

Inductive and deductive inferences in a Crowdsourced Lexical-Semantic Network

Manel ZARROUK
LIRMM

manel.zarrouk@lirmm.fr

Mathieu LAFOURCADE
LIRMM

lafourcade@lirmm.fr

Alain JOUBERT
LIRMM

alain.joubert@lirmm.fr

Abstract

In Computational Linguistics, building lexical-semantic networks and validating contained relations are paramount issues as well as adding some reasoning skills in order to enrich these knowledge bases. In this paper we devise an inference engine which aims at producing new "potential" relations from already existing ones in the JeuxDeMots network. This network is constructed with the help of a GWAP (game with a purpose) thanks to thousands of players. It handles terms and weighted relations between these terms, and currently contains over 2 million relation occurrences. Polysemous terms may be refined in several senses (*bank* may be a *bank*>*financial institution* or a *bank*>*river*) but as the network is indefinitely under construction (in the context of a Never Ending Learning approach) some senses may be missing at a given time. The approach we proposed here is founded on the triangulation method through two kinds of inference schemes: deduction (top-down from generic to specific terms) and induction (bottom-up from specific to generic terms). A blocking mechanism, whose purpose is to avoid proposing highly dubious new relations, is based on logical and statistical constraints. Automatically inferred relations are then proposed to human contributors to be validated. In case of invalidation, a reconciliation dialog is undertaken to identify the cause of the wrong inference: an exception, an error in the premises or a previously undetected confusion due to polysemy on the central term common to both premises.

1 Introduction

Developing resources in NLP is one of the crucial issue of the field. Most of the existing lexico-semantic networks have been constructed manually, like for instance the famous WordNet. Of course some tools are generally designed for consistency checking, but nevertheless the task remains time consuming and costly. Fully automated approaches are generally limited to term cooccurrences as extracting precise semantic relations between terms from text remains really difficult. New approaches involving crowdsourcing are flowering in NLP especially with the advent of Amazon Mechanical Turk or in a broader scope Wikipedia and Wiktionary, to cite the most well known examples. WordNet ((?) and (?)) is such a lexical network based on synsets which can be roughly considered as concepts. EuroWordnet (?) a multilingual version of WordNet and WOLF (?) a French version of WordNet, were built by automated crossing of WordNet and other lexical resources along with some manual checking. (?) constructed automatically BabelNet a large multilingual lexical network from term cooccurrences in the Wikipedia encyclopedia.

A highly lexicalized lexical-semantic network can contain concepts but also plain words (and multi-word expressions) as entry points (nodes) along with word meanings. The idea itself of *word senses* in the lexicographic tradition may be debatable in the context of resources for semantic analysis, and we generally prefer to consider word usages. By *word usages* we mean refinements of a given word which is clearly identified by locutors. A polysemic term has several usages that might differ substantially from word senses as classically defined. A given usage can also in turn have several deeper refinements and the whole set of usage can take the form of a de-

cision tree. In the context of a collaborative construction, such a lexical resource should be considered as being constantly evolving and a general rule of thumb is to have no definite certitude about the state of an entry.

The building of a collaborative lexical network (or any similar resource) can be devised according to two broad strategies. First, it can be designed as a contributive system like Wikipedia where people willingly add and complete entries (like for Wiktionary). Second, contributions can be made indirectly thanks to games (better known as GWAP (?) and (?)) and in this case players do not need to be aware that while playing they are helping building a lexical resource. In any case, the lexical network that is built is not free of errors which are corrected along their discovery. Thus a large number of obvious relations are not contained in the lexical network but are indeed necessary for a high quality resource usable in various NLP application and notably semantic analysis. For example, contributors do not indicate that a particular bird type can fly, as it is considered as an obvious generality. Only notable facts which are not easily deductible are naturally contributed. Well known exceptions are also generally contributed and take the form of a negative weight for the relation (for example, $fly \xrightarrow{agent:-100} ostrich$).

In order to consolidate the lexical network, we adopt a strategy based on a simple (if not simplistic) inference mechanism to propose new relations from those existing. The approach is strictly endogenous as it doesn't rely on any other external resources. Inferred relations are submitted either to contributors for voting or to expert for direct validation/invalidation. A large percentage of the inferred relations has been found to be correct. However, a non negligible part of them are found to be wrong and understanding why is both relevant and useful. The explanation process can be viewed as a reconciliation between the inference engine and the validator who is guided through a dialog to explain why he found the considered relation as incorrect. A wrong inferred relation may come from three possible origins: false premises used by the inference engine, exception or confusion due to polysemy.

