

ICT Update

Un bulletin d'alerte pour l'agriculture ACP



<http://ictupdate.cta.int>

Un synthétiseur vocal ouvre le monde de l'information aux paysans **kenyans**

La musique traditionnelle **nigériane**, source d'inspiration pour la reconnaissance vocale

Les langues locales au cœur d'un projet **panafricain** d'adaptation des technologies



Technologies du langage

- 2 **Éditorial**
Arriver à se comprendre
- 3 **Perspective**
La diversité linguistique
Adama Samassékou
- 4 **Dossier**
Accès vocal à l'information
Roger Tucker
- Étude de cas**
- 7 **Raviver une langue**
Jon Corbett, Tim Kulchyski et Tom Hukari
- 8 **Trouver le bon ton**
Túndé Adégbölá
- 10 **En idiome local, s.v.p.**
Don Osborn
- TechTip**
- 11 **Une technologie qui parle toutes les langues**
- Q&R**
- 12 **Trouver les mots justes**
Steven Bird

Arriver à se comprendre

Depuis des siècles, l'homme multiplie les techniques de communication : les tambours pour prévenir le village d'à côté, le télégraphe pour envoyer des messages aux pays voisins, la radio pour parler au-delà des océans. Ils ont aujourd'hui fait place à Internet, aux courriels et aux textos. Chacun peut ainsi communiquer malgré les distances, à la condition toutefois de comprendre la langue de l'autre et de maîtriser la technologie. De nombreux efforts sont aujourd'hui déployés pour proposer ces technologies dans d'autres langues, afin que d'autres populations puissent s'en servir et les adapter à leurs propres besoins.

Ces initiatives visent notamment à ouvrir de nouveaux horizons aux analphabètes. Plus de 40% de la population des pays ACP ne sait ni lire ni écrire, un pourcentage qui s'accroît dans les zones rurales où se pose aussi un problème d'accès aux technologies. La difficulté consiste ici à fournir des informations pertinentes – sans doute déjà imprimées ou disponibles sur le web – à des gens qui habitent dans des zones reculées sous une forme compréhensible pour eux.

Et pourquoi pas un système qui, à la demande, lirait automatiquement et à haute voix un texte pré-encodé ? La « synthèse vocale » existe déjà dans de nombreuses langues, surtout européennes et asiatiques. Une poignée d'organisations a donc décidé d'investir du temps et de l'argent dans le développement de synthétiseurs vocaux pour d'autres langues. Dans notre dossier, nous vous parlerons de la mise au point d'un système d'information en kiswahili qui guide et conseille les éleveurs de bananes kenyans. LSSTI, l'organisation à l'origine de cette initiative, espère voir son travail imité dans d'autres régions et s'orienter vers la téléphonie mobile pour transmettre les informations directement et à moindre coût.

Solutions traditionnelles

La vraie communication n'est cependant pas unidirectionnelle. Au Nigeria, Alt-I (African Language Technology) a mis au point un

système qui reconnaît le yoruba et le retranscrit automatiquement en texte. Entre autres applications, ce système de reconnaissance vocale est particulièrement utile pour consigner le savoir traditionnel des aînés du village. Fait intéressant, ce sont les « tambours parlants » de la culture yoruba qui ont inspiré le développement du système de reconnaissance vocale d'Alt-I.

Voilà un parfait exemple d'utilisation du savoir traditionnel pour développer une nouvelle technologie qui, à son tour, contribue à préserver ce savoir. On retrouve la même démarche au Canada, chez HTG (Hul'qumi'num Treaty Group). Moins d'une centaine de personnes parlent encore couramment le hul'qumi'num', mais la communauté est en train de tout mettre sur vidéo pour préserver cette langue en vue des générations futures. Le HTG espère que sa démarche en incitera d'autres, ailleurs, afin de sauvegarder les quelque 3000 autres langues appelées à disparaître au cours de ce siècle.

Tandis que certaines langues disparaissent, d'autres essaient de se frayer un chemin dans le monde des technologies. Le projet PAL (Pan-African Localisation) s'emploie à adapter ou à produire des sites web, des interfaces et des logiciels en langues africaines. Pour l'instant, la plupart des systèmes d'exploitation ou menus déroulants ne s'affichent qu'en anglais ou en français, ce qui en déconcerte plus d'un. Proposer ces outils dans une langue connue et maîtrisée est un moyen d'accroître la participation et le nombre de personnes qui se serviront demain des ordinateurs et de la toile mondiale.

Il y a belle lurette que les auteurs de science-fiction nous dépeignent des ordinateurs qui parlent et qui comprennent tout ce que nous disons. Bien que cela ne soit pas encore tout à fait le cas, il est tout de même encourageant de savoir que des travaux sont en cours pour mettre cette nouvelle technologie au service du plus grand nombre, même de ceux qui ne savent ni lire ni écrire, et surtout, dans la langue qu'ils comprennent le mieux. ■

ICT Update



ICT Update numéro 40, décembre 2007. ICT Update est un magazine multimédia disponible à la fois sur Internet (<http://ictupdate.cta.int>), en version papier et sous forme d'une newsletter diffusée par courriel. Le prochain numéro paraîtra en février.

Le CTA, Centre technique de coopération agricole et rurale (ACP-UE), est un institut du Groupe des États ACP et de l'UE, créé dans le cadre de l'Accord de Cotonou. Il est financé par l'UE. Postbus 380, 6700 AJ Wageningen, Pays-Bas. (www.cta.int).

Production et gestion du contenu Web : Contactivity bv, Stationsweg 28, 2312 AV Leiden, Pays-Bas. (www.contactivity.com)

Coordination rédactionnelle : Rutger Engelhard / Recherche et rédaction : Jim Dempsey / Copy editing : Valerie Jones / Conception de magazine : Frissewind (www.frissewind.nl) / Réalisation graphique : Robert-Jan Cornet / Traduction : Patrice Deladrière / Photo de couverture : Jorgen Schytte - Still Pictures / Conseillers scientifiques : Peter Ballantyne, Oumy Ndiaye, Dorothy Okello, Kevin Painting
Copyright : ©2007 CTA, Wageningen, the Netherlands

<http://ictupdate.cta.int>





Adama Samassékou (asamassekou@acalan.org) est secrétaire exécutif de l'Académie africaine des langues (www.acalan.org) et président du réseau Maaya (<http://maayajo.org>)

de 90 %), le manque de diversité est préoccupant, sachant qu'il existe dans le monde plus de 6 000 langues.

