



Automatic indexing and reformulation of ancient dictionaries

Abdel Belaïd, Isabelle Turcan, Jean-Marie Pierrel, Yolande Belaïd, Yves Rangoni, Hassen Hadjamar

► To cite this version:

Abdel Belaïd, Isabelle Turcan, Jean-Marie Pierrel, Yolande Belaïd, Yves Rangoni, et al.. Automatic indexing and reformulation of ancient dictionaries. First International Workshop on Document Image Analysis for Libraries - DIAL'2004, Jan 2004, Palo Alto, United States. pp.342-354, 10.1109/DIAL.2004.1263264 . inria-00100087

HAL Id: inria-00100087

<https://hal.inria.fr/inria-00100087>

Submitted on 8 Feb 2007

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.

Automatic indexing and reformulation of ancient dictionaries

A. Belaid¹, I. Turcan^{2,3}, J.M. Pierrel³, Y. Belaid¹, Y. Rangoni¹ and H. Hadjamar³

¹Loria-University of Nancy 2 France {abelaid@loria.fr, ybelaid@loria.fr, rangoni@loria.fr}

²University Jean Moulin, Lyon III France {Turcan@atilf.fr, Isabelle.Turcan@univ-lyon3.fr}

³ATILF UMR-CNRS 7118, Nancy France {Jean-Marie.Pierrel@atilf.fr, hassen.hadjamar@atilf.fr}

Abstract

This paper is related to automatic indexing and reformulation of ancient dictionaries. The objective is to make easy the access to ancient printed documents from XVI to XIX century for a diversified public (historians, scientists, librarians, etc.). Since the facsimile mode is insufficient, the aim is to look further for the use of the indexing based on the formal structure representative of some contents in order to optimize their exploration. Starting from a first indexing experiment operated on more recent documents, the TLF (“Trésor de la Langue Française”: Treasure of the French Language) in the ATILF laboratory (Nancy, France), we extended the indexing technique to automatic reformulation and reedition of ancient dictionaries. However, face to the problem extent, we limited our investigations to a very specific collections of the ATILF laboratory, the “Trévoux” dictionary (defined later).

1 Introduction

1.1 Ancient dictionary valorization

The valorization aim of ancient documents and more specially of ancient dictionaries can be summarized in the following items, already outlined in [1,2,3]:

1. Make available texts that are often presented in complex typography forms, with different levels of diffi-

culty depending on the centuries. In XVII and the beginning of XVIII century, there are problems of page-setting (columns, marginal texts), characters difficult to read (distinction between F and long S, long or double long S (see **Figure 3, 4, 5**), “esperluètes” (&), bindings, sign classifiers, the hand with stretched index (☛) which indicates a new entry in “Trévoux” dictionary and others). In the second half of the XVIII century and XIX century, some other specific difficulties appeared such as character size, reference sign increase, domain name and author of example name abbreviations, multicolumn (often 4 as in the “Grand Dictionnaire Universel” (GDU)), etc. On all the two centuries, and beyond, there is the problem of character restitution in non-Roman character sets.

2. Allow the auto-correction of the digitization errors by drawing up the typology of the errors noted on a representative sample, and to use them for the automatic correction of new occurred errors without any modification of the original text (except the characteristics of the typography: reference letter (see **Figure 6**), ancient characters, etc.).
3. Allow the fast comparison, by alignment, of texts considered “a priori” identical, but comprising differences not easily localizable with the naked eye. It is the case of many dictionary editions and republications of the XIX century, whose texts of escorts (forewords, introductions, lecturer point of view, etc.) and nomenclatures vary.
4. Finally, to define useful and effective indexing modes (definition problem of linguistic key words, lexical, and meta-linguistics, definers, etc.).

5. Unknown characters: the ancient typography presents some specimens which are not known by the OCRs. For example, the letter C represented like a big epsilon (see Figure 7).

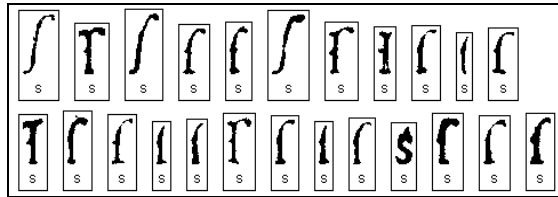


Figure 3: Long S, r and f.

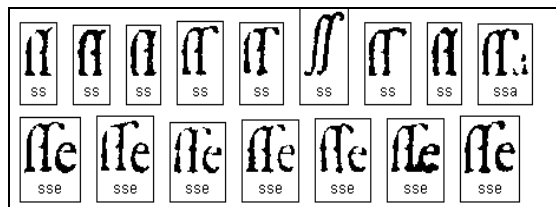


Figure 4: Examples of kerning

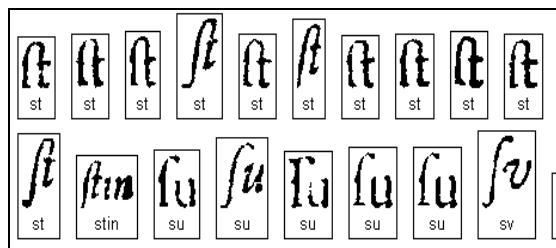


Figure 5: Examples of ligatures

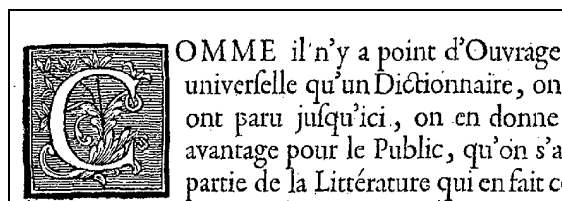


Figure 6: Example of the reference letter C in Trévoux.

