce que la substance pharmacologiquement active a une activité antidépressive.
60	A microcompartment according to claim 24 wherein the pharmacologically active substance has antibacterial activity.	Microcompartment selon la revendication 24, caractérisé en ce que la substance pharmacologiquement active a une activité antibactérienne.
61	A microcompartment according to claim 24 wherein the pharmacologically active substance has antifungal activity.	Microcompartment selon la revendication 24, caractérisé en ce que la substance pharmacologiquement active a une activité antifongique.
62	A microcompartment according to claim 37 wherein the pharmacologically active substance is Amphotericin B.	Microcompartment selon la revendication 37, caractérisé en ce que la substance pharmacologiquement active est l'amphotéricine B.
63	A microcompartment according to claim 24 wherein the pharmacologically active substance has antiparasitic activity.	Microcompartment selon la revendication 24, caractérisé en ce que la substance pharmacologiquement active a une activité antiparasitaire.
64	A microcompartment according to claim 24 wherein the pharmacologically active substance has properties beneficial in the treatment of heart disease.	Microcompartment selon la revendication 24, caractérisé en ce que la substance pharmacologiquement active a des propriétés qui sont bénéfiques dans le traitement des maladies cardiaques.

No	Anglais	Français
65	A microcompartment according to claim 24 wherein the pharmacologically active substance has immunomodulating activity.	Microcompartment selon la revendication 24, caractérisé en ce que la substance pharmacologiquement active a une activité immunorégulatrice.
66	A microcompartment according to claim 41 wherein the pharmacologically active substance has immunostimulating activity.	Microcompartment selon la revendication 41, caractérisé en ce que la substance pharmacologiquement active a une activité immunostimulante.
67	A microcompartment according to claim 24 wherein the pharmacologically active substance has immunosuppressive activity.	Microcompartment selon la revendication 24, caractérisé en ce que la substance pharmacologiquement active a une activité immunodépressive.
68	A microcompartment according to claim 41 wherein the pharmacologically active substance is a leukotriene.	Microcompartment selon la revendication 41, caractérisé on ce que la substance pharmacologiquement active est une leucotriène.
69	A microcompartment according to claim 41 wherein the pharmacologically active substance is an interleukin.	Microcompartment selon la revendication 41, caractérisé en ce que la substance pharmacologiquement active est une interleukine.
70	A microcompartment according to claim 43 wherein the pharmacologically active substance is Cyclosporin A.	Microcompartment selon la revendication 43, caractérisé en ce que la substance pharmacologiquement active est une cyclosporine A.
71	A microcompartment according to claim 24 wherein the pharmacologically active substance is a vaccine.	Microcompartment selon la revendication 24, caractérisé en ce que la substance pharmacologiquement active est un vaccin.
