



**HAL**  
open science

# Conceptual Vectors, a complementary tool to Lexical Networks

Didier Schwab, Lim Lian Tze, Mathieu Lafourcade

► **To cite this version:**

Didier Schwab, Lim Lian Tze, Mathieu Lafourcade. Conceptual Vectors, a complementary tool to Lexical Networks. NLPCS'07: 4th International Workshop on Natural Language Processing and Cognitive Science, Jun 2007, pp.10. lirmm-00200892

**HAL Id: lirmm-00200892**

**<https://hal-lirmm.ccsd.cnrs.fr/lirmm-00200892>**

Submitted on 21 Dec 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.

# Conceptual Vectors, a complementary tool to Lexical Networks

Didier Schwab<sup>1</sup>, Lim Lian Tze<sup>1</sup>, and Mathieu Lafourcade<sup>2</sup>

(1) Computer-Aided Translation Unit (UTMK) School of Computer Sciences,  
Universiti Sains Malaysia, Penang, Malaysia

{didier, liantze}@cs.usm.my

<http://www.lirmm.fr/~schwab>

(2) Université Montpellier II-LIRMM, Montpellier, France,

lafourcade@lirmm.fr,

<http://www.lirmm.fr/~lafourcade>

**Abstract.** There is currently much research in natural language processing focusing on lexical networks. Most of them, in particular the most famous, WordNet, lack syntagmatic information and especially thematic information (*"Tennis Problem"*). This article describes conceptual vectors that allows the representation of ideas in any textual segment and offers a continuous vision of related thematic, based on the distances between these thematic. We show the characteristics of conceptual vectors and explain how they complement lexico-semantic networks. We illustrate this purpose by adding conceptual vectors to WordNet by emergence.

Originally resulting from Ross Quillian's work on psycholinguistics [1], lexical networks are today object of many researches in Natural Language Processing. They are employed in many tasks (lexical disambiguation [2]) or field applications (machine translation with multilingual networks like Papillon [3] or [4], information retrieval or text classification [5]). Most of these networks and specifically the most famous, WordNet [6], miss syntagmatic information and, in particular, information concerning the domain usage of terms or at least thematically related terms. There is thus no direct relation between terms like *teacher-student* or *boat-port*. This phenomenon is called the *"tennis problem"* [[6], p. 10] because it has been noticed that it was necessary to seek *ball*, *racket* and *court* at various places of the hierarchy.

For several years, TAL team (Natural Language Processing team) from LIRMM (Montpellier Laboratory of Computer Science, Robotics, and Microelectronics) works on a formalization of the projection of the linguistic concept of semantic field in a vector space, the conceptual vectors. They allow to represent ideas contained in an unspecified textual segment and allow to obtain a continuous vision of thematic used thanks to the calculable distances between them.

In this article, we present the conceptual vectors and especially the version built by emergence. We show their characteristics and why they are complementary to lexico-semantic networks. We illustrate our purpose by an experiment done at UTMK (Computer-Aided Translation Unit), Universiti Sains Malaysia, Penang which consisted in enriching the WordNet data by conceptual vectors built by emergence.

# 1 Lexico-semantic Networks : Example of WordNet

## 1.1 Principle

WordNet is a lexical database for English developed under the direction of George Armitage Miller by the Cognitive Science Laboratory of the university of Princeton (New Jersey, USA). It aims to be consistent with the access to the human mental lexicon.

WordNet is organized in sets of synonyms called synsets. To each synset corresponds a concept. Terms meaning is described in WordNet by three means :

- their *definition*
- the *synset* to which the meaning is attached.
- the *lexical relations* which link synsets. There are, among others, hyperonymy, meronymy and antonymy.

WordNet 2.0 contains 152059 terms what constitutes a relatively broad cover of the English language. In first versions of Wordnet, the lexical relations connect only items in the same part of speech. There is thus one hierarchy for nouns, one for adjectives, one for verbs and finally one for the adverbs.

