

# Lexical sets and Text-Processing

**Christian CHAMPENDAL**

Faculté des Lettres et Sciences Humaines  
- University of Nice (UNSA)  
Nice, France  
champend@unice.fr

**Thierry PITARQUE**

Equipe ASTRE, UNSA-CNRS  
06903 Sophia Antipolis-  
Cedex France  
pitarque@unice.fr

## Abstract

The extraction of lexical sets from a corpus in Digital Signal Processing (DSP) has been detailed before on general sets, with direct ELT applications. In this contribution, a more specialized set is investigated to illustrate the possibility of actually using the results in more "intelligent" Text-Processing.

## 1 Introduction

The **Lexical Set Extraction** technique was presented before<sup>1</sup>. From a sample corpus in Digital Signal Processing, lexical sets are extracted, the most important ones are investigated to establish recurring predicative relations within the domain. In this contribution, another set from the same corpus is studied, to briefly illustrate the technique. Possible applications in Text Processing are then assessed.

### 1.1 Corpus

Results are extracted from a sample corpus of 135 articles in DSP, which represents about 500 pages, 3900 paragraphs, 16 500 lines and 350 000 word occurrences (2 million characters with a size of 3,27 Moctets). Articles were selected from a CD-Rom, ICASSP'98 International Conference on Acoustics, Speech and Signal Processing according to quality and diversity. Lexical sets are then easily obtained from the lists of lemmas and their occurrences. From all the corpus sets, about 100 are quantitatively prominent and their study will ensure a reasonable coverage of the domain.

## 2 Lexical set [FILTER]

This more specialized set is frequently found in the corpus. It is constituted of only four **lexical units** : *filter/filters* are verb-noun homographs,

---

<sup>1</sup>fully in [Champendal03] and in a previous contribution [Champendal04 Oct], both drawing on general lexical set [USE]. ELT and especially ESP classroom applications are direct: lexical acquisition, syntactic and semantic structuration.

*filtering* is a verb-noun-gerund one, *filtered* is a verb-adjective homograph. Most of the 1441 **occurrences** are nominal with few verbal combinations. Hence, the different units of the set can be predicted in the nominal **positions**: *Co*, *C<sub>1</sub>*, *C<sub>2</sub>* (complements), *Circ<sub>n</sub>* (circumstantials) or *NC* (noun complement).

[Filter] (1441/4)  
(809) filter (62) filtered  
(146) filtering (424) filters

As for **collocations**, elements of the set are found as modifiers or heads of nominal, adjectival, prepositional phrases and compounds, some of which are listed below:

adaptive filter(60)  
adaptive filtering(14)  
filter bank(151)  
fir filter(55)  
modulated filter bank(23)  
weighted median filter (13)  
polyphase synthesis filter(8)  
polyphase synthesis filter bank(6)

Some interesting side phenomena are observed for this set: **acronyms** (EKF= *extended Kalman filter*, QMFs=*quadrature mirror filters*, WM=*Wiener median filters*); **affixates** with or without dash (*subfilter*, *anti-filter*, *post-filter*); **agglutinates** (*filterbanks*), and **coinage** ("*lifered*" meaning *low-pass filtered*). Two parallel collocative **paradigms** built on *filter* 'Object' and *filtering* 'Activity' coexist (*digital filtering/filter(s)*; *kalman filtering/filter*). The quantitative importance of these phenomena underlines the necessity to process them thoroughly

## 3 Uses in Text-Processing

**Concordances** of verbal **and** nominal units and their correlations are studied with a view to extracting an **Extended Predicative Formula** that lists complements and circumstantials according to frequency for all predicative members of the set<sup>2</sup>

---

<sup>2</sup>Order of complements and circumstantials is obviously flexible, numbers correspond to statistically

Co <filter>C<sub>1</sub> C<sub>irc1</sub>C<sub>irc2</sub>C<sub>irc3</sub> C<sub>irc4</sub>

Co = {we, algorithm, filter, weight,...}

C<sub>1</sub> = {data, image, signal, noise, procedure, sample, waveforms,...}

C<sub>irc1</sub> = **from** {error, source, signal}

C<sub>irc2</sub> = **in/over** {domain, region, frequency band...}

C<sub>irc3</sub> = **to do** {imaging, noise removal...}

C<sub>irc4</sub> = **by use of /using** {statistics, optimization techniques...}

Verbal units represent only 2% of all occurrences. They are mostly passive; a small number of infinitives and gerunds are also found in the corpus.

*In WI, pitch-cycle waveforms **are filtered in** the evolution domain **to decompose** the signal **into** two waveform surfaces, one characterising voiced speech and a second representing unvoiced speech.*

*In an array processing application in [10] an additional noise determining transducer was used to cancel noise and interference, but in other circumstances it becomes appropriate **to filter noise from** signals and it has been shown in this paper that it is possible to consider **filtering in** the two dimensional third order cumulant domain*

The sets' **connectivity** (ie. affinity with other prominent sets of the corpus), is automatically detected as combined occurrences of two sets are counted in the *Associated Units* files facility of the **Z-text**® software. **Clusters** are then studied as a priority within the corpus: they will enable prediction of recurrent relations in the domain. Highest clustering obtains with the following sets: [OPTIM], [DESIGN], [FILTER] (reflexive), [RESULT], [SIGN], [SHOW], [REQUIR], [USE], [WEIGH]. Normally, varied nominal, verbal and participial solutions are present.