72	A microcompartment according to any of claims 1 to 15 wherein the membrane is composed of edible materials and encloses a foodstuff or a substance compatible therewith.	Microcompartment selon l'une quelconque des revendications 1 à 15, caractérisé en ce que la membrane est composée de substances comestibles et renferme un aliment ou une substance compatible avec lesdites substances.
73	A microcompartment according to claim 48 wherein the membrane encloses a food flavoring substance.	Microcompartment selon la revendication 48, caractérisé en ce que la membrane renferme une substance alimentaire aromatisante.
74	A microcompartment according to any of claims 1 to 17 wherein the membrane encloses a substance with herbicidal, fungicidal, acaricidal or insecticidal activity.	Microcompartment selon l'une quelconque des revendications 1 à 17, caractérisé en ce que la membrane renferme une substance ayant une activité herbicide, fongicide, acaricide ou insecticide.
75	A microcompartment according to any of claims 1 to 17 wherein the membrane encloses a substance with growth control activity.	Microcompartment selon l'une quelconque des revendications 1 à 17, caractérisé en ce que la membrane renferme une substance ayant une activité de régulation de la croissance.
76	A microcompartment according to any of claims 1 to 17 wherein the membrane encloses a pheremone.	Microcompartment selon l'une quelconque des revendications 1 à 17, caractérisé en ce que la membrane renferme une phéromone.
77	A microcompartment according to any of claims 1 to 17 wherein the membrane encloses one or more hydrophobic substances.	Microcompartment selon l'une quelconque des revendications 1 à 17, caractérisé en ce que la membrane renferme une ou plusieurs substances hydrophobes.
78	A microcompartment according to any of claims 1 to 17 wherein the membrane encloses a perfume composition or perfume concentrate composition.	Microcompartment selon l'une quelconque des revendications 1 à 17, caractérisé en ce que la membrane renferme une composition à base de parfum ou une composition à base de concentré de parfum.
79	A microcompartment according to any of claims 1 to 17 wherein the membrane encloses a particulate catalyst.	Microcompartment selon l'une quelconque des revendications 1 à 17, caractérisé en ce que la membrane renferme un catalyseur particulaire.
80	A microcompartment according to any of claims 1 to 21 and 55 wherein the membrane encloses one or more reactant macromolecular substances of molecular weight of at least about 6000.	Microcompartment selon l'une quelconque des revendications 1 à 21 et 55, caractérisé en ce que la membrane renferme une ou plusieurs substances macromoléculaires capables de réagir, d'un poids moléculaire d'au moins environ 6.000.
81	A method of effecting a chemical reaction wherein one or more reactant macromolecular substances of molecular weight of at least about 6000 which is enclosed in a microcompartment as defined in any of claims 1 to 21 and 55 is reacted in a suitable	Procédé de réalisation d'une réaction chimique, caractérisé en ce qu'une ou plusieurs substances macromoléculaires réactives, d'un poids moléculaire d'au moins environ 6.000, comprise(s) dans un microcompartment selon l'une quelconque des revendications 1 à 21 et