## 1.2 Weakness of wordnet

In [7], authors of WordNet (we were at version 1.6) record six weaknesses in the their network constitution:

1. the lack of connections between noun and verb hierarchies;
2. limited number of connections between topically related words;
3. the lack of morphological relations;
4. the absence of thematic relations/selectional restrictions;
5. some concepts (word senses) and relations are missing;
6. since glosses were written manually, sometimes there is a lack of uniformity and consistency in the definitions.

If items 3, 5 and 6 don't interest us in this article, we will show the conceptual vectors contribution to the resolution of the others, all three constitute the tennis problem.

## 1.3 Previous work to solve the problem

In this article, we will be interested only in Wordnet version 2.1 which was the last available when we carried out our experiments. A new version (3.0) was released in December 2006 but it does not seem to have some improvements compared to the previous version for what interests us here.

Since version 2, relations as *derivationally related form* makes it possible to link adjectives to verbs or adjectives to names. In the same way, an usage domain can be addressed to synsets. However, the number of these data still seem too restricted to be sufficiently relevant. Typical relations as *'teacher'-'student'* *'boat'-'port'* or *'doctor'-'hospital'*, often essential to a task of lexical disambiguation, are not still present and

the restricted number of thematic indications like domain does not make it possible to compensate this defect. Several solutions were proposed to solve whole or part of this problem.

With Extended WordNet, [7] proposes to disambiguate definitions of WordNet as a semi-automatic way. The idea is for each definition to annot each word with the number of the meaning used. One can then compare two synsets and evaluate their similarity. We will see that we use this information to manufacture the conceptual vectors of this experiment.

Others also add information to the synsets. Thus, [8] add lexical signatures resulting from tagged corpora or Web.

On the other hand, others rather seek to increase the number of arc existing. [9], for example, combines different metrics to create links between synsets from their definitions and from a thesaurus. [10] use a cocurrences network to extract typical relations like those presented in the previous section.

We can see that all these proposals have in common to belong in particular to the discrete field. Our is to introduce a continuous representation of the ideas contained into the network, conceptual vectors.

## 2 CONCEPTUAL VECTORS

### 2.1 Principle and Thematic Distance

We represent thematic aspects of textual segments (documents, paragraph, phrases, etc) by conceptual vectors. Vectors have long been used in information retrieval [11] and for meaning representation in the LSI model [12] from latent semantic analysis (LSA) studies in psycholinguistics. In computational linguistics, [13] proposed a formalism for the projection of the linguistic notion of semantic field in a vectorial space, from which our model is inspired. From a set of elementary concepts, it is possible to build vectors (conceptual vectors) and to associate them to any linguistic object. This vector approach is based on known mathematical properties. It is thus possible to apply well founded formal manipulations associated to reasonable linguistic interpretations. Concepts are defined from a thesaurus (in our prototype applied to French, we used Larousse thesaurus [14] where 873 concepts are identified) to compare with the thousand defined in Roget thesaurus [15]). Let  $C$  be a finite set of  $n$  concepts, a conceptual vector  $V$  is a linear combinaison of elements  $c_i$  of  $C$ . For a meaning  $A$ , a vector  $V(A)$  is the description (in extension) of activations of all concepts of  $C$ . For example, the different meanings of *door* could be projected on the following concepts (the  $CONCEPT[*intensity*]$  are ordered by decreasing values):  $V(\langle door \rangle) = (OPENING[0.8], BARRIER[0.7], LIMIT[0.65], PROXIMITY[0.6], EXTERIOR[0.4], INTERIOR[0.39], \dots$

### 2.2 Operations on vectors

**Angular Distance** Comparison between conceptual vectors is done using angular distance. For two conceptual vectors  $A$  and  $B$ ,

$$\begin{aligned} \text{Sim}(X, Y) &= \cos(\widehat{X, Y}) = \frac{X \cdot Y}{\|X\| \times \|Y\|} \\ D_A(A, B) &= \arccos(\text{Sim}(A, B)) \end{aligned} \quad (1)$$

Intuitively, this function constitutes an evaluation of the *thematic proximity* and measures the angle between the two vectors. We would generally consider that, for a distance  $D_A(A, B)$ :

- if  $\leq \frac{\pi}{4}$  ( $45^\circ$ ), A and B are thematically close and share many concepts;
- if  $D_A(A, B) \geq \frac{\pi}{4}$ , the thematic proximity between A and B would be considered as loose;
- around  $\frac{\pi}{2}$ , they have no relation.

