# Automatic Term Relationship Cleaning and Refinement for AGROVOC

Asanee Kawtrakul [a], Aurawan Imsombut [b], Aree Thunyakijjanukit [c], Dagobert Soergel [d],
Anita Liang [e], Margherita Sini [f], Gudrun Johannsen [g], and Johannes Keizer [h]

[a,b,c] *Department of Computer Engineering, Kasetsart University, Bangkok, Thailand, {ak, aurawani}@vivaldi.cpe.ku.ac.th*

[d] *College of Library and Information Services, University of Maryland, College Park, dsoergel@umd.edu*

[e,f,g,h] *Food and Agriculture Organization (FAO) of the United Nations, Library & Documentation Systems Division,
00100 Rome, Italy, {anita.liang, margherita.sini, gudrun.johannsen, johannes.keizer}@fao.org*

**Abstract**

AGROVOC is a multilingual thesaurus developed and maintained by the Food and Agricultural Organization of the United Nations. Like all thesauri, it contains some explicit semantics, which allow it to be transformed into an ontology or used as a resource for ontology construction. However, most thesauri, AGROVOC included, give very broad relationships that lack the semantic precision needed in an ontology. Many relationships in a thesaurus are incorrectly applied or defined too broadly. Accordingly, extracting ontological relationships from a thesaurus requires data cleaning and refinement of semantic relationships.

This paper presents a hybrid approach for (semi-)automatically detecting these problematic relationships and for suggesting more precisely defined ones. The system consists of three main modules: Rule Acquisition, Detection and Suggestion, and Verification. The Refinement Rule Acquisition module is used to acquire rules specified by experts and through machine learning. The Detection and Suggestion module uses noun phrase analysis and WordNet alignment to detect incorrect relationships and to suggest more appropriate ones based on the application of the acquired rules. The Verification module is a tool for confirming the proposed relationships. We are currently trying to apply the learning system with some semantic relationships to test our method.

*Key words*: AGROVOC, Data Cleaning, Semantic Relationship Refinement, Noun Phrase Analysis

## 1 Introduction

The performance of information processing systems may be enhanced when it is supported by ontologies, domain-specific terminologies containing rich and precise semantics. For example, in information retrieval, ontologies may be used for query expansion, for marking up documents at various levels of granularity, for knowledge discovery, etc. Research on (semi-)automatic ontology construction has been conducted using a variety of terminological resources, such as raw text (Hearst 1992, Maedche and Staab 2001, Kiet 2000, and Navigli et al. 2003), dictionaries (Janniak 1999, Kietz 2000, Kang 2001) and thesauri (Soergel et al. 2004, Clark 2000, Wielinga 2001). Each of these sources has different characteristics which require different approaches to term and relationship extraction. Raw text consists of unstructured text containing huge amounts of information that are frequently updated. Dictionaries are semi-structured resources that are infrequently updated; domain dictionaries, in particular, are suitable for extracting terms and their relationships (e.g., hyponyms, meronyms, and synonyms) as well as their definitions. Of the terminological resources considered, thesauri lend themselves most readily to ontology construction because their explicit semantic structure facilitates the natural language processing needed to extract terms and relationships. Our work is to develop and maintain the Agricultural Ontology Service, which will support the construction of an Agricultural Knowledge Portal. Therefore, we use AGROVOC as a resource to build an ontology in the domain of food and agriculture.

AGROVOC is a multilingual agricultural thesaurus developed and maintained by the Food and Agriculture Organization (FAO) of the United Nations. It is used at FAO for indexing and searching information resources within the agricultural domains. However, within AGROVOC, semantic relationships are poorly defined and inconsistently applied. For example, AGROVOC incorrectly[1] uses *NT* (*narrow term*), approximately equivalent to 'superclass of,' or 'hypernym of', in *Milk NT Milk Fat*, while a more specific, and correct, relationship could be 'containsSubstance'. In AGROVOC, *RT (related term)* is underspecified, subsuming numerous relationships; for example, it uses *RT* in *Mutton RT Sheep,* which should be refined to a more specific one, such as 'madeFrom' (Soergel et al. 2004) to distinguish from other uses of RT.