In this article, we first present the principles behind of lexical network construction with

crowdsourcing and *games with a purpose* (also know as human-based computation games) and illustrated them with the JeuxDeMots (JDM) project. Then, we present the outline of an *elicitation engine* based on an *inference engine* using deduction and induction schemes and a *reconciliation engine*. An experimentation is then reported on the performances of the system.

2 Lexical Network and Crowdsourcing

There are many ways for building a lexical network considering some crucial factors as the quality of data, cost and time. Beside manual or automated strategies, contributive approaches are more and more popular as they are both cheap to set up and efficient in quality. More specifically, there is an increasing trend of using on-line GWAPs ((?) and (?)) method for feeding such resources.

The JDM lexical network is constructed through a set of on-line associative games. In these games, players are appealed to contribute on lexical and semantic relations between terms or verbal expressions which are presented in the network by the arcs interconnecting nodes in a graph. The informations in the JDM network are gathered by an unnegotiated crowd agreement (classical contributive systems rely on a negotiated crowd agreement).

2.1 JeuxDeMots: a GWAP for Building a Lexical-Semantic Network

JeuxDeMots¹ is a two player GWAP which aims to build a large lexical-semantic network (?). The network is composed of terms (as vertices) and typed relations (as links between vertices). It contains terms and possible refinements in a similar way to the WordNet synset (?). The semantic network is constructed by connecting terms by typed and weighted relations, validated by pairs of players. These relations are labelled according to the instructions given to the players and weighted according to the number of pairs of players who choose them. Other Web-based systems exist, such as Open Mind Word Expert (?), which aims at creating large sense-tagged corpora with the help of Web users, and SemKey (?) which makes use of WordNet and Wikipedia to disambiguate lexical forms referring to concepts, thus identifying semantic keywords.

¹<http://jeuxdemots.org>

2.2 Diko as a Contributive Tool

Diko² is a web based tool for displaying the information contained in the JDM lexical network which can also be used as a contributive tool. The necessity to not rely only on the JDM game for building the lexical network comes from the fact that many relation types of JDM are either difficult to grasp for a casual player or not very productive (not many terms can be associated).

The principle of the contribution process is that a proposition made by a user will be voted pro or con by other users then included or excluded by an expert validator. What we propose in this paper falls under this type of scenario of contributions/validations.

3 Elicitation by Inference and Reconciliation

We designed a system for augmenting the number of relations in the JDM lexical network having two main components: (a) an inference engine and (b) a reconciliator. The inference engine proposes relations as if it was a contributor, to be validated by human contributors or experts. In case of invalidation of an inferred relation, the reconciliator is invoked to try to assess why the inferred relation was found wrong. Elicitation here should be understood as the process to transform some implicit knowledge of the user into explicit relations in the lexical network.

3.1 Making Inferences

The core ideas about inferences in our system are the following:

- for the engine, inferring is to derive new premises (under the form of relations between terms) from previously known premises, which are existing relations;
- candidate inferences may be logically blocked on the basis of the presence or absence of some other relations;
- candidate inferences can be filtered out on the basis of a strength evaluation.

3.1.1 Deduction Scheme

In this paper, the first type of inference we are working with is the deduction or top-down scheme, which is based on the transitivity of the ontological relation *is-a* (hypernym). If a term A

²<http://www.jeuxdemots.org/diko.php>

is a kind of B and B holds some relation R with C, then we can expect that A holds the same relation with C. The scheme can be formally written as follows:

$$\exists A \xrightarrow{is-a} B \wedge \exists B \xrightarrow{R} C \Rightarrow A \xrightarrow{R} C$$

Global processing - Let us consider a term T with a set of weighted hypernyms. From each hypernym, the inference engine deduces a set of inferences. Those inference sets are not disjoint in the general case, and the weight of an inference proposed in several sets is the incremental geometric mean of each occurrence.

Logical filtering - Of course, this scheme above is far too naive, especially considering the resource we are dealing with. In effect, B is possibly a polysemous term and ways to block inferences that are certainly wrong can be devised. If there are two distinct meanings of the term B that hold respectively the first and the second relation, as in the Figure ?? below, then most probably the inference is wrong.