La place des langues africaines (kiswahili, haoussa, yoruba, mandingue, peul...) dans les technologies est tellement marginale qu'elle n'apparaît pas dans les études. Un continent comme l'Afrique est inexistant sur Internet. L'objectif de l'accès universel, réaffirmé

En supposant que les États et les décideurs nationaux et internationaux mettent en place les stratégies pour sauvegarder ces langues, les promouvoir auprès des populations par l'alphabétisation généralisée, la technologie numérique offre un certain nombre d'opportunités pour leur intégration dans le cyberspace. Le développement de la norme Unicode favorise l'informatisation des langues et leur intégration dans les logiciels bureautiques et les progiciels de gestion. Si le monde du logiciel libre a entrepris très tôt la production de logiciels dans les langues des différentes communautés linguistiques, on peut se réjouir que des entreprises comme Microsoft commencent à s'y intéresser. Certaines de ces démarches sont sans doute commerciales, mais, pour les puissances publiques, cela devrait relever d'un devoir de souveraineté ou de service public. Par ailleurs, la possibilité de création des noms de domaine internationaux, qui permettent l'utilisation des caractères des langues reconnues sur Internet dans les noms de domaine) est un progrès indéniable. Il est heureux qu'au Forum sur la gouvernance de l'Internet, mis en place à la suite du SMSI, la question des noms de domaine internationaux cristallise souvent les débats.

Les technologies numériques peuvent également être mises à profit pour développer l'intercompréhension linguistique. Le développement des outils de traduction automatique et de reconnaissance vocale devrait davantage focaliser l'attention, comme cela a été souligné durant la dernière réunion de la Coalition dynamique sur la diversité culturelle, organisée dans le cadre du FGI, à Rio au Brésil. Les participants ont regretté que les puissances publiques et la communauté des développeurs du libre n'investissent pas assez le domaine du traitement automatique des langues, bien qu'il soit au cœur de l'accès universel. S'il est vrai que les travaux autour du langage naturel en informatique, ainsi que la mise en place de corpus et dictionnaires pour créer ces outils linguistiques numériques, sont assez coûteux, les États devraient davantage s'y intéresser.

Osons espérer que l'année 2008, qui a été proclamée « Année internationale des langues » par l'Assemblée générale des Nations unies, sera une période essentielle pour le développement des langues marginalisées dans tous les pays du monde. ■

La diversité linguistique

La diversité linguistique est à la société humaine ce qu'est la biodiversité à la nature. C'est l'enjeu fondamental de la problématique langues et technologie. Nul ne doit être exclu de cette nouvelle société de la connaissance et des savoirs partagés en construction.

La diversité linguistique demeure une quête fondamentale, déterminant le futur des cultures humaines, en particulier dans la nouvelle ère numérique qui s'amplifie. Selon des études menées par l'association Funredes en 2006 et l'Union latine, l'anglais est la langue la plus utilisée dans les pages Web (45 %), loin devant, par exemple, le français (4,95 %) et l'espagnol (4,60 %), tendance corroborée par les études de InternetWorldStats, en 2007. La diversité culturelle évolue néanmoins, puisque l'anglais, selon les études de Funredes, occupait 75 % des contenus en 1998. Toutefois, en dehors des 10 à 12 langues largement parlées sur Internet (à hauteur

lors du Sommet mondial sur la société de l'information (SMSI), ne saurait toutefois se résumer à l'accès au matériel informatique de connexion. Les contenus auxquels on accède, notamment développés par le citoyen, grâce au Web 2.0, sont la seconde face indispensable d'une même médaille. Ce qui suppose que les usagers utilisent la langue qu'ils parlent, non seulement pour exploiter les contenus en ligne, mais surtout pour en produire.

Il existe même un risque que le développement du cyberspace contribue à la disparition des langues les moins parlées, lesquelles connaissent un déficit de standardisation et de normalisation informatique ou même linguistique (sous la forme écrite). Les politiques linguistiques mises en œuvre autour de ces langues sont souvent freinées par les difficultés de financement de leur industrialisation. Tel est le cas des langues africaines.



VANESSA VICK/THE NEW YORK TIMES/PH

Un paysan kenyan qui cultive pour la première fois les bananes est tout excité au moment les récolter. Lorsqu'il a discuté d'une possible extension de cette nouvelle culture avec le responsable local de la vulgarisation agricole, ce dernier lui a donné un numéro de téléphone où il obtiendrait toutes les informations nécessaires pour le guider à chaque étape du processus. Avec l'arrivée à maturité de son premier lot, le paysan a besoin de conseils pour cueillir les régimes dans les règles de l'art et s'occuper des bananes coupées. Comme la moindre erreur pourrait lui coûter cher, il décide d'appeler le numéro.

Au bout du fil, une voix lui propose de poursuivre la conversation en kiswahili

l'élevage (NALEP) du Ministère de l'agriculture kenyan. Il a connu une phase pilote pendant plusieurs mois en 2006, afin d'étudier l'usage et les applications possibles d'un service vocal d'informations agricoles.

Fait intéressant, l'équipe universitaire ignorait pratiquement tout de la technologie vocale jusque 2004, époque à laquelle le Dr Gakuru a été contacté par la LLSTI parce qu'elle développait un synthétiseur vocal en kiswahili et cherchait un partenaire en Afrique de l'Est, où cette langue est parlée. Les premiers mois, la LLSTI a fourni les outils, la formation et les connaissances permettant au Dr Gakuru de réaliser le synthétiseur vocal utilisé par la Ligne Infos Bananes.

pertinence, enfin : il faut encourager les universitaires à se concentrer sur les problèmes qui affectent leur propre communauté.

L'institut sud-africain Meraka, un des partenaires fondateurs de la LLSTI, est le parfait exemple d'une organisation qui s'est lancée dans le développement de plusieurs nouvelles applications multilingues (comme la ligne d'assistance téléphonique en anglais et en setswana pour les prestataires de soins au Botswana, en phase pilote début 2008). Meraka est en effet devenu un pôle d'excellence en technologie vocale, dispensant des formations et un accompagnement à des chercheurs bien au-delà de ses frontières, jusqu'au Nigeria. L'institut a également élargi sa

Accès vocal à l'information

Les paysans des zones rurales, surtout celles où le niveau d'alphabétisation est faible, ont toujours eu du mal à se tenir au courant des dernières évolutions dans leur domaine. Un système de synthèse vocale accessible via le téléphone portable ouvre désormais le monde de l'information à ces communautés.

ou en anglais. Bien qu'aucune de deux ne soit sa langue maternelle, il choisit le kiswahili qu'il maîtrise mieux. En quelques minutes à peine, le paysan reçoit toutes les informations nécessaires. Au passage, il reconnaît même la voix de son interlocuteur : Ken Walibora, célèbre présentateur de la télé kenyane.

Ce qu'il ignore, c'est que Ken Walibora n'a jamais prononcé un seul mot sur la Ligne Infos Bananes ; les explications ont été générées automatiquement à partir d'un texte, grâce à un synthétiseur vocal développé par une équipe de l'université de Nairobi dirigée par le Dr Mucemi Gakuru. L'immense avantage de ce service, c'est qu'il suffit de modifier des pages web pour l'actualiser. Ken Walibora n'a dû consacrer que 45 minutes de sa vie à lire quelques phrases choisies avec soin. C'est la force de la technologie vocale.