$D_A$  is a real distance function. It verifies the properties of reflexivity, symmetry and triangular inequality. We have, for example, the following angles (values are in radian and degrees).

$$\begin{aligned} D_A(\mathcal{V}(\text{tit}), \mathcal{V}(\text{tit})) &= 0 \quad (0^\circ) \\ D_A(\mathcal{V}(\text{tit}), \mathcal{V}(\text{bird})) &= 0.55 \quad (31^\circ) \\ D_A(\mathcal{V}(\text{tit}), \mathcal{V}(\text{sparrow})) &= 0.35 \quad (20^\circ) \\ D_A(\mathcal{V}(\text{tit}), \mathcal{V}(\text{train})) &= 1.28 \quad (73^\circ) \\ D_A(\mathcal{V}(\text{tit}), \mathcal{V}(\text{insect})) &= 0.57 \quad (32^\circ) \end{aligned}$$

The first one has a straightforward interpretation, as a ‘*tit*’ cannot be closer to anything else than itself. The second and the third are not very surprising since a ‘*tit*’ is a kind of ‘*sparrow*’ which is a kind of ‘*bird*’. A ‘*tit*’ has not much in common with a ‘*train*’, which explains the large angle between them. One may wonder why ‘*tit*’ and ‘*insect*’, are rather close with only  $32^\circ$  between them. If we scrutinise the definition of ‘*tit*’ from which its vector is computed (*Insectivorous passerine bird with colorful feather.*) perhaps the interpretation of these values would seem clearer. In effect, the thematic is by no way an ontological distance.

### 2.3 Neighbourhood : a continuous vision of thematic aspects

**Principle** The thematic neighbourhood function  $\mathcal{V}$  is the function which returns the  $n$  closest LEXICAL OBJECTS<sup>1</sup> to a lexical object  $x$  according to the angular distance:

$$\begin{aligned} \sigma \times \mathbb{N} &\rightarrow \sigma^k : \\ X, k &\rightarrow E = \mathcal{V}(D_A, X, k) \end{aligned} \quad (2)$$

where  $\sigma$  the set of LEXICAL OBJECTS. The function  $\mathcal{V}$  is defined by :

$$\begin{aligned} |\mathcal{V}(D_A, Z, k)| &= k \\ \forall X \in \mathcal{V}(D_A, Z, k), \quad \forall Y \notin \mathcal{V}(D_A, Z, k), \\ D_A(X, Z) &\leq D_A(Y, Z) \end{aligned} \quad (3)$$

Thematic neighborhood function can be used for learning to check the overall relevance of the semantic base or to find the more appropriate word to use for a statement.

<sup>1</sup> We call LEXICAL OBJECT any object in the lexicon which meaning can be described. For WordNet, they are entries (called in this article lexical items) and synset.

Thus, they give new tools to access words through a proximity notion to those described in [16] and issued from psycholinguistic considerations like form, part of speech, navigation in a huge associative network. They allow to navigate in a continuous way and not in a discrete way as commonly done in semantic networks.

**Examples** For example, we can have :

$\mathcal{V}(D_A, \text{'life'}, 7) = (\text{'life'} \ 0.4) (\text{'to born'} \ 0.449) (\text{'alive'} \ 0.467) (\text{'to live'} \ 0.471) (\text{'existence'} \ 0.471) (\text{'mind'} \ 0.484) (\text{'to live'} \ 0.486)$

$\mathcal{V}(D_A, \text{'death'}, 7) = (\text{'death'} \ 0) (\text{'murdered'} \ 0.367) (\text{'killer'} \ 0.377) (\text{'age of life'} \ 0.481) (\text{'tyrannicide'} \ 0.516) (\text{'to kill'} \ 0.579) (\text{'dead'} \ 0.582)$

**Vectorial Sum** If  $X$  and  $Y$  are two vectors, their *normalised vectorial sum*  $V$  is defined as :

$$\vartheta^2 \rightarrow \vartheta : V = X \oplus Y \quad | \quad V_i = \frac{X_i + Y_i}{\|X + Y\|} \quad (4)$$

where  $\vartheta$  is the set of the conceptual vectors,  $V_i$  (resp  $X_i, Y_i$ ) is the  $i$ -th component of the vector  $V$  (resp.  $X, Y$ ).