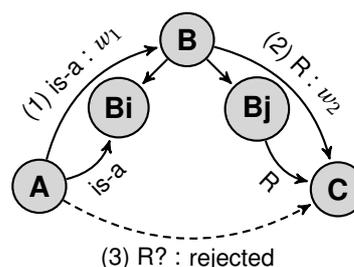


Figure 1: Triangular inference scheme with logical blocking based on the polysemy of B.

In this case, a relation R -to be inferred- must fulfill some constraints as formulated below:

$$\begin{aligned} & A \xrightarrow{is-a} B \wedge B \xrightarrow{R} C \\ \wedge & (\exists B_i \xrightarrow{meaning-of} B \wedge \exists B_j \xrightarrow{meaning-of} B) \\ \wedge & (\exists A \xrightarrow{is-a} B_i \vee \exists B_j \xrightarrow{R} C) \\ \Rightarrow & A \xrightarrow{R} C \end{aligned}$$

Moreover, if one of the premises is tagged as *true but irrelevant*, then the inference is blocked. **Statistical filtering** - It is possible to evaluate a confidence level (on an open scale) for each produced inference, in a way that dubious inferences can be filtered out. The weight w of an inferred relation is the geometric mean of the weight of the premises (relations (1) and (2) in Figure ??). If the second premise has a negative value, the weight is not a number and the proposal is discarded. As the geometric mean

is less tolerant to small values than the arithmetic mean, inferences which are not based on two rather true relations (premises) are unlikely to pass.

$$w(A \xrightarrow{R} C) = (w(A \xrightarrow{is-a} B) * w(B \xrightarrow{R} C))^{1/2}$$

$$\Rightarrow w_3 = (w_1 * w_2)^{1/2}$$

3.1.2 Induction Scheme

As for the deductive inference, induction exploits the transitivity of the relation *is-a*. If a term *A* is a kind of *B* and *A* holds a relation *R* with *C*, then we might expect that *B* could hold the same type of relation with *C*. More formally we can write:

$$\exists A \xrightarrow{is-a} B \wedge \exists A \xrightarrow{R} C \Rightarrow B \xrightarrow{R} C$$

This scheme is a generalization inference. The **global processing** is similar to the one applied to the deduction scheme and similarly some logical and statistical filtering may be undertaken.

The term joining the two premises (called *central term*, in this case term *A*) is possibly polysemous. If the term *A* is presenting two distinct meanings which hold respectively the premises (as shown in Figure ??), then the inference done from that term may be probably wrong.

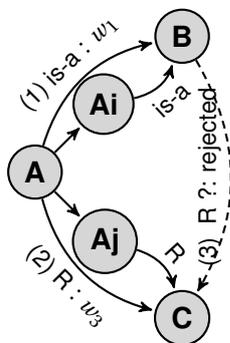


Figure 2: (1) and (2) are the premises, and (3) is the logical induction proposed for validation. Central term *A* may be polysemous with meanings holding premises, thus inducing a probably wrong relation.

Logical filtering can be formalized as follows:

$$\wedge (\exists A_i \xrightarrow{\text{meaning-of}} A \wedge A \xrightarrow{is-a} B) \wedge (\exists A_j \xrightarrow{\text{meaning-of}} A \wedge A \xrightarrow{R} C)$$

$$\wedge (\exists A_i \xrightarrow{is-a} B \vee \exists A_j \xrightarrow{R} C) \Rightarrow B \xrightarrow{R} C$$

Statistical filtering is possible, as for the deductive scheme to evaluate a confidence level.

According to the weight evaluation from the deductive diagram, the estimated weight for the induced relation is:

$$w(B \xrightarrow{R} C) = (w(A \xrightarrow{R} C))^2 / w(A \xrightarrow{is-a} B)$$

$$\Rightarrow w_2 = \frac{(w_3)^2}{w_1}$$

3.2 Performing reconciliation

Inferred relations, further to both induction and deduction, are presented to the validator to decide of their status: *rather true*, *rather true but irrelevant*, *possible* or *mostly false*. In case of invalidation, a reconciliation procedure is committed in the purpose to try to diagnose the reasons: error in one of the premises (previously existing relations are false), exception or confusion due to polysemy (the inference has been made on a polysemous central term) and initiates a dialog with the user. The latter is free to choose to pursue the dialog partially, entirely or to choose not to start it. To know in which order to proceed, the reconciliator determines if the weights of the premises are rather strong or weak. This confidence is done by comparing the relation weight to a confidence threshold which is computed as the starting point of the long tail in the distribution of the relation. For the whole set of the outgoing relations from a term the long tail starts at the point where the cumulated weights of the relations of the tail is equal the cumulated weights of the relations which do not belong to the tail (?).