Apport local

La Ligne Infos Bananes est un projet que l'initiative pour la technologie vocale en idiomes locaux (LLSTI) mène en partenariat avec le programme national de vulgarisation de l'agriculture et de

Ce projet est typique des travaux menés par la LLSTI. Cette organisation est une initiative mondiale créée en 2003 par Outside Echo, une ONG britannique qui facilite l'accès vocal aux informations et qui compte des partenaires en Inde, en Afrique du Sud, au Kenya et au Nigeria. La LLSTI aide des équipes n'ayant aucune connaissance préalable en synthèse vocale à réaliser des synthétiseurs vocaux fonctionnels, au débit naturel. Doivent impérativement figurer dans l'équipe un linguiste et un ingénieur en logiciel, de même qu'un chef d'équipe motivé, ingénieur de préférence. À la différence des projets de développement commerciaux où un locuteur rejoint temporairement une équipe d'experts dans un laboratoire européen ou américain de l'entreprise, la LLSTI travaille avec des équipes locales, qui font partie d'un institut de recherche universitaire.

Pourquoi est-il si important de favoriser l'implication locale dans le développement de la technologie vocale ? Pour une question de motivation, tout d'abord : les gens sont très attachés à leur langue. Pour une question de maintenance ensuite : il faut sans cesse adapter la technologie, compte tenu des problèmes rencontrés dans son utilisation. Pour une question de

vision originale pour y ajouter un système de reconnaissance vocale et vise, à long terme, la traduction automatisée. Le système sur lequel il travaille en ce moment, « Lwazi », est un ambitieux système d'information téléphonique en 11 langues avec apports vocaux en entrée et en sortie, commandité par le Ministère sud-africain des arts et de la culture.

Langues difficiles

Cela paraît facile, mais un synthétiseur vocal est un logiciel extrêmement intelligent, qui requiert des connaissances linguistiques particulières, beaucoup de textes et de données audio annotés à la main, et des jugements techniques avisés en fonction de la langue considérée. Au final, il s'agit de développer un système qui lit comme un être humain, ce qui existe déjà pour la plupart des langues européennes. En attendant, on peut utiliser des synthétiseurs vocaux au phrasé artificiel pour peu que le message reste intelligible, ce qui est apparemment plus facile dans certaines langues que d'autres.

Lorsque les voyelles ne sont pas écrites, comme en arabe par exemple, comment le système peut-il deviner la prononciation ? Ou lorsque la langue utilise un accent tonique libre (imprévisible) comme en anglais,



PHEMULA NASFORIC

comment le système peut-il déterminer la syllabe à accentuer ?

Pour appréhender ces problèmes, le projet LLSTI a dès l'origine mené une enquête sur 105 langues afin de recenser pour chacune d'elles les traits écrits et langagiers susceptibles de poser des problèmes à un synthétiseur vocal. Tous les traits liés à la synthèse vocale ont ainsi été repris dans une base de données multilingue, de manière à prédire les problèmes posés par chaque langue et à attribuer un coefficient de difficulté pour le développement d'un synthétiseur vocal [cf. tableau].

Langue	Synthétiseur vocal de base	Synthétiseur vocal évolué
Pashto	9	9,5
Arabe (classique)	7	8,5
Russe	6	9
Tibétain	6	7,5
isiZulu	6	8
Ibibio	5	7
Thaï	5	8
Anglais	4	6
Hindi	2	4
Gallois	1	4
Kiswahili	0	4
Tamil	0	2.5

Coefficient de difficulté pour le développement d'un synthétiseur vocal (0=facile, 10=difficile)

Du point de vue de la synthèse vocale, la langue idéale a un coefficient de difficulté égal à zéro. Il s'agit d'une langue pour laquelle le processus de synthèse vocale peut être défini par un corpus de règles simples que tout linguiste est en mesure de rédiger à partir de ce qu'il sait déjà. En pratique, ces règles ne peuvent jamais totalement définir le processus, mais il est quelques langues où elles permettent d'élaborer un système de base – c.-à-d. avec un rendu imparfait du phrasé, des mots empruntés, des abréviations et autres détails du même genre, mais néanmoins intelligible.

De manière générale, le développement d'un synthétiseur vocal se décline en plusieurs étapes :

- Définir les caractéristiques de la langue : spectre sonore (c.-à-d. l'ensemble des sons utilisés dans la langue), règles de conversion lettre-son, règles de syllabation, etc.
- Choisir un ensemble de phrases phonétiquement équilibrées, à partir d'une large base de données de textes transcrits en phonétique, couvrant l'ensemble des combinaisons sonores de la langue en aussi peu de phrases que possible. Il s'agit d'un processus automatique, mais il faut toujours trouver un compromis entre le nombre de phrases et la couverture

des combinaisons sonores rares.

- Choisir un locuteur. Le choix de la voix est sans doute l'aspect le plus crucial du système. Avec une dimension intuitive. Au fond, tout le monde parvient à se faire comprendre, non ? Mais la synthèse vocale a besoin d'un rendu clair et précis de chacun des mots prononcés dans la base de données : elle construit en effet ses messages en concaténant les mots de la base pour en faire de petits segments. Chaque mot doit donc être prononcé parfaitement. La voix doit en outre être aussi intelligible que possible à l'origine, de sorte que toute baisse de qualité due au processus de concaténation ait le moins d'incidence possible.
- Enregistrer les phrases phonétiquement équilibrées. Pour le kiswahili, cela représente quelque 400 phrases, soit environ 45 minutes d'enregistrement. La plupart des autres langues ont besoin de plus de phrases.
- Ajouter à la main des annotations phonétiques aux enregistrements. Bien que cela puisse se faire de manière automatisée, toute erreur à ce niveau pourrait engendrer des problèmes dans le message synthétisé.

Liens corrélés

Essayez le synthétiseur vocal utilisé pour la Ligne Infos Bananes (en kiswahili)
→ www.llsti.org/demos-interactive.htm.

La base de données linguistique pour synthétiseur vocal de la LLSTI
→ www.llsti.org/languages-database.htm

Le projet linguistique en Iwazi de l'institut Meraka
→ www.meraka.org.za/iwazi/

Exemples de synthétiseurs vocaux européens
→ www.nuance.com/realspeak

Dictionary Maker
→ <http://dictionarymaker.sourceforge.net>

Apprendre la morphologie des langues synthétiques complexess
→ www.cs.bris.ac.uk/Research/MachineLearning/morph



ANDREA MATONE / ALAMY

- Compiler toutes ces données dans un système de synthèse vocale à l'aide de Festival, un progiciel gratuit développé par l'université d'Édimbourg.
- Tester.

Lorsque la définition des règles n'est pas possible (étape 1), d'autres techniques guidées par les données peuvent s'appliquer. Une technique guidée par les données s'appuie sur une source considérable de données textuelles annotées, généralement prises dans le domaine public (sur Internet) lorsque cette source existe dans la langue considérée, et essaie d'en dégager mécaniquement des « règles ». L'annotation manuelle des données est laborieuse et chronophage ; aussi, au fil des ans, les chercheurs se sont-ils évertués à minimiser la part des données nécessitant une annotation précise, l'objectif final étant d'avoir un rendu acceptable sans aucune annotation manuelle. L'institut Meraka, par exemple, a développé un créateur de dictionnaires de prononciation qui utilise une technique itérative pour réaliser l'ensemble du dictionnaire avec un minimum d'interventions.