The question of reengineering AGROVOC to an ontology has recently been addressed in a few studies. Fisseha and Liang (2003) present some rough ideas for preparing AGROVOC for conversion into an ontology, such as converting BT/NT to is-a, and refining RT to more specific relationships. Soergel et al. (2004) propose the rules-as-you-go approach, where rules for semantic refinement are identified as experts work on the thesaurus and notice patterns in the occurrence of semantic relationships between terms. Since the patterns and rules are identified through intellectual work, the refinements occur gradually and can deal with only a limited number of patterns. This paper enhances the feasibility of the rules-as-you-go approach by applying machine learning to automatically extract the rules. The learning technique is based on the OntoLearn method (Navigli et al. 2003), the automatic ontology learning system that was used for extracting terms and detecting semantic relationships from a tourism text corpus. It uses inductive machine learning for extracting semantic relationships between the head word and its modifier in compound nouns.

This paper presents a hybrid approach for (semi-)automatically detecting these problematic relationships, especially BT/NT and USE/UF relationships, and suggesting more appropriate ones. In the case of RT relationships, which usually are underspecified relationships, the refinement rules, acquired from experts and machine learning, are applied. The system consists of three main modules: Rule Acquisition Refinement, Detection and Suggestion, and Verification. The Rule Acquisition module is used to train the machine based on rules specified by experts. The Detection and Suggestion module uses noun phrase analysis and WordNet alignment to detect incorrect relationships and to suggest more appropriate ones based on the application of the acquired rules. The Verification module is a tool for confirming the proposed relationships.

Section 2 describes the problems in AGROVOC. Section 3 gives an overview of the system for data cleaning and relationship refinement. Sections 4 and 5 describe the preparation of the rules and an algorithm for cleaning and refinement, respectively. Finally, the experimental results and future works are summarized in Section 6 and Section 7 gives brief conclusions.

## 2 Structural Problems in AGROVOC

AGROVOC has been found to contain relationships that are incorrectly assigned, as indicated in 2.1, and too broadly defined, as shown in 2.2.

### 2.1 Incorrectly assigned relationships

A review of the data in AGROVOC reveals that some USE/UF and BT/NT relationships are incorrect or reflect inconsistent uses of the relationships. The USE/UF relationship may link synonyms and their formal variants but also quasi-synonyms such as closely related and hierarchically related terms (Soergel et al. 2004). Likewise, the BT/NT relationship is very ambiguous (see examples in Table 1).

---

[1] Within a hierarchy based on partitivity, the use of NT would not necessarily be an incorrect one, e.g., *Milk* **NT** *Milk Fat* **NT** *Milk Fat Globule* etc. However, the refined AGROVOC is anticipated to use BT/NT to express hierarchical, super/subclass-type relationships only. And its conversion into an ontology necessitates that each relationship correspond to a unique sense.

Table 1 Examples of   inappropriately defined relationships between terms

| Relationship | Examples | Remark |
|---|---|---|
| **UF** | 1. *Locomotion* UF *Walking* | Incorrect Relationship: *Walking* is not a synonym of *Locomotion*. WordNet shows that *Walking* is the hyponym of *Locomotion*. |
| | 2. *Digestive juices* UF *Chyme* | Incorrect Relationship:  *Digestive juices* is not a synonym of Chyme, and the two terms have different hypernyms in WordNet. |
| **BT/NT** | 1. *Milk* NT *Milk fat* | Incorrect Relationship: *Milk* <containsSubstance> *Milk fat*. |
| | 2. *Portugal* BT *Western Europe* | Incorrect  Relationship: *Portugal* <spatiallyIncludedin> *Western Europe* |

## *2.2 Vaguely defined (underspecified) relationships*

The relationships used in AGROVOC consist of at least three types: UF/USE, BT/NT and RT.  Because they are very generally defined (cf. printed version of AGROVOC, Fourth Ed., pp. xv-xvii), they have been applied inconsistently.  RT, in particular, has been used to link any two, usually non-hierarchically related terms that are felt to be associated with each other.   Further refinements to this relationship are needed to reflect the more meaningful and specific associative semantics existing between terms in the thesaurus.