No	Anglais	Français
	medium with a substance of molecular weight of up to about 1000.	55, est/sont mise(s) en réaction dans un milieu approprié avec une substance d'un poids moléculaire allant jusqu'à environ 1.000.
82	A microcompartment according to claim 20 wherein the enzyme or apoenzyme is selected from the group consisting of alkaline phosphatase, asparaginase, catalase, cholesterol oxidase, cholinesterase, apo-glucoseoxidase, glucoseoxidase, peroxidase, urease, glycerolphosphate oxidase and uricase.	Microcompartment selon la revendication 20, caractérisé en ce que l'enzyme ou l'apo-enzyme est choisie dans le groupe comprenant la phosphatase alcaline, l'asparaginase, la catalase, la cholestérol-oxydase, la cholestérol-estérase, l'apoglucose-oxydase, la glucose-oxydase, la peroxydase, l'uréease, la glycérolphosphate-oxydase et l'uricase.
83	A microcompartment according to claim 58 wherein the enzyme or apo-enzyme is glucoseoxidase.	Microcompartment selon la revendication 58, caractérisé en ce que l'enzyme ou l'apo-enzyme est la glucose-oxydase.
84	A process for preparing microcompartments as claimed in any of claims 1 to 14, which comprises exposing a proteinaceous macromolecular substance, which contains a relatively hydrophilic moiety and a relatively hydrophobic moiety, to the solubilization action of a homogeneous mixture of chaotropic ions in water and a dialyzable organic solvent which is not completely miscible with water, if necessary or desired subjecting the mixture to filtration or centrifugation, forming microglobules containing the microcompartment precursors, and dialyzing out the organic solvent and the chaotropic ions.	Procédé de préparation de microcompartiments selon l'une quelconque des revendications 1 à 14, caractérisé en ce qu'il comprend les étapes consistant à soumettre une substance protéique macromoléculaire, qui comporte un motif relativement hydrophile et un motif relativement hydrophobe, à l'action de solubilisation d'un mélange homogène d'ions chaotropiques dans de l'eau et d'un solvant organique dialysable qui n'est pas complètement miscible avec l'eau, si nécessaire, ou si on le souhaite, soumettre le mélange à une filtration ou à une centrifugation, former des microglobules contenant les précurseurs du microcompartment, et éliminer le solvant organique et les ions chaotropiques par dialyse.
85	A process for preparing microcompartments as defined in any of claims 15 to 56, 58 and 59, which comprises exposing a proteinaceous macromolecular substance, which contains a relatively hydrophilic moiety and a relatively hydrophobic moiety, to the solubilization action of a homogeneous mixture of chaotropic ions in water and a dialyzable organic solvent which is not completely miscible with water, if necessary or desired subjecting the mixture to filtration or centrifugation, forming microglobules containing the microcompartment precursors, and dialyzing out the organic solvent and the chaotropic ion, the process being additionally characterized by the features that any soluble or insoluble substances which it is desired be enclosed by the said microcompartments are either present in the original reaction mixture or are added after the optional filtration or centrifugation step.	Procédé de préparation de microcompartiments selon l'une quelconque des revendications 15 à 56, 58 et 59, caractérisé en ce qu'il comprend les étapes consistant à soumettre une substance protéique macromoléculaire, qui comporte un motif relativement hydrophile et un motif relativement hydrophobe, à l'action de solubilisation d'un mélange homogène d'ions chaotropiques dans de l'eau et d'un solvant organique dialysable qui n'est pas complètement miscible avec l'eau, si nécessaire, ou si on le souhaite, soumettre le mélange à une filtration ou à une centrifugation, former des microglobules contenant les précurseurs du microcompartment, et éliminer le solvant organique et les ions chaotropiques par dialyse, le procédé étant caractérisé en outre en ce que n'importe quelles substances solubles ou insolubles qu'on souhaite inclure dans lesdits microcompartiments sont contenues dans le mélange réactionnel d'origine ou on les ajoute après l'éventuelle étape de filtration ou de centrifugation.
86	A process according to Claim 60 or Claim 61, wherein the dialysis step is effected in a dialysis bag.	Procédé selon la revendication 60 ou 61, caractérisé en ce que l'étape de dialyse est réalisée dans un sac à dialyse.
87	A process according to any one of Claims 60 to 62, wherein prior to dialysis the mixture is subjected to vigorous agitation to ensure fine division of the microglobules.	Procédé selon l'une quelconque des revendications 60 à 62, caractérisé en ce qu'on soumet le mélange, avant la dialyse, à une agitation énergique pour garantir une fine division des microglobules.
88	A process according to Claim 63 wherein the vigorous agitation is effected by ultrasonic means.	Procédé selon la revendication 63, caractérisé en ce qu'on effectue l'agitation énergique par des moyens à ultrasons.
89	A process according to any one of Claims 60 to 64, wherein the organic solvent is an aliphatic alcohol containing 4, 5 or 6 carbon atoms.	Procédé selon l'une quelconque des revendications 60 à 64, caractérisé en ce que le solvant organique est un alcool aliphatique comportant 4, 5 ou 6 atomes de carbone.
90	A process according to Claim 65, wherein the aliphatic alcohol is n-butanol.	Procédé selon la revendication 65, caractérisé en ce que l'alcool aliphatique est du n-butanol.
91	A process according to any one of Claims 60 to 66, wherein the chaotropic ion is SCN^- or CCl_3COO^- .	Procédé selon l'une quelconque des revendications 60 à 66, caractérisé en ce que l'ion chaotrope est le SCN^- ou le CCl_3COO^- .