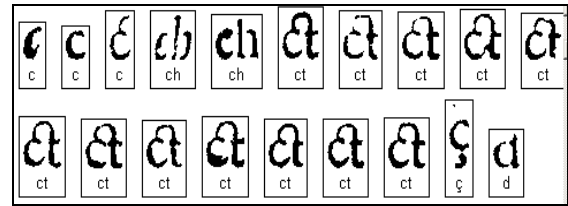


Figure 7: Example of the character C in Trévoux.

3.2 OCR Recognition

The indexing and reformulation procedures need text recognition. For this reason we have looked for an OCR solution adapted to this kind of documents. The presence of specific shapes leads us to consider them as new character classes and present them to OCRs for training. We can see in Figures 3,4,5 and 7 the different shapes.

Furthermore, we noticed that the document typography is very complex and OCRs are often not efficient individually taken. So, we proposed an adapted approach based on OCRs combination [6]. The method consists in the development of an OCRAdapter, and an integration procedure of commercial OCRs as plug-ins. A ground truth is first established from a representative extract of the document class. Then, the first step is dedicated to the individual evaluation of each OCR. A specific alignment algorithm performs the evaluation. We have used the Myers's algorithm, based on an optimal dynamic programming matching [7,8,9].

This evaluation step allowed us to characterize for each OCR its performances and more precisely its drawbacks: bad recognized characters, recognition context considered, etc.

This knowledge on each OCR is taken into account in a second phase during the combination. The combination's philosophy is to consider only the two best OCRs and keep the better according to each kind type of error. As for evaluation, the two OCRs output files are linearized and aligned. The alignment technique is described in [9]. The principle is to choose a reference chain (also called chain of consensus) from which we will initialize the alignment. Then, this alignment is rectified continuously by comparison between the other chains and this reference chain.

The experiments realized on the Trévoux dictionary are still in progress, but similar works, operated on more recent documents, show that, for a basic combination, a score of 20/10000 errors is obtained. Furthermore, considering some specific heuristics, we can easily reach 1/10000 error. We hope that, from the first results ob-

tained on ancient dictionaries, we will obtain identical results. The heuristics employed are issued from the observation of the confusions provided by the OCRs in competition and attempt to correct them. For example, each time that OCR1 is giving "I" (always confused by "1"), we replace this answer by that one given by OCR2 more efficient on the recognition of digits.

4 Dictionary indexing

4.1 Dictionary indexing problems

The ancient dictionaries were neither written nor composed as the current dictionaries [4,5]. They don't follow the same systematic structure principles for the entry definitions. These definitions are made up of informational fields that designate various indication categories which constitute the article microstructure. This, for example, gives grammatical category indications (noun, verb, adjective, article, adverb, conjunction, etc.), pronunciation category indications, orthography difficulties, information relating to the etymology, definitions organized or not according to the semantic logic of the proper sense illustrated or derived, examples, synonyms and antonyms, particular definitions, etc.

4.2 Data indexing difficulties

The systematic absence of the speech structure in these dictionaries causes two major difficulties for the automatic processing [6]:

1. On one hand, the geographic repartition informational fields is irregular, often random, unequally differentiated by formal criteria such as the typography (for example, semiotics of the italic or capitals) or the punctuation (for example, semiotics of the point not followed by a capital letter in an enumeration) and the position in the article layout (information structure according to the subparagraphs);
2. On the other hand, the polysemous statute of certain statements presents a functional ambiguity of the fields. In fact, we can wonder to which informational fields belong the synonymic definitions or the synonyms capable to be interpreted like definitional. Similarly for the definitional etymologies or the etymology definitions, can we know if they concern the definition field or the etymology field?

This structural blur led us to create the concept of keywords and meta-linguistic sequence-keys to facilitate the information access without using a finest labeling. These tools showed their effectiveness in exploration procedures of lexicographical textual bases once imple-

mented in a rigorous way, not by applying arbitrary labels to the text, but by extracting objective information (combinatory of typographical criteria, positional and lexical or phraseological).

In spite of the two difficulties previously stated, it is not utopian to envisage some automatic processing modalities even for a part of the information. This can help to derivate the principle of an automatic content reformulation.

In fact, a first set of investigations carried out on our corpus shows that it is possible to identify within the definitions some specific structures. These structures that correspond to specific information in the microstructure, are easily detectable thanks to the typography and positional indices of some elements.

As it is pointed out previously, this kind of formal identification is however limited to the physical aspect of the information. In fact, it is very difficult to identify the logical structures because this needs interpretation.

Nevertheless, in spite of this reserve, we consider that it will be possible to proceed to an automatic reformulation approach, based only on physical criteria. We are however conscious that in this case, a user can exploit only a part of the articles, such as the synonyms, the examples, the named sources or the grammatical terminologies.

4.3 Indexing methodology

The goal of the reformulation is thus to be able to extract, from a textual set, definite parts of the content, i.e. fields or information in "coded" position, and to reorganize them in different manner than in the original text [10,11].