La majorité des langues exige une analyse morphologique (AM), même pour des synthétiseurs vocaux de base. L'AM consiste à analyser un mot pour en extraire la part vocale. L'AM est propre à chaque langue et généralement très complexe. C'est ce qui fait l'essence même de la plupart

de nos langues. C'est pour cette raison que la LLSTI participe en ce moment à un grand projet de recherche de l'université de Bristol visant à développer un système d'AM automatique, applicable aux langues pour lesquelles on dispose de peu de données linguistiques.

Les mots parlent

Dans le cas de la Ligne Infos Bananes, une évaluation officielle a été menée auprès d'une sélection de dix paysans du district de Kirinyaga. Elle a révélé quelques problèmes intéressants. Sept des dix paysans, par exemple, ont choisi l'anglais, mais ont ensuite eu des problèmes à cause de l'accent britannique. Ceux qui ont choisi le kiswahili ont aimé la voix, mais la grammaire kiswahili formelle utilisée dans la traduction leur a posé problème. Tous ont dit aimer le système vocal et le préférer à des documents écrits, mais il fallait manifestement régler les questions d'accent et de traduction avant d'en élargir l'usage. Le Dr Gakuru a fait écouter aux paysans quelques échantillons du synthétiseur vocal en anglais kenyan qu'il était en train de développer : tous l'ont trouvé plus clair que la version originale en kiswahili. Ce synthétiseur est aujourd'hui pratiquement achevé et sera bientôt mis en service.

Les consultations menées avec le NALEP après ce projet pilote ont permis de dresser un tableau complet des services que devrait offrir une

ligne d'informations agricoles. Des services qui permettraient aux paysans d'obtenir toutes les informations nécessaires, au moment voulu, et dans un langage compréhensible.

Avec l'arrivée des services mobiles de transmission des données en Afrique, il serait désormais possible de faire tourner le synthétiseur de la ligne d'information sur le téléphone proprement dit. Seules les données seraient transférées et c'est le téléphone qui les convertirait en paroles. Tout le monde ne dispose pas d'un téléphone aussi perfectionné, mais pour ceux qui en ont un, ce serait une option très attrayante. Primo, le coût de transmission de l'information est beaucoup plus faible : les appels mobiles sont très chers au Kenya. Secundo, on peut accéder à l'information via un menu visuel, avec recherche de texte pour ceux que cela tente. Tertio, des images et des mots / chiffres clés peuvent être affichés en complément de la voix, ce qui facilite l'assimilation et la mémorisation des informations.

Les technologies d'aujourd'hui le permettent. Il ne manque qu'un synthétiseur vocal dans les langues que les gens pratiquent couramment. Certaines de ces langues posent de réelles difficultés. Il y a donc un vrai challenge à relever, mais une fois ce travail accompli, ces systèmes seront opérationnels. Les Sud-Africains y consacrent des fonds publics, d'autres pays feraient bien de suivre leur exemple. ■

Raviver une langue

Alors que la moitié des langues sont appelées à disparaître au cours de ce siècle, une initiative canadienne les remet au goût du jour grâce au multimédia.

Au Canada, comme dans d'autres anciens pays coloniaux, de nombreuses langues indigènes sont sur le point de disparaître à tout jamais. L'une d'elles est le hul'q'umi'num', une langue parlée par plusieurs communautés du sud de l'île de Vancouver, notamment par les tribus Cowichan, la première nation Chemainus, la tribu Penelakut, les premières nations Lyackson, la première nation Halalt et la première nation Lake Cowichan. Comme il reste aujourd'hui moins d'une centaine de locuteurs hul'q'umi'num', en majorité des personnes du troisième âge, c'est la survie même de la langue qui est en jeu.

Hul'q'umi'num Treaty Group (HTG), une organisation agissant au nom de ces tribus et premières nations, et des chercheurs des universités de Victoria et de Colombie britannique ont uni leurs efforts pour raviver cette langue. Les travaux de leur projet sont supervisés par un conseil consultatif des anciens, composé d'acteurs intéressés et de personnes des six communautés parlant couramment le hul'q'umi'num'. Les chercheurs universitaires s'emploient à créer un module d'apprentissage basé sur le multimédia numérique. Ensemble, ils incitent les membres de la communauté à revenir au hul'q'umi'num' pour la vie quotidienne.

Le multimédia numérique, qui s'appuie sur des ressources audio, visuelles et texte, offre un attrait visuel considérable. Cet atout incite les anciens et les jeunes à créer et à évaluer ensemble les outils d'apprentissage de la langue. La production et la distribution de médias numériques comme les DVD étant désormais plus facile, ceux-ci peuvent largement contribuer à étendre l'usage du hul'q'umi'num' dans les écoles, les cours de langues et les communautés.

Implication

La vidéo participative est une des techniques utilisées par le projet. Celui-ci inclut donc un volet de formation pour

apprendre aux membres de la communauté à filmer et à produire des DVD. Ces membres sont associés à la conception du film, au choix du lieu de tournage et à l'élaboration du contenu. Le conseil consultatif des anciens a tenu à ce que les sujets des films et des DVD ne contribuent pas qu'à revitaliser leur langue, mais aussi les pratiques culturelles traditionnelles. Les DVD abordent par conséquent divers thèmes, dont l'archivage des formes traditionnelles des discours de circonstance et des savoir-faire comme la vannerie à base d'écorce de cèdre.

Lancé en mai 2004, le projet entendait fournir une information linguistique très interactive et d'autres ressources que les membres de la communauté pourraient utiliser à la maison. La plupart des services offerts sur Internet s'accompagnent de contenus multimédias téléchargeables et en vidéo continue, qui obligent les utilisateurs à disposer de connexions à haut débit. Or la plupart des membres de la communauté hul'q'umi'num' vivent dans des réserves où l'accès au haut débit est rare. À l'époque, seuls 10 à 15% des membres de la communauté disposaient d'une connexion à haut débit à la maison, alors que 90% d'entre eux avaient accès à un lecteur DVD. D'où le choix évident du DVD plutôt que d'une technologie Internet.

Les DVD ont ceci de commun avec les technologies Internet qu'ils prévoient une part d'interactivité avec l'utilisateur. On y trouve des exercices de grammaire, de vocabulaire et d'autres outils, au choix de l'utilisateur. Au travers du menu interactif, l'utilisateur choisit le type d'information auquel il veut accéder sur le DVD. Cette dimension de choix est un atout des DVD car il est peu probable que les gens accèdent à l'ensemble du contenu en une fois. De multiples séances sont à prévoir avant d'en faire le tour. Par leur souplesse et leur robustesse, les DVD conviennent donc parfaitement à l'apprentissage des langues.

Nouvelle orientation

Il a récemment été proposé de mettre des séquences vidéo des premiers films sur Internet, afin de mettre la langue hul'q'umi'num' à la portée du plus grand nombre. Cela favoriserait notamment des apports des jeunes, qui se servent de plus en plus des sites sociaux pour se divertir et communiquer. L'augmentation du nombre d'internautes est notamment due à l'augmentation des connexions à haut débit à domicile de même qu'à une offre de services informatiques accessibles au public.