Table 2 Examples of the use of RT to represent different semantic relationships

| Relationship | Examples | **Remark** ( More Appropriate Relationship) |
|---|---|---|
| **RT** | 1. *Mutton* RT *Sheep* | *Mutton* <madeFrom> *Sheep* |
| | 2. *Rice* RT *Rice flour* | *Rice* <usedToMake> *Rice flour* |
| | 3. *FAO* RT *UN* | *FAO* <memberOf> *UN* |

## 3 A Hybrid Approach to the Process of Cleaning and Refining Term Relationships

As mentioned in Section 2, AGROVOC data should be cleaned before using it for ontology construction. We have divided the process of data cleaning and semantic relationship refinement into three main steps: Refinement Rule Acquisition, Detection and Suggestion, and, finally, Verification. The system overview is shown in Fig. 1.
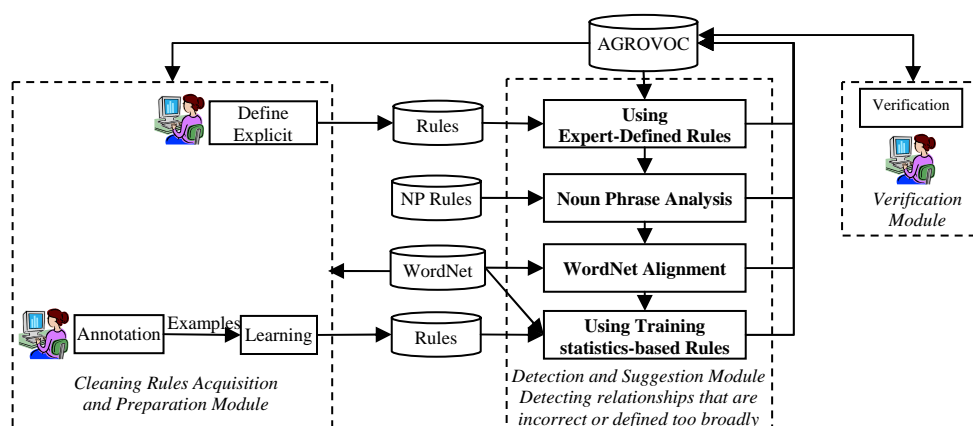


Fig. 1 The System Overview

# 4 The Rule Acquisition Module: Expert-defined Rules and Learning by Example

This module acquires the semantic relationship rules that are used to suggest the appropriate relationships when the AGROVOC relationship is underspecified (defined too broadly), especially RT. The rules will be provided by experts and by machine learning.

## 4.1 Expert-defined rules

As shown in Fig.2, some relationships between terms are presented in AGROVOC completely and consistently, if not as precisely as required for an ontology. In this case, the experts can simply define the rules for systematically revising the inappropriate relationships to the new one. The expert will observe the AGROVOC data and define rules using data on concept types given in AGROVOC as shown in Fig.2. For example, the rules constraint consists of the data in 'concept type data', the category of term such as GC (Geographic term: Country level), GG (Geographic term: above country level), TA (Taxonomic term: Animal), TP (Taxonomic term: Plant).

Based on the given rules, the relationship that satisfies the rule will be revised automatically. For example, consider the following rule:

If X and Y are marked as "**T\***" in the concept type field, and X **BT** Y then X<*subclassOf*> Y

From AGROVOC data, the concept types of *Rosaceae* and *Malus* are TP and they are related by **BT**. Then, the original relationship BT of "Malus BT Rosaceae" will be replaced by <*subclassOf*>.