This is why, considering the work realized these last years concerning the electronic publishing methodology of ancient dictionaries and the possible procedures of exploration and exploitation of their content, we consider that it is important to take into account the following elements:

1. The explicit structure of the lexicographic speech.
2. The main informational fields: typological proposition.

A) The explicit structure of the lexicographic speech

In a first time, we discard all the informational fields functional ambiguity cases whose the logic design is

implicit or registered in an editorial process implying an interpretation accompanied with demonstrations. This excludes any prospect of automatic treatment based on formal criteria specific to the considered texts except if we admit a lexicographical bases preliminary treatment by a fine labeling based on specialists meta-lexicographic competences of these texts.

From a strict formal point of view, three series of combined criteria must be taken into account to determine the first principles of fields location. This is a necessary precondition to any automatic extraction of their contents:

1. **Typographical criteria:** character size (large / small right Roman, large / small capitals, etc.), character status (upper case / lower case), font style (italic / straight Roman, etc).
2. **Positional criteria** of the subparagraphs within the same article and of the information present within the same subparagraph:
 - In the article area: identification and delimitation of the subparagraphs. Each subparagraph must be numbered and profit from a mark of subparagraph beginning / end, for each subparagraph withdrawal / absence of withdrawal.
 - In the area of a subparagraph: all the punctuation marks must be codified (specially the points not followed by a capital letter, commas followed by italic and certain semicolons). The typography (in particular upper case / lower case, straight / Italic), when having a diacritic function, must be codified because it allows to determine the beginning and the end of an informational field or particular sequence identifiable within a field.
3. **Lexical or phraseological criteria:** we find here the logic of keywords and metalinguistic sequence-keys which allow to isolate the grammatical category marks (often shortened in coded position), the field marks (sometimes with alternatives), the synonyms with definitional function (especially locatable from a positional point of view), the formulas specific to the definitional or etymological speech, etc.

B) The main informational fields: typological proposition

The main fields encountered in the ancient dictionaries are as follows:

1. Entries or "vedette".
2. Grammatical category marks.
3. Indications of pronunciation and of graphic variant.

4. Usage, domain, technical, language level marks, etc.
5. Synonymous with definitional function.
6. Elaborated definitions.
7. Translations of foreign languages.
8. Etymologies.
9. Marked bibliographic references (author names cited and concerned entitled texts)
10. Texts cited with and without reference.
11. Intra-textual returns (to other articles within the same piece of work) or inter-textual (to a set of text sources).

We propose a modeling in the form of a table as following:

Criteria	Typographical	Positional	Phascoligical
1.Entry	Capitals	Beginning subparagraph	
2.Pronunciation	Sometimes in Italic or between parentheses	Just after the entry or the field 3	Key-words «on / se prononce»
3. Grammar Cat.	Italic	Just after the entry or the field 2	Typology of the abbreviation of metalinguistic Keywords
4.Domain Cat.	Straight Roman	After fields 2 or 3	Metalinguistic keywords: «term(s) de», «mot de», «en» + domain name
5.Definition or definitional synonyms	Straight Roman	After fields 3 or 4	a) Metalinguistic keywords, «sorte de», «spèce de», «signific, veut dire, comme quand on dit, etc.» b) Def. by synonymous (hyperonymous), or derivated.
6.Examples / ref. & quotations	Sometimes in Italic or between brackets	After field 5	The author names and work titles can be listed
7.Etymology	The etymology is sometimes in italic or in straight roman in small capitals	Random, sometimes in position 1, sometimes after fields 2-6	Metalinguist. Keywords or keyword sequences: «vient de», «formé de», «d'origine»
8.References	Sometimes in Italic	Position at the end of an article or a subparagraph	Metalinguistic keywords. «voyez», «voy.»

A number is given to each automatizable field category. We associate to this number the principal criteria, according to the corpus on which the work is applied (French language dictionaries of XVII, XVIII and XIX centuries), knowing that, it is obviously imperative to

distinguish the fields numbers of those corresponding to the order in which appear the fields within an article. The proposed table is, of course, only an indication.

Illustration of the typologies:

We are going to detail the previous presented typologies with some extracts chosen for representative dictionaries:

Extract from « Dictionnaire françois de Richelet, 1680 »
 « Dictionnaire, s. m. Livre qui contient les mots d'une langue, d'un art, ou d'une science par ordre alphabétique. [Un bon dictionnaire est tres-dificile à faire. Un dictionnaire de droit, un dictionnaire de médecine, un dictionnaire de rimes.] »

The corresponding modeling table is as follows:

1. Entry	<i>Dictionnaire</i>
2. Grammatical category	s.m.
3. Pronunciation	0
4. Domain Category	0
5. Definition	« <i>livre qui ...</i> »
6. Example	« [Un bon dictionnaire est tres-dificile à faire. Un dictionnaire de droit, un dictionnaire de médecine, un dictionnaire de rimes.] »
7. Etymology	0
8. Reference to author and work	0

Extract from « Dictionnaire de l'Académie française, 1694 »

« DICTIONNAIRE, s. m. Vocabulaire. Recueil de tous les mots d'une Langue, mis par ordre. Dictionnaire François. dictionnaire Latin. dictionnaire François-Latin. dictionnaire par ordre alphabétique. dictionnaire par l'ordre des racines, par racines. bon dictionnaire. ample dictionnaire... »

The corresponding modeling table is as follows:

1. Entry	<i>DICTIONNAIRE</i>
2. Grammatical category	s.m.
3. Pronunciation	0
4. Domain Category	0
5. Definition	« <i>Vocabulaire. Recueil de tous les mots d'une Langue</i> »
6. Example	« <i>Dictionnaire François. Dictionnaire Latin. dictionnaire François-Latin...</i> »
7. Etymology	0
8. Reference to author and work	0

Extract from « Dictionnaire étymologique de Ménage, 1694 »

« *AMARER. Terme de Marine, qui signifie attacher, ou lier. Voyez le Sr Guillet dans son Dictionnaire de la Marine : Et Mrs de l'Academie dans leur Dictionnaire de la Langue Française.* »

The corresponding modeling table is as follows:

1. Entry	<i>AMARER</i>
2. Grammatical category	0
3. Pronunciation	0
4. Domain Category	« <i>Terme de Marine</i> »
5. Definition by synonymous	« <i>qui signifie attacher, ou lier</i> »
6. Example	0
7. Etymology	0
8. Reference to author and work	« <i>Voyez le Sr Guillet dans son Dictionnaire de la Marine : Et Mrs de l'Academie dans leur Dictionnaire de la Langue Française.</i> »

Extract from « Dictionnaire de Ménage, 1694 » and « Dictionnaire de Trévoux, 1743 » the same word: « *Ménage* »

« *LARGUE. Terme de Marine. Vent largue. De l'Italien largo. On a dit largue, pour large ; comme cargue, pour charge.* »

The corresponding modeling table is as follows:

1. Entry	<i>LARGUE</i>
2. Grammatical category	[<i>adj</i>]
3. Pronunciation	0
4. Domain Category	« <i>Terme de Marine</i> »
5. Definition	0
6. Example	« <i>Vent largue</i> »
7. Etymology	« <i>De l'Italien largo</i> »
8. Reference to author and work	0

« *LARGUE, s. m. Quelquefois on donne un article féminin à ce mot, & on dit la largue. Terme de Marine. Haute mer. Altum mare, altum. Il n'a guère d'usage qu'en ces phrases ; Prendre le largue, tenir le largue, faire largue ; pour dire, Prendre la haute mer, tenir la haute mer, aller en haute mer. On dit aussi adverbialement qu'ils se sont mis à la largue, qu'ils se sont mis en haute mer, de peur d'être jettés sur les côtes. Tous les autres vaisseaux qui étoient dans le port s'étant mis à la largue, saluerent ces nouveaux venus de toute leur artillerie. Du Loir, Voyage du Lev. page 107. Ce mot est aussi adjectif ; ainsi on appelle vent largue, ou vent de quartier, Obliquus ventus, l'aire de vent qui est comprise entre le vent*

arrière, & le vent de bouline. C'est le plus favorable des vents pour le sillage, car il donne dans toutes les voiles : au lieu que le vent en poupe ne porte que dans les voiles d'arrière, qui dérobent le vent aux voiles des mâts d'avant. Un vaisseau qui fait trois lieues par heure de vent large, n'en fait que deux de vent en poupe. Au lieu de nous tenir au plus près, nous courions vent large de deux aires de vent, afin d'être mieux en ligne. M. le C. de Toulouse. Nous avons un peu de vent, mais il est large. Toutes nos voiles portent, & nous ne roulons plus. L'Ab. de Choisi. Ce mot est la même chose que large, il n'y a que la prononciation de la dernière syllabe qui soit différente ; mais il ne faut s'en servir qu'en termes de Marine. »

The corresponding modeling table is as follows:

1. Entry	LARGUE
2. Grammatical category	s.m. « Quelquesfois on donne un article féminin à ce mot, & on dit la large » « On dit aussi adverbialement » « Ce mot est aussi adjectif »
3. Pronunciation	« il n'y a que la prononciation de la dernière syllabe qui soit différente »
4. Domain Category	Terme de Marine
5. Definition by synonymous	« Haute mer » « pour dire » « on appelle »
6. Example	« Prendre le large, tenir le large, faire large; » « ils se sont mis à la large » « vent large, ou vent de quartier »
7. Etymology	0
8. Reference to author and work	« Du Loir, Voyage du Lev. page 107 » « M. le C. de Toulouse » « L'Ab. de Choisi »

We will thus remember that:

1. It is necessary to label the fields that will enter in the reformulation functioning logic, therefore identifiable in an absolutely objective manner without lending to discussion because of semantic or functional ambiguities subjected to a subjective interpretation.
2. It is necessary to define, for each specific corpus, the criteria set, combinatory or not, that allows the implementation of an automatic labeling based only on formal reference marks, therefore objectively defined.

However, according to the apprehended corpus, the formal criteria will be able to vary. This implies to determine, in addition to the minimal formal criteria, the facts of random distribution and different variant typologies.

Fields isolated according these minimal formal criteria

The fields that can be automatically recognized are:

1. Entries or “vedettes” => position + typography (subparagraph beginning, shifting, big capitals).
2. Grammatical category marks => position + typography + meta-linguistic keywords (just after the entry, sometime in italic, abbreviation system).
3. Indication of pronunciation and of graphic variants => position + typography + meta-linguistic keywords (before or after the mark of grammatical category, pronunciation indication often between parentheses, keyword: the verb « prononcer »).
4. Usage marks, domain marks, technical marks, language level marks ... => position + typography + metalinguistic keywords (« terme de », « en termes de », « en » + domain name, « mot de », ...) provided taking into account the trade names (« les tonneliers disent », « se servent de »...).
5. Synonyms with definitional functions => position + typography (“par exemple dans l’Académie”, capital).
6. Intra-textual references (to other articles within the same work) or inter-textual (to a set of text sources) => position + typography + metalinguistic words (“voyez”).