Le projet est en train de développer un blogue, en réalité une application composite (mash-up), qui conjugue des textes, des photos, des vidéos et des bandes son hébergées sur Picasa Web Albums et YouTube. Ce blogue peut facilement être mis à jour et enrichi de nouveaux documents en hul'q'umi'num' pour inciter chacun à y revenir régulièrement. Les utilisateurs peuvent également laisser leurs commentaires, ce qui contribue bien évidemment à l'éclosion d'une communauté en ligne « d'étudiants » en hul'q'umi'num'.

La facilité d'accès à la technologie, la baisse du prix du matériel et la convivialité grandissante des logiciels se traduisent par une prolifération de groupes d'intérêts aussi divers que variés qui se servent des médias numériques pour faire passer des informations culturelles et partager leurs points de vue. Les multimédias numériques se prêtent particulièrement bien à l'apprentissage des langues. Vu le rythme rapide auquel les langues mineures disparaissent, elles offrent une occasion unique de raviver l'intérêt pour une langue et, dans le cas du hul'q'umi'num', de lui éviter de disparaître à jamais. ■



Jon Corbett (jon.corbett@ubc.ca) est professeur assistant à l'université de Colombie britannique, Canada (<http://web.ubc.ca/okanagan/ccgs/welcome.html>). Tim Kulchyski (hemutth@gmail.com) est le coordonnateur hul'q'umi'num' des tribus Cowichan, Colombie britannique, Canada (www.cowichantribes.com). Tom Hukari (hukari@uvic.ca) est professeur émérite à l'université de Victoria (<http://web.uvic.ca/ling>).



DRONE ZECU / ALAMY

Trouver le bon ton

Les « tambours parlants » traditionnels des Yoruba, au Nigéria, sont à l'origine d'une démarche unique dans la reconnaissance vocale. Une technique qui pourrait s'appliquer à d'autres langues africaines et à des programmes d'alphabétisation.

Les gens sont généralement plus efficaces et plus performants lorsqu'ils s'expriment dans leur langue maternelle. Raison pour laquelle les nations qui ont enseigné les sciences et les technologies dans leurs langues maternelles sont celles qui ont le plus profité de l'ère industrielle. A l'opposé, les nations qui continuent de vivre avec les séquelles sociales et culturelles de la colonisation et d'enseigner les sciences et les technologies dans des langues étrangères ont bien du mal à

être compétitives.

Si la langue a eu tant d'importance au cours de l'ère industrielle, elle n'en aura que plus encore à l'ère de l'information. C'est ce constat qui est à l'origine de l'initiative Alt-I (African Languages Technology Initiative), laquelle vise à faire rentrer les cultures africaines dans l'ère de la technologie en développant les ressources nécessaires à l'utilisation de toutes les technologies de l'information et de la communication (TIC) dans les langues africaines.

L'affaire est loin d'être simple. Sur les 2 000 langues parlées en Afrique, il y en a moins de 10 qui comptent plus de 10 millions de locuteurs. Les autres ne sont généralement parlées par quelques dizaines de milliers de personnes.

L'ampleur de la tâche ne doit cependant pas nous décourager. Faire le tour du monde commence toujours par un premier pas.

Source originale

Alt-I a commencé par développer divers aspects de la technologie vocale en prenant le yoruba comme langue pilote parce qu'il est parlé par plus de 25 millions de personnes. Malgré de nombreux dialectes, il a une forme standard, qui a donné naissance à une orthographe standard, qui s'est muée en un système d'écriture stable au cours du dernier siècle.

Notre organisation avait déjà développé une disposition de clavier efficace et ergonomique pour le yoruba,

Túndé Adégbólá (tadegbola@alt-i.org) est directeur exécutif d'Alt-I (www.alt-i.org) et maître de conférences associé à l'Africa Regional Centre for Information Science, Université d'Ibadan, Nigéria (www.arcisng.net)

Nous avons également adapté un kiosque d'information sur le VIH/sida développé en yoruba par l'AIISI (Africa Information Society Initiative, émanation de la Commission économique pour l'Afrique des Nations unies). Ce kiosque prévient les gens qui ne savent ni lire ni écrire des risques liés au VIH/sida et fournit toutes ces informations oralement, en yoruba. Il est utilisé par le personnel des ONG qui travaillent dans le secteur des soins de santé.

Le haut fait d'Alt-I aura été de développer un système de reconnaissance vocale unique et plus efficace pour les langues africaines à ton, où le sens des mots varie en fonction de la « hauteur » (tonalité) de la voix du locuteur. L'identification des tons, en plus de celle des consonnes et des voyelles, est généralement considérée comme un fardeau supplémentaire pour les systèmes de reconnaissance vocale.

Les recherches menées par Alt-I ont démontré au contraire qu'il était possible d'utiliser la tonalité pour accroître la vitesse et l'exactitude des systèmes de reconnaissance vocale. L'implémentation de la reconnaissance vocale, non seulement sur des ordinateurs portables ou de bureau, mais aussi sur des PDA ou des consoles de jeu converties devient ainsi envisageable. Un outil particulièrement utile pour recueillir de précieuses connaissances locales auprès des aînés de la communauté qui n'ont jamais appris à lire ni à écrire. Conjugué à d'autres technologies vocales, comme la synthèse vocale, ces connaissances pourront ensuite être partagées et affinées par d'autres.

Il tapent sur des tambours...

Nous nous sommes rendu compte que la tonalité pouvait être un élément important pour les technologies vocales et surtout pour la reconnaissance vocale après avoir observé des joueurs de « tambours parlants » yoruba. Héritiers d'un savoir ancestral, ces joueurs utilisent leurs tambours comme moyen de communication fonctionnel, en modulant la hauteur du son. Nous en avons déduit que les tons yoruba contenaient suffisamment d'informations pour que ceux qui écoutent le « tambour parlant » comprennent le message, sans consonnes ni voyelles. Pour approfondir la question, nous avons collaboré avec plusieurs membres du clan Ayan, dépositaire de ce savoir. Leur enseignement fut pour nous une expérience enrichissante et d'humilité.

Les connaissances acquises au travers

de ce qui fut essentiellement un apprentissage, furent ensuite passées au crible de théories scientifiques modernes aussi diverses que l'ethnomusicologie, la psychologie, la linguistique, l'informatique et la théorie du codage de l'information. Forts de ces connaissances, nous avons développé un système de reconnaissance vocale en yoruba plus efficace que ceux présentés dans les ouvrages traitant des techniques de reconnaissance des langues humaines.

La science de la reconnaissance vocale ayant débuté par l'anglais, langue atonale, ses efforts se sont essentiellement concentrés sur les reconnaissances des consonnes et des voyelles. Chez Alt-I, nous nous sommes néanmoins aperçu qu'à partir du moment où les tons étaient précisément reconnus, ils pouvaient suggérer les voyelles qui les soutenaient et les consonnes qui les accompagnaient. Fait intéressant, les tons sont nettement plus faciles à reconnaître que les consonnes et les voyelles.

Ces caractéristiques étant communes à de nombreuses langues tonales africaines, des outils et des techniques pourront probablement être développés et partagés pour d'autres langues, y compris asiatiques.