Fig. 2 Examples of term relationships in AGROVOC that could be handled by revision rules formulated by experts

## 4.2 Learning-from-Examples Rules

Many terms in AGROVOC Database do not have enough information for defining the rule. Moreover, some relationships, especially the relationship named **RT**, could be refined more precisely, as shown in Fig.3. In this case, the rules are prepared by learning from examples.

To prepare the learning examples, we provide an annotation tool that allows the domain expert to manually tag term senses (labelled by a sense id number in WordNet) and to specify the appropriate semantic relationship between them. For example, *(Mutton#1* <madeFrom> *Sheep#1)*.

In the case of compound nouns, only the noun heads are used. For example: *Rice* and *Rice Flour* will be annotated as follows: (*Rice#1* <usedToMake> *Flour#1)*

After preparing the examples, the complete hypernym path of each term will be extracted from WordNet as in the following examples:

{*sheep#1*, *bovid#1, ruminant#1, mammal#1,vertebreate#1, animal#1, organism#1, living_thing#1, object#1,entity#1*}

{*mutton#1*, *meat#1, food#2, solid#1, substance#1, entity#1*}

Fig. 3 Some examples of appropriate relationships for learning the revision rules by examples

The hypernym list, given above, will be used as the basis of the features vector, i.e.

features_vector{{list of hypernym class of term1},{ list of hypernym class of term2}}

The features will be converted into binary representation for obtaining vectors of equal length. The learning system, C4.5, will be applied to learn the common ancestral concept for term1 and term2, and then generate the rules. Fig. 4 shows the example of the data set for training the <madeFrom> relationship. Table 3 shows the revision rules learnt from the training examples.



Fig. 4 Examples of hierarchical data used for training the 'usedToMake' relationship

Table 3 Examples of training statistical-based rule.

| | Rule | Example |
|---|---|---|
| 1 | If class X is *animal#1* and class Y is *meat#1*, and X RT Y Then X <UsedToMake> Y | *Sheep* RT *Mutton*, *Swine* RT *Pork*, *Calf* RT *Veal* |
| 2 | If class X is *plant#2* and class Y is *food#1*, and X RT Y Then X <usedToMake> Y | *Rice* RT *Rice flour*, *Oat* RT *Oatmeal* , *Sugar Cane* RT *Cane Sugar* |
| 3 | If class X is *fruit#1* and class Y is *oil#3*, and X RT Y Then X <usedToMake> Y | *Castor beans* RT *Castor oil*, *Cottonseed* RT *Cottonseed oil* |

By applying the Rule 1, the original relationship RT of "Chicken *RT* Chicken meat " will be replaced by <usedToMake>.

# 5 The Detection and Suggestion Module: An Algorithm for Term Relationship Revision

## 5.1 Overview of the algorithm

In this module, the system detects incorrect and inconsistently applied relationships and suggests the appropriate relationships for expert confirmation. We propose three techniques to handle this process: semantic relationship rules, noun phrase analysis, and WordNet alignment.