No	Anglais	Français
92	A process according to any one of Claims 60 to 67, using substantially equal parts by volume of the organic solvent and of an aqueous solution of chaotropic ions of concentration in the range 20 to 50% w/v.	Procédé selon l'une quelconque des revendications 60 à 67, caractérisé en ce qu'on utilise des parties en volume substantiellement égales du solvant organique et d'une solution aqueuse d'ions chaotropiques à une concentration comprise entre 20 et 50 % p/v.
93	A process according to any one of Claims 60 to 68, wherein, depending upon whether the organic solvent is one the water-solubility of which decreases or increase with a rise in temperature, microglobule formation is effected by respectively raising or lowering the temperature, and/or by diluting the solution with water.	Procédé selon l'une quelconque des revendications 60 à 68, caractérisé en ce que, selon que le solvant organique est un solvant dont la solubilité dans l'eau augmente ou diminue lorsque la température augmente, on réalise la formation des microglobules respectivement en augmentant ou en abaissant la température, et/ou en diluant la solution avec de l'eau.
94	A process for preparing microcompartments in which the peripheral fabric is formed of proteinaceous macromolecules having their relatively hydrophobic and relatively hydrophilic moieties orientated in the reverse sense from that defined in Claim 1, which comprises exposing a macromolecular substance that contains a relatively hydrophilic moiety and a relatively hydrophobic moiety to the solubilization action of a homogeneous mixture of chaotropic ions in water and a dialyzable organic solvent which is not complete miscible with water, if necessary or desired subjecting the mixture to filtration of centrifugation, forming microglobules containing the microcompartment precursors by the addition of a substantially water-immiscible organic solvent and the subsequent creation of a water-in-oil emulsion, and dialyzing out the organic solvent and the chaotropic ions into a substantially immiscible organic solvent.	Procédé de préparation de microcompartiments dans lesquels le tissu périphérique est formé de macromolécules protéiques ayant leur motif relativement hydrophobe et leur motif relativement hydrophile orientés dans le sens inverse à celui défini dans la revendication 1, caractérisé en ce qu'il comprend les étapes consistant à soumettre une substance macromoléculaire, qui contient un motif relativement hydrophile et un motif relativement hydrophobe, à l'action de solubilisation d'un mélange homogène d'ions chaotropiques dans de l'eau et d'un solvant organique dialysable qui n'est pas complètement miscible avec l'eau, si nécessaire ou si on le souhaite, soumettre le mélange à une filtration ou à une centrifugation, former des microglobules contenant les précurseurs du microcompartment, par addition d'un solvant organique substantiellement non miscible avec l'eau et par formation subséquente d'une émulsion eau-dans-huile, et éliminer par dialyse le solvant organique et les ions chaotropiques dans un solvant organique substantiellement non miscible avec l'eau.
95	A process according to Claim 70 wherein the substantially immiscible organic solvent is n-decanol.	Procédé selon la revendication 70, caractérisé en ce que le solvant organique substantiellement non miscible est le n-décanol.
96	A microcompartment produced by the process of any one of Claims 60 to 71.	Microcompartment préparé au moyen du procédé selon l'une quelconque des revendications 60 à 71.
97	A pharmaceutical composition comprising as active ingredient one or more pharmacologically active substances enclosed in microcompartments as defined in any one of Claims 24 to 47, together with a carrier or diluent.	Composition pharmaceutique comprenant, comme ingrédient actif, une ou plusieurs substances pharmacologiquement actives qui sont contenues dans des microcompartiments selon l'une quelconque des revendications 24 à 47, conjointement avec un excipient ou un diluant.
98	A pharmaceutical composition according to Claim 73, wherein at least part of the carrier or diluent is enclosed in the microcompartments together with the one or more pharmacologically active substances.	Composition pharmaceutique selon la revendication 73, caractérisée en ce qu'au moins une partie de l'excipient ou du diluant est contenue dans les microcompartiments, conjointement avec la ou les substances pharmacologiquement actives.
99	A pharmaceutical composition according to Claim 73 or Claim 74, in unit dosage form.	Composition pharmaceutique selon la revendication 73 ou 74, caractérisée en ce qu'elle se présente sous forme de doses unitaires.
100	A pharmaceutical composition according to any one of Claims 73 to 75, which is constituted in a form suitable for oral, parenteral or rectal administration, or for administration by insufflation.	Composition pharmaceutique selon l'une quelconque des revendications 73 à 75, caractérisée en ce qu'elle se présente sous une forme appropriée à une administration par voie buccale, parentérale ou rectale, ou à une administration par insufflation.