Conscients des besoins pressants de nombreuses autres langues africaines, nous abordons le développement de ce produit de manière à en faciliter l'adaptation pour d'autres langues tonales africaines dans un proche avenir.

À voix haute

Ce travail sur les systèmes de reconnaissance vocale a été entrepris dans le cadre d'un plus gros projet d'Alt-I, baptisé « Redefining Literacy » et financé par l'OSIWA (Open Society Initiative for West Africa, une des composantes du George Soros' Open Society Institute). Ce projet cherche à utiliser la reconnaissance et la synthèse vocales pour permettre à tous les africains, alphabétisés ou non, d'interagir avec des documents écrits. Au travers de cette démarche, Alt-I espère changer la définition de l'alphabétisation, non plus comprise comme « la capacité de lire et d'écrire », mais comme « la capacité d'interagir avec l'écrit ». Redéfinir le degré d'alphabétisation en ces termes permet d'inclure des personnes qui seraient considérées comme illettrés. Elles passent ainsi du statut d'illettré à celui d'e-lettré.

Sur un plan plus personnel, je suis un de ces villageois africains ordinaires dont la famille est passée d'une culture

Le clavier yoruba

Une des premières réalisations d'Alt-I aura été de développer une disposition de clavier efficace et ergonomique pour le yoruba. Pour ce faire, un large corpus de yoruba a été analysé pour déterminer la fréquence d'occurrence de chaque caractère. Partant de là, les caractères les plus fréquents ont été placés aux endroits adéquats du clavier. Celui-ci a été présenté aux imprimeurs locaux et a connu un franc succès. Auparavant, il fallait trois fois plus de temps pour taper une page en yoruba qu'une page en anglais. Cet écart est pratiquement résorbé. Des claviers d'ordinateur ergonomiques comme le nôtre réduisent les problèmes de microtraumatismes répétés consécutifs à un usage intensif de l'ordinateur. Alt-I soutient des démarches identiques pour toutes les langues qui ne disposent pas de clavier.



totalement orale au monde de la technologie en à peine deux générations ; je connais bien les problèmes d'adaptation. Et comme je vis à l'heure de l'informatique, j'en connais toutes les possibilités. Les TIC peuvent permettre à l'Afrique de rester en prise avec l'ère de la connaissance. Il serait dommage de ne pas tirer pleinement profit de cette opportunité. ■

Liens corrélés

International Institute for Communication and Development (IICD)

→ www.iicd.org

Commission économique pour l'Afrique des Nations unies (CEA/UNECA)

→ www.uneca.org

Africa Information Society Initiative (AIISI)

→ www.uneca.org/aisi

Open Society Initiative for West Africa (OSIWA)

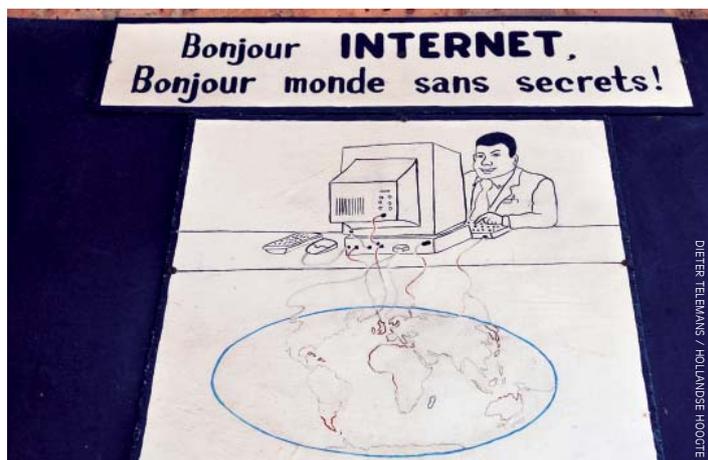
→ www.osiwa.org

George Soros' Open Society Institute

→ www.soros.org

En idiome local, s.v.p.

Diverses initiatives s'emploient à rendre le contenu et les logiciels plus accessibles et plus pertinents pour les populations du continent africain en les adaptant aux langues locales.



L'Afrique est une mosaïque linguistique ; dans la plupart des pays africains, plusieurs langues locales et nationales cohabitent. D'après certaines estimations, on parlerait plus de 2 000 langues en Afrique (selon *Ethnologue: Languages of the World*). Les langues européennes, héritées de l'époque coloniale et conservées pour des usages officiels, compliquent encore la situation.

Dans un contexte multilingue, la (ou les) langue(s) utilisée(s) est (sont) toujours une affaire de choix, mais qui n'est pas sans conséquences. Décider de travailler dans une langue plutôt que dans une autre peut réduire le cercle participatif au sein d'une communauté ou influencer l'usage des savoirs locaux. Une bonne partie de l'action de vulgarisation et de développement dépend en outre d'expatriés ou d'autochtones multilingues pour la traduction, quoique le recours à des intermédiaires ne soit pas toujours pratique pour des projets TIC.

Je dirais même qu'il n'est pas nécessaire. En principe, les TIC peuvent fonctionner dans toutes les langues humaines. Les systèmes doivent en effet pouvoir gérer des caractères spéciaux ou des alphabets non latins. Ajouter des interfaces et des contenus en langues africaines s'avère non seulement possible, mais souhaitable dans la mesure où les gens peuvent plus facilement accéder à la technologie et aux informations dans les

langues qu'ils maîtrisent le mieux. Ces langues ne sont toutefois pas celles qui dominent généralement les TIC.

Soutien panafricain

C'est pour attirer l'attention sur l'usage des langues africaines dans les TIC que l'initiative Bisharat a vu le jour au Mali en 2000. L'embryon d'agenda de ce qu'on appelle aujourd'hui les TIC au service du développement faisait en effet peu de cas des contenus en langues africaines. Pour des langues comme le bambara, dont l'alphabet inclut des lettres latines modifiées, il fallait en outre s'accommoder du problème des « polices spéciales » développées localement. L'exigence était donc double : montrer l'importance de contenus et d'interfaces informatiques en langues africaines et trouver les moyens de surmonter les problèmes techniques liés à l'utilisation d'alphabets tels que celui du bambara.

À l'époque, les forums et le site web de Bisharat ont concentré la majorité des efforts déployés pour résoudre les problèmes liés à l'utilisation des langues africaines dans les TIC. Il s'en est suivi en 2004 une discussion entre le Centre canadien de recherches pour le développement international (CRDI), Bisharat et l'ONG Kabissa. Ces pourparlers ont abouti à un nouveau projet, baptisé PAL (PanAfrican Localisation), dont l'objectif était de promouvoir l'adaptation des technologies au contexte africain en mettant plus particulièrement l'accent sur le développement et l'éducation.

Le projet PAL a débuté en avril 2005, au moment où plusieurs contenus Internet en langues africaines ainsi que divers projets africains d'adaptation faisaient leur apparition. Un atelier PAL tenu à Pretoria, en Afrique du Sud, au mois de novembre 2007, a réuni divers intervenants à la pointe de la localisation des technologies. Les projets présentés allaient de l'adaptation des logiciels en wolof et en kinyarwandais au renforcement des contenus Internet (avec notamment des blogues et un dictionnaire en swahili), en passant par l'élaboration de terminologies dans plusieurs langues comme le lingala et des recherches sur la reconnaissance vocale, notamment en yoruba et en somali.