The normalized vectorial sum of two vectors gives a vector equidistant according to the angle of the first two vectors. It is in fact an average the summoned vectors. As an operation on the conceptual vectors, one can thus see the normalized vectorial sum as the union of the ideas contained in the terms.

**Normalised Term to Term Product** If  $X$  and  $Y$  are two vectors, their *normalised term to term product*  $V$  is defined as :

$$\vartheta^2 \rightarrow \vartheta : V = X \otimes Y \quad | \quad v_i = \sqrt{x_i y_i} \quad (5)$$

The  $\otimes$  operator can be interpreted as an operator of intersection between vectors. If the intersection between two vectors is the null vector, then they do not have anything in common. From the point of view of the conceptual vectors, this operation thus makes it possible to select the ideas common to terms involved.

## 2.4 Construction of vectors by emergence

The approach by emergence is free from any thesaurus and vectors of concept as bases departure. Only  $d$  the vector size is fixed *a priori*. The construction method of the vectors is identical to the traditional model with the difference that if one of the vectors entering the sum is non-existent, because not yet calculated, then this vector is drawn randomly. The computing process is reiterated until convergence of each vector.

As we show in a more detailed way in [17], there is a certain number of advantages to use this model. The first of them is to be able to freely choose the quantity of resources which one wishes to use by choosing the size of the vectors in a suitable way. To give an idea of the importance of this choice, a base of 500000 vectors of dimension 1000 is approximately 2Go, of size 2000, 4Go, ... As it would not be then reasonable

nor easy to define a concept set of the size chosen, It is easier to seek an approach which enables us to avoid it. Moreover, what can seem a makeshift or at least a compromise proves to be an advantage because the lexical density in space of the words calculated by emergence is much more constant than in a space where concepts are predefined. Indeed, the resources (dimensions of space) have tendency to be harmoniously distributed according to the lexical richness.

### **3 Hybrid modelisation of meaning : conceptual vectors and lexical networks**

#### **3.1 Contribution of the lexical networks to the conceptual vectors**

As shown in [18], distances computed on vectors are influenced by shared components and/or disinct components. Angular distance is a good tool for our aims because of its mathematical characteristics, its simplicity to understand and to linguistically interpret and futhermore it is effective for computational process. Whatever is the chosen distance, used on this kind of vectors (representing ideas and not term occurrences), the lower the distance is, the more the lexical objects are in the same semantic field (isotopy as said by Rastier [19]).

In the framework of semantic analysis as the one which interests us, we use angular distance to benefit from mutual information carried by conceptual vectors to make lexical disambiguation on words whose meanings are in close semantic fields. Thus, "*Zidane scored a goal.*" can be disambiguated thanks to common ideas about sport while "*The lawyer pleads at the court.*" can be disambiguated thanks to those of justice. Furthermore, for prepositionnal attachments, vectors can permit in "*He saw the girl with the telescope.*" to attach "*with a telescope*" to the verb "*saw*" due to ideas about vision.

On the contrary, conceptual vectors can't be used to disambiguate terms which are in different semantic fields. We can even note that an analysis only based on them can lead to misinterpretation. For example, the French noun '*avocat*' has two meanings. It is the equivalent of '*lawyer*' and the equivalent of '*avocado*'. In the French sentence "*L'avocat a mangé un fruit.*", "*The lawyer has eaten a fruit*", '*to eat*' and '*fruit*' carry idea of '*food*' then the acception computed by conceptual vectors for '*avocat*' will be '*avocado*'. It would have been necessary that the knowledge "*a lawyer is a human*" and "*a human eats*" can be identified, something that is not possible with only conceptual vectors. Alone, they are not sufficient to exploit lexical functions instanciations in the texts, a lexical network can thus contribute to correct these shortcomings. These limits were shown in experiments for the semantic analysis using ant algorithms in [20].