Difficult fields

In the case of the difficult fields that are very delicate to mark out in an automatic way, we will retain, as an example, the principal complexities specific to the various texts for some categories of fields only. Are concerned the following fields:

7. The elaborated definitions.
8. The translations in foreign languages.
9. The etymologies.
10. The marked bibliographic references (author names cited and concerned text entitled).
11. The cited texts with or without reference.

The complexity report of the identifying information of the fields is as follows:

12. The elaborated definitions: the metalinguistic keyword of reference « signifie » sometimes knows complex alternatives (« on dit de ... que... parce que... ») or is not always expressed, or is taken again during the same article in the case of multiple references words treatment.
13. The translations in foreign languages are often in italic, but as the italic is polysemous in many an-

cient dictionaries, it is necessary to associate another criterion. Indeed, the labeling of the translation languages is very random. Therefore, except considering works announcing a translation of principle (like the “Dictionnaire Universel François & Latin de Trévoux” = DUFLT), only the linguistic ability of the dictionaries reader-consultant authorizes an identification of the translation languages. Knowing that there is often absolute ambiguity of the forms belonging to the whole of the Romance languages, and that in any event even in a dictionary like the DUFLT, it is difficult automatically to select the shapes in italics corresponding to Latin translations.

14. The etymologies are not easily automatically isolable not only because of the variety and complexity of the metalinguistic key-sequences «a pour origine, est formé de, est composé de, est fait de, est emprunté de, vient de, etc...”, but also especially because of polysemous operations of the formulations... to what we can add the randomly presence of an etymological speech, the relative blur statute semantic of certain etymologist speeches, without counting the alternatives of positions and the delimitation difficulties of the fields of a speech relating at the same time to the etymology, the history of the directions, and the socio-cultural evolutions of notions...
15. Marked bibliographical references (names of cited authors and headings of the concerned texts). Without a precise knowledge of the authors, we can trust neither the typography semiotics (proper names authors are always in capital letters or in italic), neither with the position at the end of the quoted text, provided that this last either identifiable like such by italics, nor with the association of a proper name and a more or less identifiable reference because of the abbreviations, etc.
16. Texts cited with or without reference: identification difficulty of the precise limits of a citation, even when the authors and printers of dictionaries choose the italic for the quoted texts, without neglecting the complex facts of citations inserted in quotations, quotations not marked and sometimes easily locatable by a reader accustomed to the reading of certain texts.

5 Automatic Indexing

5.1 Entry Recognition

Table 1 shows in the first column the different kind of entries, in the second column some examples occurred,

and in the last column some rules used for the entry extraction. As the entries are very specific in terms of typography and composition, their extraction doesn't give major errors.

Entry	Example	Characteristics
Word	ABANDON, f. m. Mépris; délaissement. <i>Derelictio, desitium. Neglectus rui aliquis. usage. On ne le trouve guère que dans Mo</i>	Upper case bancicent, spaced, followed by a coma, left justified
Subword with ≠ meaning	ABANDON, fe dit d'ordinaire son à l' <i>abandon</i> , au pillage. j	Small capitals similar typo, but included in bloc def.
Multiple words: several orthog. for same word	ABASSI, ou ABASSIS, subst. m. qui est ronde & qui a cours en Pérse viron dix-huit fols six deniers. Il fa	Similar single word, separated by « ou » and a coma
Particular case	AbBREUVOIR, f. m. Lieu où <i>riam. Mener les chevaux à l'abbr. d'un glacis le plus souvent payé</i>	2 ^d character in lower case (sometimes in small capital) if main entry
Sub-definition	AbBREUVOIR, en terme de Maçonner. que les Maçons laissent entre les joints c entrer du mortier. En ce lens l'on se fet	In a secondary entry, small capitals

Table 1: Entry characterization.

5.2 Information field identification

Concerning the identification and the tagging of the informational fields, an optimal solution could consist in implementing a procedure similar to that performed for the computerization of the TLFi (Treasure of the French Language computerized: www.atilf.fr/tlfi) [18] which is composed into various semi-automatic steps [14]. The first step is dedicated to the typography style labeling. **Figure 9** shows the result of this labeling procedure on the example of the **Figure 8**. The text is cut according to the punctuation. Each text chain is surrounded by a specific XML tag corresponding to the font style, for example <I> ... </I> for italic. The second step is dedicated to the semantic labeling **Figure 10**. By using some dictionaries and by the help of the operator, some semantic tags are assigned to the same text chains previously labeled. Among the semantic labels used, we can find: article (art), entry (ved), word (mot), domain (dom), etc. The last step is related to the structural labeling **Figure 11**. The structure identifies the definition levels (main definition, sub-definitions, etc.).

6 The reedition chain

Once the indexation phase is finished, the following phase concerns the reformulation and the reformatting taking into account the user wishes, the electronic support limits (physical and software characteristic) and the usage rules.

6.1 Reformulation

Here we are interested in the document logical structure, which we wish to modify in order to satisfy a personal use (expressed by a scenario). This modification concerns at the same time the document organization (layout) and its content. For example, the user can ask for a synthetic view of the document, to see only some definitions, by changing the original order, etc. In the case of the dictionaries, the user can select only the synonymous, the grammatical forms, etc.