- If $w(A \xrightarrow{is-a} B) \geq conf - thr(A) \Rightarrow$ trusted relation
- If $w(A \xrightarrow{is-a} B) < conf - thr(A) \Rightarrow$ dubious relation

In the case we have both relations (1) and (2) as trusted, the reconciliator tries, by initiating a dialog with the validator, to check at first if the relation inferred is an exception. If not, it proceeds by checking if term *B* is polysemous and finally checks if it is an error case. We check the error case in the final step because the confidence level of relations (1) and (2) made them trusted.

In the case of having a dubious relation either for (1) and (2), the reconciliator suspects that it is an error case and this relation was the cause of a wrong inference. So, the validator is asked to confirm or to disprove it. In case of refutation of one of the relations, we have an error. If not, we proceed with checking if it's an exception case or a polysemy.

3.2.1 Errors in the premises

In this case, suppose that relation (1) has a relatively low weight. The reconciliator asks the validator about the relation (1).

- If is false, a negative weight is attributed to (1) and the reconciliation is completed. As such, this relation will not be used later on as premises on further inferences;
- If it is true, ask if relation (2) is true and proceed as above if the answer is negative;
- Otherwise, move to checking the other cases (exception, polysemy).

3.2.2 Errors as exceptions

For the *deduction*, if the validator indicates that the inferred relation is an exception relatively to the term B , the relation is stored in the lexical network with a negative weight along with a meta-information which indicates that it is an exception.³

For the *induction*, if the alidator indicates that the relation ($A \xrightarrow{R} C$) (which served as premise) is an exception relatively to the term B , in addition to storing the false inferred relation ($B \xrightarrow{R} C$) in the network with a negative weight, the relation ($A \xrightarrow{R} C$) is tagged with a meta-information indicating it as an exception. In the induction case, the exception is a true premise which leads to a false induced relation.⁴

In both cases of induction and deduction, the *exception* tag concerns always the relation ($A \xrightarrow{R} C$). Once this relation is tagged as an exception, it will not participate as a premise in inferring generalized relations (bottom-up model) but can still be used in inducing specified relations (top-down model).

3.2.3 Errors due to Polysemy

In this case, if the middle term (B for deduction and A for induction) presenting a polysemy is mentioned as polysemous in the network, the refinement terms $term_1, term_2, \dots, term_n$ are presented to the validator so he can choose the

³For example, suppose we have ($ostrich \xrightarrow{agent} fly$) inferred by *deduction* with the central term B . In this case, it's true that an ($ostrich \xrightarrow{is-a} bird$) and that a ($bird \xrightarrow{agent} fly$), but the inferred relation an *ostrich can fly* is *false* and it is considered as an *exception* considering the central term "*bird*".

⁴As for the relation ($fish \xrightarrow{agent} fly$) which is a false inferred relation based on the central term *exocet*. The ($exocet \xrightarrow{is-a} fish$) and ($exocet \xrightarrow{agent} fly$) are true but the latter one is an *exception* in the form of a *true* relation.

appropriate one. The validator can propose new terms as refinements if he is not satisfied with the listed ones (inducing the creation of new appropriate refinements). After this procedure, two new relations ($A \xrightarrow{is-a} B_i$ and $B_j \xrightarrow{R} C$ in the case of deduction, or $A_i \xrightarrow{is-a} B$ and $A_j \xrightarrow{R} C$ in the induction case) will be included in the network with positive values and the inference engine will use them later on as premises.

4 Experimentation

We made an experiment with a unique run of the engine over the lexical network of JDM. The purpose is to measure the production of the inference engine along with the blocking and filtering. From the set of supposedly valid inferred relations (both by induction and deduction), we took a random sample of 400 propositions for each relation type and undertook the validation/reconciliation process. The experiment conducted is for evaluation purpose only, as actually the system is running iteratively along with contributors and games.