#### **3.2 Contribution of conceptual vectors to lexical networks**

If they benefit of an unquestionable precision, the recall of networks is poor. It is, indeed, difficult to think that one could represent all the relations between the terms. Indeed, how can we represent the fact that two terms are in the same semantic field? They may be absent from the network because they may not be connected by "traditional"

arcs. The introduction of arcs of the type "semantic field" also seems problematic for us because of two reasons implicated the fuzzy and flexible of this relation :

- the first one is related to the database conceptor's idea on this relation, when to consider that two synsets are in the same semantic field? In an unfavourable case, there would be very few arcs while in an opposite case we could have a combinative explosion of the number of arc;
- the second problem, more fundamental, is related to representation itself. How to plan to represent by a discrete element a fuzzy relation by essence of the continuous field?

Thus, the continuous domain offered by conceptual vectors gives flexibilities that the discrete domain offered by the networks cannot. They are able to bring closer words on minority ideas but however common what it is not possible with a network. The conceptual vectors and the operation of thematic distance can correct the weak recall inherent of the lexical networks. The defects of the ones are thus mitigated by qualities of the others what makes therefore conceptual vectors and lexical networks complementary tools.

## 4 Expérience on WordNet : usage of data

### 4.1 Exploitation of definitions

*EXtended WordNet* [21] is a project carried out to *Southern Methodist University* of Dallas (Texas, USA) which has two aims:

- to disambiguate terms used in the definitions of the synsets i.e. to indicate which are the synsets employed in the definition;
- to transform these definitions into logical form to allow more easy calculations.

These data were built semi-automatically using information from the network (for example if the genus of the definition, within the meaning of Aristote, has a meaning which is also an hyperonym of the defined synset, it is considered that the meaning of the genus is this hyperonym), of distances between definitions or information about the domain. These data are partly manually controlled and the rate of precision of more than 90%.

For the conceptual vectors construction, we used these data as logical form because they make it possible to locate the most important elements of the definition in particular the genus. Calculation is done thus on a dependency tree manufactured starting from pretreated definition to remove the metalanguage not easily exploitable for a thematic analysis. In our explanations, we will use the example of the logical form of the definition of *ant*.

$ant : NN(x1) - > social : JJ(x1) insect : NN(x1) live : VB(e1, x1, x3)$   
 $in : IN(e1, x2) organized : JJ(x2) colony : NN(x2)$

There is 3 sets :  $x1 = \{social, insect\}$ ,  $x2 = \{organised, colony\}$  and  $e1 = \{live\}$ . This last and *in* make it possible to organise the sets as a hierarchy. The vector of each one of these sets is calculated making the vectorial sum of the item which carry most of



the meaning of this set (verbs, VB; nouns, NN) and half of the ones of the dependents (adverbs, RB; adjectives, JJ). The computation of the global vector is done then by weighted vectorial sum of the various sets in the tree in starting with the lowest part. This mode of calculation makes it possible to consider in a dominating way the genus on the other terms of the definitions and in a more general ways the heads on their syntactic dependent. The figure 1 synthesizes this calculation. No predicate is the set x3 then it does not appears on the figure.