Annexe 14 : 20 bisegments anglais-français en format TMX

```
<?XML VERSION="1.0" ENCODING="UTF-8" ?>
<!DOCTYPE tmx SYSTEM "tmx14.dtd">
<tmx version="1.4">
  <header creationtool="LF Aligner" creationtoolversion="3.11" datatype="patent document"
  segtype="sentence" adminlang="EN" srclang="EN" o-tmf="TW4Win 2.0 Format">
  </header>
  <body>
    <tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>
    <tuv xml:lang="EN"><seg>An additive for a lubricating oil or hydrocarbon fuel obtainable by a process which comprises reacting a polyamino alkenyl or alkyl succinimide with a compound of the general formula:
    wherein W is oxygen or sulfur; X is oxygen or sulfur; R4 is an alkylene group of 2 or 3 carbon atoms optionally substituted by from 1 to 3 alkyl groups of 1 or 2 carbon atoms each; and R5 is hydrogen or alkyl of from 1
    to 20 carbon atoms.</seg></tuv>
    <tuv xml:lang="FR"><seg>Additif pour une huile lubrifiante ou un combustible hydrocarboné, pouvant être obtenu par un procédé qui consiste à faire réagir un polyamino-alcényl- ou alkyl-succinimide avec un
    composé de formule générale dans laquelle W représente l&apos;oxygène ou le soufre ; X représente l&apos;oxygène ou le soufre ; R4 représente un groupe alkylène ayant 2 ou 3 atomes de carbone, facultativement
    substitué avec 1 à 3 groupes alkyle ayant chacun 1 ou 2 atomes de carbone ; et R5 représente l&apos;hydrogène ou un groupe alkyle ayant 1 à 20 atomes de carbone.</seg></tuv> </tu>
    <tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>
    <tuv xml:lang="EN"><seg>An additive as claimed in Claim 1, wherein R4 is alkylene of 2 or 3 carbon atoms.</seg></tuv>
    <tuv xml:lang="FR"><seg>Additif suivant la revendication 1, dans lequel R4 représente un groupe alkylène ayant 2 ou 3 atomes de carbone.</seg></tuv> </tu>
    <tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>
    <tuv xml:lang="EN"><seg>An additive as claimed in Claim 1 or 2, wherein R5 is hydrogen or alkyl of from 1 to 10 carbon atoms.</seg></tuv>
    <tuv xml:lang="FR"><seg>Additif suivant la revendication 1 ou 2, dans lequel R5 représente l&apos;hydrogène ou un groupe alkyle ayant 1 à 10 atomes de carbone.</seg></tuv> </tu>
    <tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>
    <tuv xml:lang="EN"><seg>An additive as claimed in Claim 1, 2 or 3, wherein W and X are both oxygen.</seg></tuv>
    <tuv xml:lang="FR"><seg>Additif suivant la revendication 1, 2 ou 3, dans lequel W et X représentent l&apos;un et l&apos;autre l&apos;oxygène.</seg></tuv> </tu>
    <tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>
    <tuv xml:lang="EN"><seg>An additive as claimed in Claim 1, 2 or 3, wherein W is sulfur and X is oxygen.</seg></tuv>
    <tuv xml:lang="FR"><seg>Additif suivant la revendication 1, 2 ou 3, dans lequel W représente le soufre et X représente l&apos;oxygène.</seg></tuv> </tu>
    <tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>
    <tuv xml:lang="EN"><seg>An additive as claimed in Claim 1, 2 or 3, wherein W and X are both sulfur.</seg></tuv>
    <tuv xml:lang="FR"><seg>Additif suivant la revendication 1, 2 ou 3, dans lequel W et X représentent l&apos;un et l&apos;autre le soufre.</seg></tuv> </tu>
    <tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>
    <tuv xml:lang="EN"><seg>An additive as claimed in any preceding claim, wherein the reaction is conducted at from 0°C to 250°C.</seg></tuv>
    <tuv xml:lang="FR"><seg>Additif suivant l&apos;une quelconque des revendications précédentes, dans lequel la réaction est conduite à une température de 0°C à 250°C.</seg></tuv> </tu>
    <tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>
    <tuv xml:lang="EN"><seg>An additive as claimed in any preceding claim, wherein the molar charge of the compound of Formula I to the basic nitrogen of the polyamino moiety of the polyaminoalkenyl or alkyl
```