4.1 Unleashing the Inference Engine

We applied the inference engine on around 23 000 randomly selected terms having at least one hypernym or one hyponym and thus produced by deduction 1 484 209 inferences (77 089 more were blocked). The threshold for filtering was set to a weight of 25. This value is relevant as when a human contributor proposed relation is validated by an expert, it is introduced with a default weight of 25. For induction, the inference engine produced 353 371 relation candidates. The table ?? presents the number of relations proposed by the inference engine through deduction. The different types for the second premise are variously productive. Of course, this is mainly due to the number of existing relations and the distribution of their type in the network.

The transitive relation *is-a* is the less productive which might seems surprising at first glance. In fact, this relation is already quite populated in the network, and as such, fewer new relations can be inferred. The figures are inverted for some other relations that are not so well populated but still are potentially valid. The agent semantic role (the *agent-1* relation) is by far the most productive, with 30 time more propositions than what currently exists in the lexical network.

Relation type	Proposed	Blocked	Filtered
is-a	91k (6,1)	4 k (5.2)	53 k (26,3)
has-parts	372k (25.1)	31 k (40.7)	100 k (49.3)
holonym	108k (7.2)	17 k (23.3)	26 k (13.2)
place	271k (18.3)	11 k (15)	14 k (7)
charac	203k (13.7)	2 k (3.4)	6 k (3.2)
agent-1	198k (13.3)	9 k (11.7)	1122 (0.5)
instr-1	24k (1.7)	127 (0.2)	391 (0.2)
patient-1	14k (1)	7 (0.01)	13 (0)
place-1	145k (9.8)	129 (0.2)	206 (0.1)
place >action	50k (3.4)	91 (0.1)	132 (0.06)
obj >mater	4k (0.3)	135 (0.2)	262 (0.1)
Total	1 484k	77 k	203 k

Table 1: Numbers and percentages for inferences (proposed, blocked or filtered) by the deduction.

4.2 Figures on Reconciliation

Table ?? contains some evaluation of the status of the inferences proposed by the inference engine through deduction. Inferences are valid for an overall of 80-90% with around 10% valid but not relevant (like for instance *dog* $\xrightarrow{has-parts}$ *proton*). We observe that error number in premises is quite low, and nevertheless errors can be easily corrected. Of course, not all possible errors are detected through this process. The reconciliation allows in 5% of the cases to identify polysemous terms. Globally false negatives (inferences voted false but are true) and false positives (inferences voted true but are false) are evaluated to less than 0,5%. For the induction process (table ??), the relation *is-a* is not obvious (a lexical network is not reducible to an ontology and multiple inheritance is possible). Result seems about 5% better than for the deduction process: inferences are valid for an overall of 80-95%. The error number is very low. The main difference with the deduction process is on errors due to polysemy which is lower with the induction process.

5 Conclusion

We presented some issues about inferring new relations from existing ones in a contributed lexical-semantic network in which word usages are discovered incrementally along its construction. Errors are naturally present as they might originate from games played on difficult relations, but they are usually spotted and corrected by contributors for terms they are interested in. To be able to enhance the network quality, we proposed an elicitation engine based on inferences and reconciliations. Inferences are here proposed with two different schemes (induction and deduction), along with a logical blocking and statistical filtering. If an inferred relation is proven wrong, a reconciliation is conducted

to identify the underlying cause. As global figures, we can conclude that inferred deductive relations are correct and relevant in about 78% of the cases and correct but irrelevant in 10% of the case. Overall wrong deductive inferences is about 12% with at least one error in the premises of about 2%, exceptions about 5% and polysemy confusion about 5%. Induction is naturally less productive but more reliable. Beside a tool for increasing relations in a lexical network, the elicitation engine is both an error detector and a polysemy identifier. Actions taken during the reconciliation forbid an inference proven wrong or exceptional to be proposed again. Such an approach should be pushed forward with other types of inference scheme like abduction, and possibly with distribution evaluation of term semantic classes on which inferences are conducted. Indeed, some classes like concrete objects or living beings may be substantially more productive for certain relation types than abstract nouns of processes or events. Anyway, such discrepancies of inference productivity between classes are worthy to investigate further.

Deduction Relation type	% valid		% error		
	rlvt	¬rlvnt	prem	excep	pol
is-a	76%	13%	2%	0%	9%
has-parts	65%	8%	4%	13%	10%
holonym	57%	16%	2%	20%	5%
typical place	78%	12%	1%	4%	5%
charac	82%	4%	2%	8%	4%
agent-1	81%	11%	1%	4%	3%
instr-1	62%	21%	1%	10%	6%
patient-1	47%	32%	3%	7%	11%
typical place-1	72%	12%	2%	10%	6%
place >action	67%	25%	1%	4%	3%
object >mater	60%	3%	7%	18%	12%

Table 2: Results of the validation/reconciliation according to relation types in the deduction. Valid relations can be relevant or not, and errors can be in **premises**, **exceptions** or **polysemy**.