The outline of this algorithm is illustrated in Fig. 5, where $T_1$, $T_2$ and Rel denote, respectively, Term1, Term2, and the AGROVOC relationship between them.

```
AGROVOC Cleaning_& Refinement (T₁, T₂, Rel)                ;Return new__relationship
Input: Term1, Term2, Relationship
Output: New Relationship
1. If (Rel = BT or Rel = NT)
   Then If Agree_Expert_defined_Rules (T₁, T₂, Rel)
        Then return new_refined_relationship.              ; following the rules
        Else If Headword-Is-Compatible (T₁, T₂)
            Then return subclass/superclass relationship.
            Else If Is_Wordnet_HypernymPath (T₁,T₂)
                Then return subclass/superclass relationship.
                Else If Agree_Revision_Rules (T₁, T₂, Rel)
                    Then return new_relationship            ; following the rules
                    Else return U.                         ; Un-refined
2. Else If (Rel=UF or Rel = USE)
        Then If Is_Wordnet_Synset (T₁, T₂)
            Then return synonym relationship.
            Else If Agree_Revision_Rules (T₁, T₂, Rel)
                Then return new_relationship.               ; following the rules
                Else return U.                             ; Un-refined
3.      Else If (Rel=RT)
            Then If Agree_Revision_Rules (T₁, T₂, Rel)
                Then return new_relationship.               ; following the rules
                Else return U.                             ; Un-refined
```
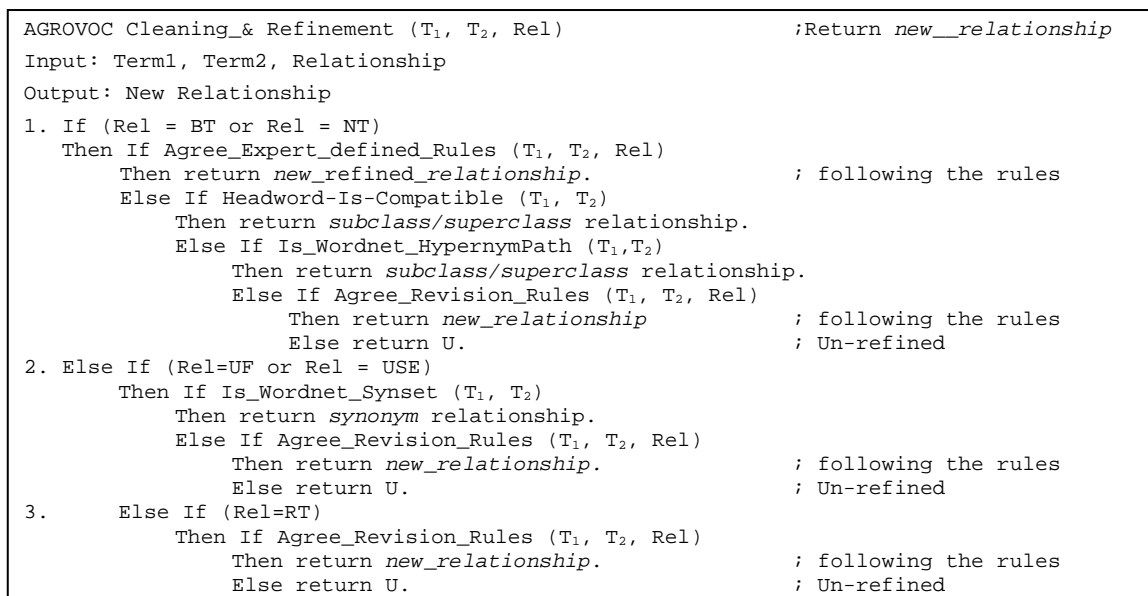
Fig. 5 An Algorithm for Data Cleaning and Relationship Refinement

The relationship revision rules have been discussed in Section 4. Section 5.2 briefly describes the procedures based on noun phrase analysis and WordNet alignment, and Section 5.3 describes the verification tool.

## 5.2 Noun phrase analysis and WordNet alignment

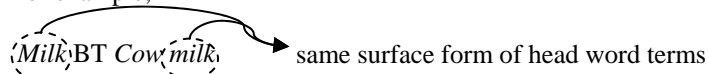- *Using noun phrase analysis*

The noun phrase analysis technique is used to analyze the surface form of a compound term's head word. If the head word of a term has the same surface form as its broader term, the system will apply the 'subclassOf'/ 'superclassOf' relationship to them. The system analyzes compound nouns using the following rule:
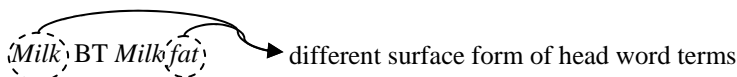
NP ->  MOD NCN

MOD ->  NCN, NPN, ADJ, …

Where MOD is a modifier, NCN is Common Noun, NPN is a proper name, ADJ is an adjective

For example,

*Milk* BT *Cow milk*       same surface form of head word terms

From the compound noun analysis, the head word of *Cow milk* is *milk* which has the same surface form as *Milk*, the broader term of *Cow milk*. Then, the system will apply the <subclassOf> relationship to *Cow milk* and *Milk*.