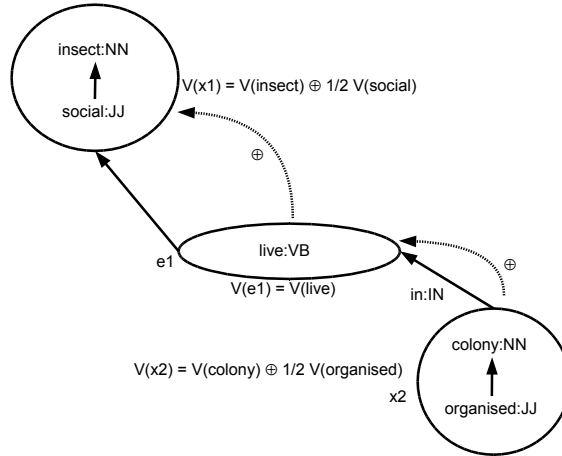


Fig. 1. Construction of a conceptual vector from a definition : example of ant

## 4.2 Exploitation of relations

The exploitation of the relations is done at two level : (1) for the vector construction, they build in a complementary way to the definitions the vector of a synset; (2) to avoid phenomenon of regrouping of distinct sets.

**Vectors Construction** The construction of a conceptual vector is done for each node of the network by simple weighted normalised sum of the vectors of the linked nodes. If  $N$  is a node linked to  $k$  nodes  $N_1 \dots N_k$ , the vector of  $N$  is

$$V(N) = p_1 V(N_1) + p_2 V(N_2) + \dots + p_k V(N_k) \quad (6)$$

This approach naturally involves an agglomeration of the vectors. It is thus necessary to increase the contrast of one vector following its computation. With this intention, one calculates the coefficient of variation <sup>2</sup> of  $V$ . If this last is not around 10% of the average CV the vector undergoes a nonlinear operation of amplification (exponentiation of each component then normalisation), and this in a reiterated way until obtaining a coefficient of variation in the acceptable values. This last was estimated starting from predefined concepts.

<sup>2</sup> the coefficient of variation CV is given by the formula  $\frac{EC(V)}{\mu(V)}$  with  $EC(V)$  the standard deviation of the vector  $V$  and  $\mu(V)$  the arithmetic mean of the components of  $V$ .

**Phenomenon of regrouping of distinct sets** A last potential problem is that the vectors of two distinct sets (at the same time for the lexical network and for thematic) of terms can occupy the same area of space. Computation is done by activation and vectors are randomly drawn at initialization, then that can occur by accident. It is thus necessary to "separate" the close vectors but corresponding however to very different parts of the lexical network and of thematic.

The phenomenon detection is done by examination of the neighbourhood of a conceptual vector. If among the  $N$  first neighbors, the density of words with no correlation with the target word is important then an action of separation must be undertaken.

This action of separation consists in plunging the whole network in field where the nodes tend to be pushed back. In directly being inspired by physics, a force of repulsion in  $1/d^2$  is calculated iteratively between nodes. For a given node, one can thus calculate a vector displacement which will move away it from nodes to which it is too near. Nodes not bringing closer by thematic neighbourhood (at the time of the first phase of calculation cf. section 4.1) but being close "accidentally" end thus naturally up separating.

## 5 Conclusion

In this article, we presented the conceptual vectors built by emergence. We showed in what they can help to solve the tennis problem from their character complementary to the lexico-semantic networks one whose most famous example in current research is WordNet. Indeed, the recall of the networks is weak, they easily do not make it possible to represent the semantic fields contrary to the vectors while the latter are not sufficient to represent relations like hyperonymy or meronymy.

Our proposal is to benefit from this complementarity while adding to WordNet conceptual vectors built starting from definitions and relations contained in this base. The method suggested here holds of the continuous field contrary to the whole methods we studied in the literature which belong to the discrete field (addition of arcs for the relations, symbols about the domain, etc).

We are aware that this method only makes it possible to solve part of the problem of tennis. Indeed, the conceptual vectors do not allow to exhibit not-thematic collocational relationship between items. They are primarily the relations that Igor Mel'čuk models with his syntagmatic lexical functions [22] like the intensification ("great fear"; *Magn* (‘fear’) = ‘great’)), name of center ("crux of the problem"; *Centr* (‘problem’) = ‘crux’) or even the confirmator ("legitimate excuse"; *Ver* (‘excuse’) = ‘legitimate’). As notices [10], these relations belong to those which would probably be necessary to have in a lexical base. We share this point of view, some tracks were explored in [23] and currently continue to be followed.