The scenario is a collection of queries. Considering the document as a structure hierarchy of elements (XML, see **Figure 13**), the queries express two types of operations: tree element selection and tree element reorganization.

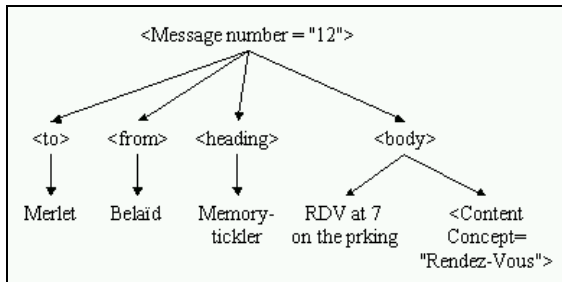


Figure 13: Logical description XML of the indexed document.

There are two kind of selection: direct and indirect. The direct selection permits to indicate an element starting from its absolute position in the logical structure (for example, the sub-definition of the 2nd dictionary entry). The indirect selection (or conditional) gives access to all the elements satisfying certain conditions. Those relate either to the content or its meaning (for example, to reach all the definitions speaking about the concept of “navy”), or function of the elements relative position to be selected in the document logic structure (to reach all the definitions of an entry). We defined an element designation language, which is detailed in **Table 2** and illustrated by examples in **Table 3**.

select_direct(list of target nodes)

<i>select_remove_direct(list of target nodes)</i> direct selection/unselection of nodes and their sub-structures
<i>select_prune()</i> Pruning of empty branches
<i>select_keep_parent_from_tag(node, list of tags to search)</i> <i>select_remove_parent_from_tag(node, list of tags to search)</i> selection/unselection of a node and its substructure according to the presence of specific descendants;
<i>select_in_parent_keep_tag(node, list of tags to search)</i> and <i>select_in_parent_remove_tag(parent, node, list of tags to search)</i> selection/unselection of the part of the node substructure answering to some conditions

Table 2: Selection operators.

<i>select_direct([Entry/Entry definition entry/synonymous])</i> to keep only the definition and the synonymous of a definition
<i>select_remove_direct([entry/definition [2], entry/definition [4])</i> to keep only the definitions 2 and 4 of an entry
<i>select_keep_parent_from_tag(entry/definition, [synonymous])</i> to keep only the entries of a n entry containing a synonymous
<i>select_remove_parent_from_tag(entry/definition, [content [:@concept = "navy"]])</i> to not keep the entry definitions speaking about the concept "navy"
<i>select_in_parent_keep_tag(definition, [synonymous, illustration])</i> to keep, in each definition, only synonymous and illustrations
<i>select_in_parent_remove_tag(definition, [illustration])</i> to erase, in each definition, all the illustrations

Table 3: Example of selection operators.

The purpose of the reorganization operators is to authorize the original document reading order modification. For example, the user can wish to read synonymous before definition, or to replace all the illustrations at the end of a definition. We can see some example of such operators in **Table 4**.

<i>arrange_change_order(common parent, tag element1, tag element2)</i> to change the position of two elements (regarding the corresponding tree root, or regarding a specific node)
<i>arrange_move_after(tag element to move, tag element used as a reference, list of eventual tags to add)</i> to move an element and to situate it after an other
<i>arrange_up_one_level(element to climb up, list of tags to add)</i> <i>arrange_change_level(element to climb up, new parent, list of tags to add)</i> to « climb up » an element from one or more levels in tree structure

arrange_reference(tag denoting a link, attribute of this tag allowing to treat the link, path where one can find the link tag) in order to put together the link tags and references

Table 4: Reorganization operators.

<i>arrange_change_order</i> (//entry, definition[1], definition[2]) to inverse the definitions 1 and 2
<i>arrange_move_after</i> (entry//illustration, entry//definition [last ()], [all_illustration]) to place all the entry illustrations in a structure tagged by <all_illustration> and situated after the last definition
<i>arrange_up_one_level</i> (entry/definition[1]/sub-definition[3], [definition]) to move the third sub-definition of the first definition in a new definition (which will be the second definition, the following definitions are "delayed")
<i>arrange_change_level</i> (entry/definition[1]/sub-definition[3], entry/definition[2], []) to move the third sub-definition of the first definition of the entry in the second definition.
<i>arrange_reference</i> (link, ref_id, entry/citation/ reference) to replace all the type tags <link ref_id = "****"> by the reference elements placed in citation part of the entry and having the same value for ref_id

Table 5: Reorganization examples.

To carry out these reformulation actions, we choose to use XSL¹[7]. It is a style sheets expression language, defining a transformation language for XML documents (XSLT²) [16,17]. The user wishes are retransmitted in a series of transformations (XSLT) and enunciated using the selection and reorganization operators, and XPath, a language of object designation adapted to the tree structure of XML. This language allows expressing conditions on the node, its position in the tree and on the presence (or the value) of an attribute. For example, "*definition [not (contains (@domain, 'navy'))]*" selects the nodes "*definition*" of the node root having an attribute "*domain*" which doesn't comprise the chain "navy".

The call of these operators generates an XSL that is applied to the document. **Table 5** provides an example of the XSL file generated by a selection request carried out by the user.