*Milk* BT *Milk fat* ➝ different surface form of head word terms

The result of the analysis shows that the head word of *Milk fat* is *fat*, which is not compatible with the broader term, *Milk*. In this case, other techniques must be used to refine the relationship.

- *Using WordNet Relationships*

In this step, the hypernym/hyponym relationships of WordNet are used to align the BT/NT relationship in AGROVOC, and the synset of a term in WordNet is used to align the UF/USE relationship in AGROVOC. Since the relationships in WordNet are verified by experts and WordNet contains a great number of general domain terms including agricultural terms, WordNet is a good resource for aligning some AGROVOC relationships such as taxonomic and synonym relationships. (Other verified sources could be used as available, individually or in combination.) The process of this step starts with the system retrieving the synset offset number of the AGROVOC UF/USE term in WordNet. If the system finds these terms and they have the same synset id number, the system will apply the 'synonym' relationship to them. The system will also query the AGROVOC broader term and narrower term in WordNet. If the system finds that the broader term is the ancestor of the narrower term in the WordNet hierarchy, the system will apply the 'subclassOf'/'superclassOf' relationship to them. For example,

*Cabbage* BT *Vegetable*

Query results for *Cabbage* and *Vegetable* in WordNet show that *Cabbage* is a hyponym of *Cruciferous vegetable* and *Cruciferous vegetable* is a hyponym of *Vegetable*. Fig. 6 shows the relationship of *Vegetable* and *Cabbage* in WordNet and AGROVOC.

Since *Vegetable* is an ancestor of *Cabbage*, the system will define *Vegetable* as superclassOf *Cabbage*. In the case of *Milk* NT *Milk fat*, the relationship is not refined by this technique because *Milk* and *fat* are in different hypernym paths in WordNet.
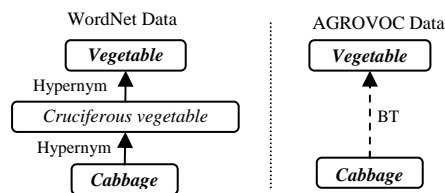


Fig.6 The relationship between Vegetable and Cabbage in WordNet and AGROVOC

### 5.3 The Verification Tool

After the system has suggested the new relationships for terms, the expert will verify the semantic relationship refinement results and also define the appropriate relationship for the cases that cannot be handled by the system. Fig.7 is the user interface for verifying the output of the system. The expert can verify the terms and relationships by querying by

1. Term to verify each term and its relationships to other terms e.g., rice

2. Semantic relationship e.g., <containsSubstance>

3. Rule e.g., 'If class X is meat#1 and class Y is animal#1, and X RT Y then X <madeFrom> Y'.

Fig.7 Verification Tools

## 6 Experimental Results and Future work

We ran an experiment testing the training rules technique using 100 examples for 5 semantic relationships. It produced around 10 classification rules. The experimental results using these rules as well as expert-defined rules, noun phrase analysis, and WordNet Alignment are shown in Table 4.

Table 4 The experimental results classified by relationship

| Relation-ship | No. | No. of refinement | Expert-defined rules | | NP Analysis | | WordNet Alignment | | Training Rules | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | No. | PC(%) | No. | PC(%) | No. | PC(%) | No. | PC(%) |
| BT/NT | 32176 | 21072 | 16587 | 100% | 2062 | 95% | 2423 | 95% | ** | ** |
| USE/UF | 21605 | 3553 | - | - | - | - | 3553 | 70% | ** | ** |
| RT | 27589 | 1420 | 622* | 100%* | - | - | - | - | 798* | 72%* |
| Total | 81370 | 26045 | 17209 | 100% | 2062 | 95% | 5976 | 80% | 798* | 72%* |

Remarks: - indicates this technique can not revise this relationship, * indicates the experiment is run with some data, ** indicates the experiment is in initial state.