Induction Relation types	% valid		% error		
	rlvt	¬rlvnt	prem	excep	pol
has-parts	78%	10%	3%	2%	7%
holonym	68%	17%	2%	8%	5%
typical loc	81%	13%	1%	2%	3%
carac	87%	6%	2%	2%	3%
agent-1	84%	12%	1%	2%	1%
instr-1	68%	24%	1%	4%	3%
patient-1	57%	36%	3%	2%	2%
typical loc-1	75%	16%	2%	5%	2%
lieu-action	67%	28%	1%	3%	1%
object mater	75%	10%	7%	5%	3%

Table 3: Results of the validation/reconciliation according to relation types in the induction. The relation *is-a* is inappropriate for Induction.

References

- von Ahn, L. and Dabbish, L. (2008) *Designing games with a purpose*. Communications of the ACM, number 8, volume 51. pp. 58-67.
- Chamberlain, J., Fort, K., Kruschwitz, U., Lafourcade, M. and Poesio, M. (2013) (2013) Using Games to Create Language Resources: Successes and Limitations of the Approach. Gurevych, Iryna; Kim, Jungi (Eds.), Springer, ISBN 978-3-642-35084-9, 2013, 42 p.
- Fellbaum, C. (1988, ed.) *WordNet: An Electronic Lexical Database*. The MIT Press.
- Law, E., Luis von Ahn, L. and Mitchell, T. (2009) *Search war: a game for improving web search*. KDD Workshop on Human Computation 2009. 31p.
- Law, E., Luis von Ahn, L., Dannenberg, R. B. and Crawford, M.. (2007) *TAgATune: A Game for Music and Sound Annotation*. ISMIR 2007. pp. 361-364.
- Lafourcade, M. and Joubert, A. (2012) *Long Tail in Weighted Lexical Networks*. In proc of Cognitive Aspects of the Lexicon (CogAlex-III), COLING, Mumbai, India, December 2012. 16 p.
- Lafourcade, M. (2007) *Making people play for Lexical Acquisition*. In Proc. SNLP 2007, 7th Symposium on Natural Language Processing. Pattaya, Thailande, 13-15 December. 8 p.
- Marchetti, A., Tesconi, M., Ronzano, F., Mosella, M and Minutoli, S. (2007) *SemKey: A Semantic Collaborative Tagging System*. in Procs of WWW2007, Banff, Canada. 9 p.
- Mihalcea, R. and Chklovski, T. (2003) *Open Mind-Word Expert: Creating large annotated data collections with web users help..* In Proceedings of the EACL 2003, Workshop on Linguistically Annotated Corpora (LINC). 10 p.
- Miller, G.A., Beckwith, R., Fellbaum, C., Gross, D. and Miller, K.J. (1990) *Introduction to WordNet: an on-line lexical database*. International Journal of Lexicography. Volume 3, pp. 235-244.
- Miller, G.A. (1995) *WordNet: A Lexical Database for English*. Communications of the ACM Vol. 38, No. 11, pp. 39-41.
- Navigli, R. and Ponzetto, S. (2010) *BabelNet: Building a very large multilingual semantic network*. In Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics, Uppsala, Sweden, 11-16 July 2010. pp: 216-225.
- Navigli, R. and Ponzetto, S. (2012) *BabelNet: The Automatic Construction, Evaluation and Application of a Wide-Coverage Multilingual Semantic Network*. , Artificial Intelligence. 193: pp. 217-150.
- Sagot, B. and Fier, D. (2010) *Construction d'un wordnet libre du français à partir de ressources multilingues*. In Proceedings of TALN 2008, Avignon, France, 2008.12 p.
- Thaler, S., Siorpaes, K., Simperl, E. and Hofer, C. (2011) *A Survey on Games for Knowledge Acquisition*. STI Technical Report, May 2011.19 p.
- Vossen, P. (1998) *EuroWordNet: a multilingual database with lexical semantic networks..* Kluwer Academic Publishers.Norwell, MA, USA. 200 p.