6.2 Reformatting

There are three constraints for the document displaying on an electronic support: user wishes, support capacities and usage rules. The choice of a support produces

¹ eXtensible StyleSheet Language

² XSL Transformation

physical constraints (size of the screen, navigation tools available, etc.) and displayable formats (HTML, pdf, Bitmap, etc.). The usage rules are provided by experts: ergonomists and pedagogues concerning the facility of training and text reading. It is necessary to test the coherence between the constraints emanating from these various sources and, if necessary, to solve the conflicts.

Accordingly, we introduced a priority associated with each constraint. We classified these constraints in three categories:

- 1) The intra-blocks constraints concerning the specific posting of a definition block (font style, size, alignment, line space, hyphenation management, widows management and orphans, dimensions, typographic grey, border...)
- 2) The inter-blocks constraints relating to the sequence of the blocks during their posting (with left/right of, page break before/after, to keep on the same page that next/previous)
- 3) The constraints of page setting on the final support (margins, share-cropping, navigation, distinction peer/odd pages...). Initially, only the constraints with intra and inter blocks were studied.

```

1 <?xml version="1.0" encoding="ISO-8859-1"?>
2 <xsl:stylesheet
3 xmlns:xsl = http://www.w3.org/1999/XSL/Transform
4 version = "1.0">
5   <xsl:output method = "xml" encoding="ISO-8859-1"/>
6   <xsl:output indent = "yes"/>
7
8   <xsl:template match = "*" priority="-1">
9     <xsl:copy>
10       <xsl:apply-templates select = "*" />
11     </xsl:copy>
12   </xsl:template>
13
14   <xsl:template match = "entry/definition">
15     <xsl:copy-of select = "." />
16   </xsl:template>
17
18   <xsl:template match = "entry/synonymous">
19     <xsl:copy-of select = "." />
20   </xsl:template>
21
22 </xsl:stylesheet>

```

Table 6: XSL generated by the operator "select_direct([entry/definition, entry/synonymous])".

The constraints intra-blocks are expressed as a quintuplet [*Constraint, Target, Value, Priority, Identifier*]. *Constraint* corresponds to the attribute (font style, size...), *Target* indicates the concerned elements ("en-

try/synonymous", "entry/definition"...). *Value* indicates the value associated with the constraint (for "font style", the value is the name the font style: Times, Arial...). *Priority* is an integer and *Identifier* is a single number. Here are some examples of intra-blocks constraints expressed according to this syntax: [font style, book/title, times, 2, 1], [withdrawal, book/paragraph, 10, 5, 2] and [size, entry/definition, \geq , 20], 4, 3].

The inter-block constraints are specified in the form of another quintuplet [*Font style*, *Target_1*, *Target_2*, *Priority*, *Identifier*]. *Constraint*, *Identifier* and *Priority* have the same meaning as for the constraints intra-blocks. *Target_1* indicates the elements that correspond to the first argument intervening in the inter-blocks constraint and *Target_2* indicates the elements of the second argument so necessary [keep_together, drawing, legend, 1,4] and [page_break, entry/definition, $_$, 5, 5] are two examples of inter-block constraints.

In order to combine these various display constraints and detect the possible conflicts effectively, we use an expert system. Only constraints belonging to the same category (intra-blocks, inter-blocks and page setting) can have a conflict. For example, if two constraints relate the same attribute and if the targets are identically referred (for example: [*font style*, *author*, *Times*, 4, 6] and [*font style*, *author*, *Arial*, 5, 7]) or if they refer to different attributes being able to enter in conflict (for example: [*page_break*, *legend*, $_$, 2,8], and [*keep_together*, *figure*, *legend*, 4, 9]). We allowed the user (or more exactly an expert of the edition field) to express the couples of different constraints in conflict. We start by locating "the wrong" conflicts: it is the case of two rules of which one is more restrictive than the other. For example, [*size*, *entry/definition* \geq , 20], 4, 3] and [*size*, *entry/definition*, \geq , 16], 5, 11] will be replaced by [*size*, *entry/definition*, \geq , 20], 5, 3]. One preserves the most restrictive rule (in the example: \geq 20) but with the highest priority (here: 5 and either 4). "True" conflicts will be solved according to the priorities attached to each constraint.

To express the formatting constraints, we use the formatting objects (FO). They are XML tags defining the visual aspect of a document. They constitute a description language and allow to describe all the problems relating to the typography, as well as the hyphenation constraints, the widows and the orphans, all the inter-block constraints, etc. For each document element we associate its page-setting. For example, we can define a formatting object for the element TITLE in order to display his content in underlined or italic. The use of the FO allows being independent from the layout format.

The expert system was performed in Prolog. First of all, it ensures the detection and resolution stages of the conflicts then it generates the FO. The fact basis counts all the constraints and the rule base translated the practices and the expert knowledge. The obtained document XML + FO can either be visualized directly in a navigator XSL FO, or transformed into a document of another nature by a processor XSL FO. The project Apache [19] delivered the processor FOP³ to realize the rendered of a XML document with "Formatting Objects" in various formats: PDF, PCL⁴, PS⁵, SVG⁶, XML, direct impression, AWT (presents the result as a Java graphic in a window), MIF⁷ and TXT.

The electronic support occurs in the composition stage giving the output format. The high number of existing output formats allow the same document to be portable on multiple supports, and justifies the use of the XSL FO language.