succinimide is in the range from 0.2:1 to 5:1.</seg></tuv>

<tuv xml:lang="FR"><seg>Additif suivant l'une quelconque des revendications précédentes, dans lequel le rapport molaire de la charge de composé de formule I à l'azote basique du groupement polyamino du polyamino-alcényl- ou alkyl-succinimide est compris dans l'intervalle de 0,2:1 à 5:1.</seg></tuv> </tu>

<tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>

<tuv xml:lang="EN"><seg>A lubricating oil composition comprising an oil of lubricating viscosity and 0.2 to 10 percent by weight of an additive as claimed in any preceding claim.</seg></tuv>

<tuv xml:lang="FR"><seg>Composition d'huile lubrifiante, comprenant une huile de viscosité propre à la lubrification et 0,2 à 10 % en poids d'un additif suivant l'une quelconque des revendications précédentes.</seg></tuv> </tu>

<tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>

<tuv xml:lang="EN"><seg>A fuel composition comprising a hydrocarbon boiling in the gasoline or diesel range and an additive as claimed in any one of Claims 1 to 8.</seg></tuv>

<tuv xml:lang="FR"><seg>Composition de combustible comprenant un hydrocarbure bouillant dans la plage de l'essence ou du combustible diesel et un additif suivant l'une quelconque des revendications 1 à 8.</seg></tuv> </tu>

<tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>

<tuv xml:lang="EN"><seg>A fuel composition as claimed in Claim 10, wherein the additive is present in an amount of from 10 to 10,000 weight parts per million.</seg></tuv>

<tuv xml:lang="FR"><seg>Composition de combustible suivant la revendication 10 dans laquelle l'additif est présent en une quantité de 10 à 10 000 parties par million, en poids.</seg></tuv> </tu>

<tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>

<tuv xml:lang="EN"><seg>An output apparatus for outputting information to a recording medium of the kind which comprises an ink sheet (9) and output means for outputting information to a recording medium (2) through the ink sheet (9), a first feed mechanism (33,24,1) for feeding said recording medium (2), a second feed mechanism (14, 25, 30, 31, 53, 75) for feeding the ink sheet (9), a third feed mechanism (14, 25, 30, 31, 54, 58, 113-115, 67) for feeding an ink correction sheet (11), means for driving the first, second and third feed mechanisms, characterised in that the said driving means comprises a single motor there being first transfer means (20, 24, 25, 27) for transferring the driving force generated by said motor selectively to (a) said first feed mechanism (33, 24, 1) and (b) to one of said second feed mechanism (14, 25, 30, 31, 53, 75) and said third feed mechanism (14, 25, 30, 31, 54, 58, 113-115, 67); and second transfer means (73, 57b) for switching the driving force generated by said motor (19) between said second feed mechanism and said third feed mechanism in response to a change in direction of rotation of said motor (19).</seg></tuv>

<tuv xml:lang="FR"><seg>Appareil périphérique de sortie destiné à délivrer en sortie une information vers un support d'enregistrement du type qui comprend une feuille encreuse (9) et des moyens de sortie destinés à délivrer en sortie une information vers un support d'enregistrement (2) à travers la feuille encreuse (9), un premier mécanisme d'alimentation (33, 24, 1) destiné à l'alimentation dudit support d'enregistrement (2), un deuxième mécanisme d'alimentation (14, 25, 30, 31, 53, 75) destiné à l'alimentation de la feuille encreuse (9), un troisième mécanisme d'alimentation (14, 25, 30, 31, 54, 58, 113-115, 67) destiné à l'alimentation d'une feuille encreuse de correction (11), des moyens destinés à entraîner les premier, deuxième et troisième mécanismes d'alimentation, caractérisé en ce que lesdits moyens d'entraînement comprennent un moteur unique, des premiers moyens de transfert (20, 24, 25, 27) étant destinés à transférer la force d'entraînement générée par ledit moteur sélectivement vers (a) ledit premier mécanisme d'alimentation (33, 24, 1) et (b) l'un dudit deuxième mécanisme d'alimentation (14, 25, 30, 31, 53, 75) et dudit troisième mécanisme d'alimentation (14, 25, 30, 31, 54, 58, 113-115, 67) ; et des seconds moyens de transfert (73, 57b) destinés à commuter la force d'entraînement générée par ledit moteur (19) entre ledit deuxième mécanisme d'alimentation et ledit troisième mécanisme d'alimentation en réponse à un changement de sens de la rotation dudit moteur (19).</seg></tuv> </tu>

<tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>

<tuv xml:lang="EN"><seg>An output apparatus according to claim 1, wherein said first transfer means comprises a clutch mechanism (20, 64, 25, 27) with a serrated engagement member (25b).</seg></tuv>

<tuv xml:lang="FR"><seg>Appareil périphérique de sortie selon la revendication 1, dans lequel lesdits premiers moyens de transfert comprennent un mécanisme d'embrayage (20, 64, 25, 27) comportant un élément strié de prise (25b).</seg></tuv> </tu>

<tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>

<tuv xml:lang="EN"><seg>An output apparatus according to claim 1, wherein said third feed mechanism also is operable to shift the correction ink sheet (11) vertically relative to the recording medium (2).</seg></tuv>

<tuv xml:lang="FR"><seg>Appareil périphérique de sortie selon la revendication 1, dans lequel ledit troisième mécanisme d'alimentation peut également être mis en oeuvre pour décaler la feuille encreuse de correction (11) verticalement par rapport au support d'enregistrement (2).</seg></tuv> </tu>

<tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>

<tuv xml:lang="EN"><seg>An apparatus according to claim 1, further comprising key means (5) for inputting data and control means (201) for controlling said apparatus in accordance with the data input by said key means.</seg></tuv>

<tuv xml:lang="FR"><seg>Appareil selon la revendication 1, comportant en outre des moyens à touche (5) destinés à introduire des données et des moyens de commande (201) destinés à commander ledit appareil conformément aux données introduites par lesdits moyens à touche.</seg></tuv> </tu>

<tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>

<tuv xml:lang="EN"><seg>An apparatus according to claim 4, wherein said control means (201) controls said motor (19) to make an additional rotation by a predetermined amount of angle of rotation immediately after the change in the direction of rotation of said motor.</seg></tuv>

<tuv xml:lang="FR"><seg>Appareil selon la revendication 4, dans lequel lesdits moyens de commande (201) commandent ledit moteur (19) pour produire une rotation supplémentaire d'un angle de rotation d'une valeur prédéterminée immédiatement après le changement du sens de rotation dudit moteur.</seg></tuv> </tu>

<tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>

<tuv xml:lang="EN"><seg>An apparatus according to claim 4, wherein said control means (201) is divided into a keyboard controller (202), an apparatus controller and a print controller (204).</seg></tuv>

<tuv xml:lang="FR"><seg>Appareil selon la revendication 4, dans lequel lesdits moyens de commande (201) sont divisés en un dispositif de commande (202) de clavier, en un dispositif de commande d'appareil et en un dispositif (204) de commande d'impression.</seg></tuv> </tu>

<tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>

<tuv xml:lang="EN"><seg>An apparatus according to claim 1, wherein said first transfer means comprises a clutch mechanism (24, 25, 27) and a solenoid (20) and wherein when said solenoid is in its off state, said clutch mechanism is connected to said second and third feed mechanism but not to said first feed mechanism.</seg></tuv>

<tuv xml:lang="FR"><seg>Appareil selon la revendication 1, dans lequel lesdits premiers moyens de transfert comprennent un mécanisme d'embrayage (24, 25, 27) et une bobine (20) et dans lequel, lorsque ladite bobine est dans son état hors circuit, ledit mécanisme d'embrayage est relié auxdits deuxième et troisième mécanismes d'alimentation, mais non audit premier mécanisme d'alimentation.</seg></tuv> </tu>

<tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>

<tuv xml:lang="EN"><seg>An apparatus according to claim 5, wherein the predetermined amount of angle of rotation is about 15°.</seg></tuv>

<tuv xml:lang="FR"><seg>Appareil selon la revendication 5, dans lequel la valeur prédéterminée de l'angle de rotation est d'environ 15°.</seg></tuv> </tu>

<tu creationdate="20131025T135933Z" creationid="clef-ip-2011"><prop type="Txt::Note">CLEF-claim_en-fr.p1.en-CLEF-claim_en-fr.p1.fr</prop>

<tuv xml:lang="EN"><seg>An apparatus according to claim 5, wherein the direction of rotation of said motor (19) is a direction in which the ink sheet (9) is fed in response to said second transfer means (73, 57b) and said second feed mechanism (14, 25, 30, 31, 53, 75).</seg></tuv>

<tuv xml:lang="FR"><seg>Appareil selon la revendication 5, dans lequel le sens de la rotation dudit moteur (19) est un sens dans lequel la feuille encreuse (9) est avancée en réponse auxdits seconds moyens de transfert (73, 57b) et audit deuxième mécanisme d'alimentation (14, 25, 30, 31, 53, 75).</seg></tuv> </tu>