Based on an expert's review of a small sample of data, some initial rough estimates were made regarding the precision of the methods. The precision of the Expert-defined Rules technique was estimated to be around 100% and 95% correctness for NP Analysis. The WordNet Alignment technique was estimated to be lower, about 94% precision, because some synonym relationships in WordNet should be replaced with the 'abbreviation_of' relationship. For example, in *AMP* <synonym> *Adenosine monophosphate*, <abbreviaton_of> should be used. The precision of the Training Rules technique was estimated to be about 72%. Sources of error include ambiguity in concept classes used as arguments for a given rule, such as the following, 'If class X is *food#1* and class Y is *food#1*, and X RT Y, then X <usedToMake> Y' where, because X and Y belong to the same concept class, the system cannot distinguish between X and Y and may generate erroneous relationships, e.g., *pork* <usedToMake> *hams*, and *hams* <usedToMake> *pork*. These cases can be revised only by the expert.

There are remain around 55325 unrevised relationships and we will revise only half of them because the inverse relationships will be automatically set. We plan to finish revision in one year with four experts.

## 7 Conclusion

This paper presents the three methodologies for data cleaning and semantic relationship refinement to solve the problem of producing well-defined semantics from poorly defined or underspecified semantics in a thesaurus. The system refines the semantic relationships though noun phrase analysis, WordNet alignment, and semantic relationship rules, some generated by experts and others generated from annotated examples by an inductive statistical machine learning system. Finally, the relationships were verified by the expert. Initial results are promising.

Ontologies with precise semantic are important for improving retrieval systems, for automating processes through machine reasoning, and for the Semantic Web. Developing ontologies is labor-intensive and time-consuming. This paper contributes to solving the ontology development bottleneck by exploiting the enormous intellectual capital amassed over many years in classification schemes and thesauri.

## Acknowledgements

## References

Clark P., et al. 2000. Exploiting a Thesaurus based Semantic Net for Knowledge-based Search", In Proceedings of IAAI-2000

Fisseha, F. and A. C. Liang 2003. Reengineering AGROVOC to Ontologies: Steps towards better semantic structure. NKOS Workshop, 31 May 2003. Rice University, Houston, Texas.

Hearst, M. 1992. Automatic acquisition of hyponyms from large text corpora, In Proceedings of the 14th International Conference on Computational Linguistics.

FAO. 1999. AGROVOC: Multilingual agricultural thesaurus. Rome: Food and Agricultural Organization.

Jannink, J., 1999. Thesaurus Entry Extraction from an On-line Dictionary, In Proceedings of Fusion '99

Kang S. J. and J. H. Lee. 2001. Semi-Automatic Practical Ontology Construction by Using a Thesaurus, Computational Dictionaries, and Large Corpora , ACL 2001Workshop on Human Language Technology and Knowledge Management.

Kietz J. U., A. Maedche, and R. Volz., 2000. A method for semi-automatic ontology acquisition from a corporate intranet, In Proceedings of Workshop Ontologies and Text, EKAW'2000

Maedche B. and S.Staab., 2000. Discovering conceptual relationships from text, In Proceedings of ECAI-2000

Navigli R., Velardi P., Gangemi A. Ontology Learning and its application to automated terminology translation. *IEEE Intelligent Systems*, vol. 18:1, January/February 2003.

Navigli, R., et al., 2003. Ontology Learning and its application to automated terminology translation. IEEE Intelligent Systems, vol. 18, n.1, January February 2003

Soergel, D., B. Lauser, A. Liang, and F. Fisseha. 2004. Reengineering thesauri for new applications. The AGROVOC example. Journal of Digital Information, Volume 4 Issue 4, Article No. 257, 2004-03-17. http://jodi.ecs.soton.ac.uk/Articles/v04/i04/Soergel.

Wielinga, B.J.; Schreiber, A. Th.; Wielemaker, J. ; Sandberg, J. A. C. From Thesaurus to Ontology. En Proceedings of the International Conference on Knowledge Capture, 2001. ACM Press, p. 194-201.