Si l'on y ajoute les efforts consentis

par Microsoft pour proposer la plupart de ses logiciels Windows et Office en plusieurs langues africaines et l'engagement du projet « One Laptop Per Child » à s'adapter aux langues locales, on voit clairement une tendance se dessiner. En 2008, un nouveau projet PAL financé par le CRDI va s'amorcer sous la direction de Translate.org, afin de développer les éléments clés de la localisation africaine : adaptation des fichiers de données, des dispositions de clavier et de la terminologie, et examen des relations entre langues, adaptation et politiques en matière de TIC.

Boucler la boucle

Tous ces efforts seront jugés à l'aune de leur utilisation sur le terrain. Il faut que les anciens comme les nouveaux projets visant à mettre les TIC au service du développement incorporent ces schémas. Plus il y aura de logiciels disponibles en différentes langues grâce à des initiatives locales ou à des projets comme PAL, moins on pourra trouver de justifications techniques au court-circuitage des langues africaines dans les projets TIC. Les polices de caractères, les claviers voire, pour certaines langues africaines, l'ensemble des applications bureautiques, sont déjà disponibles. Autant d'éléments à valeur ajoutée, gratuits ou à faible coût, qui permettent aux systèmes informatiques de s'exprimer dans la même langue que leurs utilisateurs présumés.

Restent deux obstacles à surmonter : la commercialisation de produits en langues africaines et l'idée reçue selon laquelle pour ajouter une nouvelle capacité linguistique à un ordinateur, il faut forcément en supprimer une autre, généralement l'anglais ou le français. Au final, l'incorporation des langues africaines dans les TIC s'avérera sans doute moins difficile que de changer les mentalités vis-à-vis de l'emploi de ces langues sur la toile ou sur des ordinateurs et que de faire connaître ce qui existe déjà. L'informatique multilingue est une réalité ; comment l'exploiter à bon escient et de manière optimale, telle est la question. ■

Liens corrélés

PanAfrican Localisation

→ www.panafril10n.org

Translate.co.za

→ www.translate.co.za

KiLinux

→ www.kilinux.udsm.ac.tz

Don Osborn (dzo@bisharat.net) est le fondateur et le directeur de Bisharat (www.bisharat.net)

Une technologie qui parle toutes les langues

Plusieurs applications liées aux langues sont gratuites et disponibles sur Internet. Les plus utilisées sont les outils de traduction.

Traduction

Les traducteurs professionnels appliquent la devise de St Jérôme, non verbum e verbo sed sensum exprimere de sensu, ne pas traduire des mots, mais des idées. Leur tâche ne consiste pas à traduire littéralement, mais à restituer le sens du texte original. Il n'en va toutefois pas de même pour les outils de traduction en ligne. Ces applications ne « comprennent » pas le sens du texte. Elles se contentent de le traduire mot à mot. Voyez la citation ci-dessous et comparez les différentes traductions françaises obtenues :

Kindness is the language which the deaf can hear and the blind can see.
Mark Twain, auteur et satiriste américain (1835–1910)

La bonté est la langue que le bidon sourd entend et l'aveugle peut voir.
AltaVista – Babel Fish (<http://babelfish.altavista.com>)

Gentillesse est la langue qui peut entendre les sourds et les aveugles peuvent voir.
Google Translate (www.google.com/translate_t)

La gentillesse est la langue que les sourds peuvent entendre et l'aveugle peut voir.
PROMT online (<http://translation2.paralink.com>)

La bonté est la langue que le sourd peut entendre et les aveugles peuvent voir.
Dictionary.com

La bonté est un langage qu'un sourd peut entendre et qu'un aveugle peut voir.
Traducteur professionnel.

Seules deux traductions donnent un résultat acceptable (PROMT et Dictionary). Les autres rencontrent essentiellement des problèmes d'analyse et de syntaxe: ordre des mots, fonction de that (sujet ou complément), nature des mots (can), erreur de vocabulaire (blind)...

Cet exemple montre qu'on ne peut trop se fier à des logiciels de traduction pour nous restituer le véritable sens d'un texte ; il nous montre aussi combien la réalisation d'applications liées aux langues est difficile. Les langues ne se résument pas à la juxtaposition de mots, et il est bien difficile d'inculquer ce sens de la nuance à un ordinateur.

Lecture et écriture automatisées

Les analphabètes peuvent se servir de divers progiciels (lecteurs d'écran, systèmes de reconnaissance vocale) pour accéder à Internet, rédiger des courriels, voire même créer des documents texte.

Lecteurs d'écran

Les lecteurs d'écran utilisent un synthétiseur vocal pour lire le texte affiché à l'écran d'un ordinateur. Bien que développés initialement pour les aveugles et les malvoyants, ils peuvent également venir en aide à ceux qui ne savent ni lire ni écrire. Dans la plupart des cas, il faut se servir de la souris ou du clavier pour déplacer le curseur vers le texte à lire (y compris des adresses web). Il existe pas mal de lecteurs d'écran de marque, généralement très coûteux. Quelques-uns peuvent être téléchargés gratuitement. Votre ordinateur doit néanmoins être équipé d'une carte son et de haut-parleurs et/ou d'un casque.

NonVisual Desktop Access (NVDA)

NVDA est un gratuiciel qui tourne sous Windows XP et Vista. Il lit à haute voix les textes affichés par Mozilla Firefox, Internet Explorer, Word, Excel et Outlook Express.
www.nvda-project.org

Thunder

Thunder tourne sous Windows XP, 2000 et Vista. Il reconnaît notamment les affichages Word, Excel et Internet Explorer. Son téléchargement est gratuit pour un usage individuel. Les associations et les entreprises doivent négocier un prix d'utilisation.
www.screenreader.net

Reconnaissance vocale

Avec un logiciel de reconnaissance vocale et un micro, les utilisateurs peuvent commander certains programmes, dicter des courriels, des textes et naviguer sur la toile. Les commandes sont généralement simples,

du style « ouvrir programme » ou « enregistrer document », mais le logiciel doit les apprendre et met parfois du temps à s'y faire. La plupart des programmes de reconnaissance vocale sont disponibles en plusieurs langues (européennes), mais certaines voix ou accents les indisposent.