7 Conclusion

This work explored a generic way of ancient documents automatic valorization. The valorization is based on an adaptive indexing that highlights the main physical and logical structure components. Then, a reedition chain is proposed based on structure reformulation and reformatting. All these phases obey to three kinds of constraints coming from the user, the usage and the electronic support for which the system has to resolve the conflicts. The reformatting phase leads to a structure section and reorganization while the reformatting leads to the creation of formatting objects. Considering the document in an XML format, all the chain is seen as an XSL transform. All the constraints are transformed in XSL orders.

The application of this chain on ancient dictionaries allowed us to express the main indexes useful for dictionary reformulation and the main queries that we can ask for dictionary consulting. The work done in this first step was limited to some simple pages of the dictionary. We plan to enlarge the procedure to all the dictionary and to automate the index searching procedure. This step constitutes the main challenge of the work because this needs to extract ontology and to be able to interpret some defini-

³ FOP: Formatting Objects Processor

⁴ Format for Hewlett-Packard Printers

⁵ PostScript

⁶ Scalable Vector Graphics

⁷ Maker Interchange Format, used for Adobe Framemaker

tional fields to extract the synonyms and the references.

8 References

[1] I. Turca, "Les requêtes prédictives à l'impact croissant de l'inférence pour les bases de données électroniques de l'Université de la Sorbonne", in Actes du colloque TIC-SHS, L'ors, mars 2002, p. 25-31.

[2] I. Turca "L'édition scientifique d'ouvrages anciens sur support électronique : l'éditique méthodologique du traitement numérique des originaux et leurs typographies des dictionnaires dans le programme de numérisation des collections d'ouvrages anciens de la bibliothèque TILF", Actes de la XIVe Conférence Européenne TEX (EuroTeX20/3).

[3] I. Turca, Travaux, oct. 1999 : "Les différentes éditions du Dictionnaire de Trivoux : l'impact des éditions et prises en compte de l'éditique pour une nouvelle approche de l'ordre chronologique des éditions, problèmes et perspectives", in Actes du colloque international sur la "Rayonnement du Dictionnaire Universalis : le Trésor" (Télex, octobre 1999).

[4] I. Turca "Basage XML ou autres logiciels pour la diffusion des archives formatées : bibliographie et applications méthodologiques. Le cas de l'éditique de l'Encyclopédie Française (1694) et des bases électroniques de l'Encyclopédie Française (1694) et de l'Encyclopédie Française (1694)", in Actes du colloque international sur les problèmes de l'éditique de l'Encyclopédie Française, Limoges, 11-20 novembre 1998.

[5] I. Turca, "Principes de la numérisation des collections d'ouvrages anciens de la bibliothèque TILF de Nancy, UMR 1118 : procédures de numérisation, suivi de l'indexation et de la structuration de l'information de diffusion", ICHIM/03, sur la numérisation du patrimoine et de l'Écolle de Louvain.

[6] A. Belaid and L. Pierrat. A general approach to OCR performance evaluation. Electronic Imaging / San Jose, California, 2002.

[7] L. Pierrat and A. Belaid, "XML/SGML platform for document analysis". In DLI / Seattle, WA, USA, 2001.

[8] J. V. Rice, F. R. Jenkins, and T. S. Parker. "Textual analysis of OCR accuracy". Technical report, Information Science Research Institute, Univ. of Nevada at Las Vegas, Las Vegas, Nevada, 1997.

[9] E.W. Myers, "On O(N) difference algorithms and its variation", Algorithmica, 1, (1986), pp. 251-261.

[10] D. Dorian, "Textual analysis and retrieval of document images: A survey", Computer Vision and Image Understanding, 70(3), June 1998. Special Issue on Document Image Understanding and Retrieval," L. Kuhl and H. J. Baird (Eds.).

[11] R. Aram, B. K. Bhargava, and Y. Yesha, eds, "Digital Libraries: Current Issues". Springer-Verlag, December 1997. Lecture Notes in Computer Science 1111.

[12] J. Baird, "Recommendation of Tables of Contents for Electronic Library Consulting". International Journal on Document Analysis and Recognition, 10(1), vol. 1, n° 1, 11 p.

[13] D. Beaulieu, J. Belaid and N. Benet, "L'éditique pour la bibliographie : références typographiques et fichiers". I/D/R'03, Éditions, 3-1 / Juin 2003.

[14] P. Bernard, J. Lecomte, J. Deleu and J.-M. Pierrat, "Computerized linguistic resources of the research laboratory ATILF for lexical and textual analysis: Facts, TLF, and the software Shell". LREC 2002, Proceedings of the Third International Conference on Language Resources and Evaluation, Las Palmas, Spain, 7 May - 2 June 2002, Vol. 1, p. 1091-1097.

[15] A. Todirascu, F. De Beuvoir and F. Ruscel, "A new Description Logic for Intelligent Documents". Proceedings of the 14th IIR Annual Meeting, 18-19 November 1999, Innsbruck, France, pp. 179-186.

[16] E. R. Harter and J. S. Mans, "XML is a Natural". January 2001, R/V Edition.

[17] O. Hitz, Lys Ruffad and R. In, "Using XML for document recognition". Proceedings of DLIA/9, Bengali, India.

[18] J. Deleu and J.-M. Pierrat, "La Téléologie de la Langue Française informatique : un exemple d'informalisation d'un thème de l'enseignement de l'écriture". TAL, vol 14 n° 2/03, Hermès Éditions, pp. 28.

[19] <http://xml.org>