[FILTER+OPTIM]

filter~optimal:14 filters~optimal:5

filter~optimization:7 filters~optimization:3

filter~optimize:2 filtering~optimal:4

filter~optimum:2

established positions. Co is 'human', 'object' or 'abstract', C<sub>1</sub> 'object' or 'abstract'; Circ<sub>n</sub> can be expressed or implied and recoverable in context

The studied set [FILTER] is mostly nominal. Nominal clusters (ie. connection of mostly nominal sets) produce mainly Noun-Noun or Adjective-Noun associations, but this set also clusters with verbal sets ([SHOW], [USE], [REQUIRE]...), generating predicative relations where it occupies subject or direct object positions. Various connections with [USE] are presented below:

[FILTER+USE]

filter~use:8 filter~useful:3

filter~uses:2 filter~using:27

filters~use:5 filters~using:14

**Extended Predicative Formula: <Use>**

C<sub>0</sub>=Subject <Use>C<sub>1</sub>=Complement(Od)

Circ<sub>n</sub>=Circumstantial(Circ/Adv)

*The algorithm **uses** two neighbouring pixels, one left and the other top, **as** a pioneering block **to search** for the best matched blocks **inside** a pre-defined window.*

a) Verbal use

C<sub>0</sub>=**filter** <use>C<sub>1</sub> (Circ<sub>n</sub>)

C<sub>0</sub> <use>C<sub>1</sub>=**filter**  
(Circ<sub>n</sub>)

*The echo-path model **uses** a single pole single zero digital **filter***

C<sub>01</sub>=**filter** <be-used>(C<sub>10</sub>) (Circ<sub>n</sub>)

*In prediction, a **filter** is **used** to estimate future values of a signal from prior observations.*

b) Nominal use

<use> of C<sub>1</sub>=**filter** (Circ<sub>n</sub>)

*This standard allows the **use** of any wavelet **filter** up to a length of 32 taps.*

c) Adjectival-participle use

C<sub>01</sub>= **filter** <used> (Circ<sub>n</sub>)

<modif+used>C<sub>01</sub>=**filter** (Circ<sub>n</sub>)

*Block diagram of the constrained adaptive **filter** used to determine the Wiener solution.*

*Our iterative algorithm for state estimation is based on the Expectation- Maximization (EM) algorithm and outperforms the widely used Extended Kalman **filter** (EKF)*

d) Circumstantial link <Use>

C<sub>0</sub> <use>C<sub>1</sub> Circ<sub>3</sub>=to filter

*These weights, in turn, were used to filter the multi-tone signal resulting in the estimate shown in Fig1* 5

Circumstantial link <Filter>

C<sub>0</sub> <filter>C<sub>1</sub>  
Circ<sub>4</sub>=(by) using/by use of Z

*Digital image enhancement and noise filtering by use of local statistics.*

**-Ing** participle forms generate typical ambiguities where a similar surface portion of the sentence, can refer to different underlying relations<sup>3</sup>

C<sub>0</sub>=filter <using>C<sub>1</sub> (Circ<sub>n</sub>)

*A filter using an 11 8 grid filter for F1 and a 5x5 average for F2 was trained on images of faces degraded by AWGN with  $\sigma^2 = 200$ .*

#### 4 Conclusions

Combined with reduction of predicative relations to their core relations, the technique helps optimize various NLP applications: data and text-mining, computer-aided generation of abstracts, multilingual corpus alignment (an equivalent sample corpus was easily obtained in French, using the same technique and the sets connected together with their results: **English:** [FILTER]<>**French:**[FILTRE]). In the next stage, top-down expert knowledge will help to shape legal ontologies of the domain.

#### References

CHAMPENDAL C. , 2004 October; Language teaching and corpus linguistics, IATEFL – CALL Review, ISSN 1028-428

CHAMPENDAL C., 2004 November ; Outils linguistiques pour Ingénieurs, Actes du Colloque TICE Méditerranée – ISDM ([www.isdm.org](http://www.isdm.org))

CHAMPENDAL C. , 2003 ; Analyse d'un corpus d'articles en Traitement Numérique du Signal en vue de modélisation linguistique et

d'application TIC ; Thèse de Sciences du langage, directeur: Henri Zingle.

CHAMPENDAL C. 1999 Morphosyntaxe du verbe anglais; contribution au JILA'99 ; Nice 24-25 juin 1999.

MEL'UK Igor A., CLAS A, POLGUÈRE A, 1995 *Introduction à la lexicologie explicative et combinatoire*; Ed Duculot.

TRIBBLE C. 2000, Practical uses for language corpora in ELT in A Special Interest in Computers, P. Brett and G. Motteram Eds, IATEFL.

*TdL99a Travaux du LILLA* , 1999 ; Henri Zingle Numero special– La modélisation des langues naturelles – Aspects théoriques et pratiques.

ZINGLE H, 1998 ZTEXT : un outil pour l'analyse de corpus in *Travaux du LILLA n°3*, pp 69-78.

---

<sup>3</sup>Occurrences being limited, some sentences are manipulated to obtain the required string **filter+using**. Even marginal topicalized possibilities are accounted for: C<sub>0</sub> <filter>C<sub>1</sub> Circ<sub>4</sub>=<using>C<sub>1</sub>.  
*This entire data file we filter (by) using a Hamming weighted (mis)matched filter*