E-speaking

E-speaking tourne sous Windows 2000 et XP, et fonctionne avec les applications Microsoft Office. Après une période d'essai gratuit de 30 jours, vous pouvez acquérir pour 14 dollars une licence de mise à niveau qui vous permet d'ajouter ou de supprimer des commandes utilisateur.
www.e-speaking.com

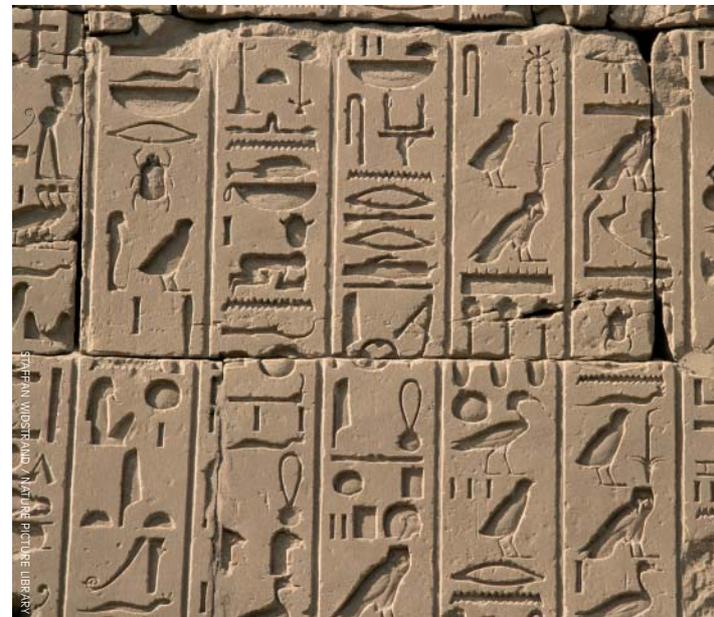
SpeechVibe

Ce programme revient à 15 dollars après une période d'essai gratuit de 30 jours. Il tourne sous Windows 2000, XP et Vista.
<http://speechvibe.com>

Simply Web 2000

Simply Web ne fonctionne qu'avec Internet Explorer (version 4.01 ou supérieure) et permet de naviguer sur la toile grâce à un synthétiseur vocal.
www.econointl.com/sw

Remarque: Il existe un synthétiseur vocal gratuit pour les Mac qui, associé à un plug-in logiciel pour Netscape, est capable de lire les pages web à haute voix.





Steven Bird (sb@csse.unimelb.edu.au) est professeur associé à la Faculté d'informatique et d'ingénierie logicielle de l'université de Melbourne (www.csse.unimelb.edu.au) et chercheur émérite au Consortium des données linguistiques de l'université de Pennsylvanie (www ldc.upenn.edu).

phrases ci-après et déterminez à quoi se rapporte le pronom **eux** :

- Certains d'entre **eux** ont par la suite été vendus.
- Certains d'entre **eux** ont par la suite été arrêtés.
- Certains d'entre **eux** ont par la suite été retrouvés.

Pour comprendre le texte, l'ordinateur doit déterminer qui a fait quoi à qui. Est-ce que ce sont les malfrats ou les bijoux qui ont été vendus, arrêtés ou retrouvés ? Pour répondre à cette question, nous devons nous appuyer sur notre connaissance du monde et sur notre capacité de déduction.

pour analyser plusieurs langues tonales non écrites et créer le dictionnaire de l'une d'entre elles. À partir du moment où une langue locale a un système d'écriture solidement établi, une tradition écrite et l'appui du gouvernement, elle est parée pour les technologies du langage.

On ne trouve généralement que des outils de base (claviers et correcteurs orthographiques) car l'élaboration d'un synthétiseur vocal, d'un moteur de recherche, d'un système de traduction requièrent des investissements considérables. Il y a des technologies utiles qui permettent à des expatriés d'apprendre une langue locale. Transcriber, par exemple, (trans.sourceforge.net) permet d'enregistrer, de transcrire et d'étudier des dialogues entre locuteurs indigènes.

Trouver les mots justes

Quelles sont les problèmes inhérents au développement d'une technologie qui comprend le langage humain ?

→ Des tâches simples, qui sont une seconde nature pour nous et que nous maîtrisons dès notre plus jeune âge, sont étonnamment difficiles à automatiser. La première étape consiste à déterminer les mots prononcés, à « saisir les messages ». Dites à toute vitesse « César l'aimait sage » et vous verrez que la tâche de l'ordinateur n'est pas si simple. Ajoutez-y une variation dialectale et du bruit de fond, et le problème se complique singulièrement. La deuxième étape consiste à déterminer la relation entre les mots et le monde qui nous entoure. Supposons que vous lisiez un article dans lequel apparaît la phrase : « Les malfrats ont dérobé des bijoux ». Imaginez qu'elle soit suivie d'une des trois

ou si s'avère extrêmement compliqué d'inculquer nos connaissances et nos capacités de raisonnement à un ordinateur. La troisième étape consiste à réagir de manière appropriée. Cela va du plus simple (jouer notre air favori ou répondre à une simple question factuelle) au summum de la difficulté (engager une longue discussion ou raisonner de manière critique).

Y a-t-il plusieurs types de technologies du langage ?

→ Il y a essentiellement deux types de technologies : celles qui comblent l'écart entre les hommes et les ordinateurs en proposant des interfaces plus naturelles. C'est le cas de la prédiction de la frappe, de la reconnaissance vocale et de la détection des émotions. Et celles qui comblent l'écart entre la multitude d'informations disponibles sur le web et celles dont nous avons effectivement besoin. Il s'agit notamment des résumés de texte et de la traduction automatique. Ces technologies existent en anglais et pour les langues les plus courantes, grâce à de nombreuses recherches et au financement de leur développement.

Ces types de technologie existent-ils pour la plupart des langues du monde ?

D'après ethnologue.com, il y aurait environ 7 000 langues parlées à travers le monde. Seule la moitié d'entre elles ont une écriture standardisée, et encore moins une tradition écrite. Certains instituts de recherche et organisations de développement participent activement au « développement des langues », en commençant par créer une écriture puis en dispensant des cours d'apprentissage de la lecture et de l'écriture dans la langue maternelle. Dans les années '90, j'ai moi-même travaillé dans l'ouest du Cameroun

Peut-on envisager une approche peu coûteuse permettant de développer ces technologies pour d'autres langues ?

→ Le principal défi consiste à rassembler une masse énorme de textes et de transcriptions audio pour chaque langue, couvrant l'expression libre de nombreux locuteurs. Je m'imagine une sorte de YouTube, assorti de traductions et de transcriptions temporisées, largement utilisé par les locuteurs des langues locales pour consigner leur héritage linguistique. Une sorte de Wikipédia dans chaque langue nous procurerait aussi un gisement de textes couvrant une large gamme de sujets. Ces ressources s'ajouteraient aux nombreuses sources de données indispensables au développement de technologies dans chaque langue. Also, a type of Wikipedia for each language could provide a useful body of text covering a wide range of topics. These materials would add up to a rich source of data to be used in developing technologies for each language. ■

Open Language Archives Community (OLAC)

OLAC est un partenariat international de particuliers et d'institutions qui est en train de créer une librairie virtuelle mondiale des ressources linguistiques : dictionnaires, textes annotés et grammaires, notamment. À ce jour OLAC dispose de plus de 30 dépôts d'archives et de quelque 30 000 items référencés correspondant à plus de 3 000 langues du monde entier. Certains items sont numériques et disponibles en ligne ; d'autres sont des objets physiques auxquels on ne peut accéder qu'en se rendant personnellement dans un dépôt d'archives.

www.language-archives.org



DETER TELEMANS / HOLLANDSE HOOFDTE