## References

1. Quillian, R.: Semantic memory. In: Semantic Informatic processing. MIT Press (1968) 227–270
2. Mihalcea, R., Tarau, P., Figa, E.: Pagerank on semantic networks, with application toward sense disambiguation. In: COLING'2004 : 20th International Conference on Computational Linguistics, Geneva, Switzerland (2004) 1126–1132

3. Mangeot-Lerebours, M., Sérasset, G., Lafourcade, M.: Construction collaborative d'une base lexicale multilingue : Le projet papillon. *TAL (Traitement Automatique des langues) : Les dictionnaires électroniques* **44** (2003) 151–176
4. Knight, K., Luk, S.: Building a large-scale knowledge base for machine translation. In: *AAAI'1994 : National Conference on Artificial Intelligence*, Stanford University, Palo Alto, California (1994)
5. Harabagiu, S., Chai, J., eds.: *Usage of WordNet in Natural Language Processing Systems*, Université de Montréal, Montréal, Canada (1998)
6. Fellbaum, C., ed.: *WordNet: An Electronic Lexical Database*. The MIT Press (1988)
7. Harabagiu, S.M., Miller, G.A., Moldovan, D.I.: Wordnet 2 - a morphologically and semantically enhanced resource. In: *Workshop SIGLEX'99 : Standardizing Lexical Resources*. (1999) 1–8
8. Agirre, E., Ansa, O., Martinez, D., Hovy, E.: Enriching wordnet concepts with topic signatures. In: *NAACL workshop on WordNet and Other Lexical Resources: Applications, Extensions and Customizations*, Pittsburg, USA (2001)
9. Stevenson, M.: Augmenting noun taxonomies by combining lexical similarity metrics. In: *COLING'2002 : 19th International Conference on Computational Linguistics*. Volume 2/2., Taipei, Taiwan (2002) 953–959
10. Ferret, O., Zock, M.: Enhancing electronic dictionaries with an index based on associations. In: *Proceedings of the 21st International Conference on Computational Linguistics*, Sydney, Australia, Association for Computational Linguistics (2006) 281–288
11. Salton, G., McGill, M.: *Introduction to Modern Information Retrieval*. McGrawHill, New York (1983)
12. Deerwester, S.C., Dumais, S.T., Landauer, T.K., Furnas, G.W., Harshman, R.A.: Indexing by latent semantic analysis. *Journal of the American Society of Information Science* **41** (1990) 391–407
13. Chauché, J.: Détermination sémantique en analyse structurelle : une expérience basée sur une définition de distance. *TAL Information* **31/1** (1990) 17–24
14. Larousse, ed.: *Thésaurus Larousse - des idées aux mots, des mots aux idées*. Larousse (1992)
15. Kirkpatrick, B., ed.: *Roget's Thesaurus of English Words and Phrases*. Penguin books, London (1987)
16. Zock, M.: Sorry, what was your name again, or how to overcome the tip-of-the tongue with the help of a computer? In: *SemaNet'02: Building and Using Semantic Networks*, Taipei, Taiwan (2002)
17. Lafourcade, M.: Conceptual vector learning - comparing bootstrapping from a thesaurus or induction by emergence. In: *LREC'2006*, Genoa, Italia (2006)
18. Besançon, R.: *Intégration de connaissances syntaxiques et sémantiques dans les représentations vectorielles de texte* (2001)
19. Rastier, F.: *L'isotopie sémantique, du mot au texte* (1985)
20. Lafourcade, M., Guinand, F.: Ants for natural language processing. *International Journal of Computational Intelligence Research* (2006) À paraître.
21. Mihalcea, R., Moldovan, D.: extended wordnet: progress report. In: *NAACL 2001 - Workshop on WordNet and Other Lexical Resources*, Pittsburgh, USA (2001)
22. Mel'čuk, I., Clas, A., Polguère, A.: *Introduction à la lexicologie explicative et combinatoire*. Duculot (1995)
23. Schwab, D.: Approche hybride - lexicale et thématique - pour la modélisation, la détection et l'exploitation des fonctions lexicales en vue de l'analyse sémantique de texte. (2005)