

An Analysis of Knowledge Collected from Volunteer Contributors

Timothy Chklovski, Yolanda Gil

Information Sciences Institute
University of Southern California
4676 Admiralty Way, Marina del Rey, CA 90292, USA
{timc, gil}@isi.edu

Abstract

A new generation of intelligent applications can be enabled by broad-coverage repositories of knowledge. One emerging approach to constructing such repositories is proactive knowledge collection from large numbers of volunteer contributors. In this paper, we study the coverage and quality of a representative collection of part-of information contributed by volunteers. We analyze growth of coverage over time, redundancy of the collected knowledge, and the effect of the coverage and redundancy on the quality of the collection. We also present initial comparisons with collections created by ontology engineering and text extraction approaches. Our analysis reveals that redundancy of contribution helps identify high quality statements, but that some of the statements also have overly high redundancy, drawing contributor effort away from areas where they are needed more. We suggest possible ways to address these issues in future collection efforts.

Introduction

Broad-coverage knowledge repositories stand to enable a new generation of intelligent applications and natural language understanding systems (Chklovski, 2003; Lenat, 1995). The variety of tasks and applications which can benefit from broad-coverage semantic resources are exemplified by uses of WordNet (Miller, 1990), a broad-coverage semantic resource which emphasizes lexical semantics. The WordNet bibliography (Mihalcea, 2004) illustrates hundreds of uses in research.

One approach to constructing broad-coverage semantic (and lexical) resources is by employing a relatively small team of highly trained ontology and knowledge engineers. This approach has been taken by WordNet, CYC (Lenat, 1995), and DOLCE (Gangemi et al., 2003). This approach faces issues stemming from shortage of person-hours available, which can limit the coverage of facts and even limit which semantic relations are included (Lenat, 1995; Miller, 1990). This shortage can also lead to encoding viewpoints or statements that may require later reengineering or refinement (Gangemi et al., 2003; Friedland et al, 2004).

Another approach to constructing broad-coverage resources is text mining (Hearst, 1992; Berland and Charniak, 1999; Riloff & Jones 1999; Schubert, 2002; Girju, Badulescu, & Moldovan 2002; Etzioni et al., 2004). Through sophisticated statistical analysis and training algorithms, these approaches extract entities and discover useful lexical and semantic relations. While the level of precision and recall varies, the extraction of semantic relations remains a challenging area of research.

An emerging approach that we are exploring is to collect knowledge from a multitude of minimally instructed volunteers. The approach can be traced back to at least 1857, when many volunteers aided the construction of the Oxford English Dictionary by mailing in knowledge about earliest known word usages. The recent advent of the Web has greatly simplified distributed contribution of knowledge, attracting a growing amount of research, including *Open Mind Common Sense (OMCS)*, (Singh et al. 2002), *LEARNER* (Chklovski 2003a, 2003b), *LEARNER2* (Chklovski, 2005), the *Fact Entry Tool (FET)* for CYC (Belasco et al, 2002), *Open Mind Word Expert (OMWE)* (Mihalcea & Chklovski 2004), and *Open Mind Indoor Common Sense (OMICS)* (Gupta & Kochenderfer, 2004). Handling quality variation in mass collaboration settings is also being looked at (Lam & Stork 2003) and (Richardson & Domingos 2003). A key benefit of the mass collaboration approach is its inherent ability to bring orders of magnitude more effort to the construction process, since the approach can tap volunteers with minimal or no training. These volunteers also can be prompted with extensively conditioned questions, providing answers which may be challenging to automatically extract from bare text. Also, because different contributors may have different backgrounds and contexts, the collection gathered from them is likely to include statements which are rare. Practical uses of broad-coverage knowledge collections collected from volunteers are also being developed (Lieberman et al, 2004; Gupta and Kochenderfer, 2004).

This paper presents an analysis of the statements collected with one such system in terms of its coverage and quality. This analysis was done on a representative corpus, specifically statements about parts of everyday objects collected by the *LEARNER2* system (Chklovski, 2005). Our analysis shows that if statements are spontaneously contributed, achieving broad coverage is unlikely since

coverage grows ineffectively over time and over number of contributors. Our analysis also shows that some of the collected statements should clearly be discarded, and that such statements can be detected when several humans agree on discarding a statement. We also observe from the analyzed data that for many contributed statements, there is disagreement regarding their acceptability, and that the disagreement can have a variety of sources. Our analysis suggests that there is a role for volunteer contributors in evaluating and qualifying knowledge contributed by others.

The next section introduces LEARNER2 and the knowledge collection studied. After that, we motivate our analysis with examples of deficiencies in coverage and quality from the collection analyzed. Next, we present our analysis of the coverage and the quality of the collected statements. We close with a discussion of aspects of the collection process responsible for the identified strengths and weaknesses in coverage and quality and propose how the collection process can be improved to address these challenges.

Knowledge Collected by LEARNER2

LEARNER2 (Chklovski, 2005) has been deployed for six months as an interactive kiosk at a science museum as part of a traveling exhibit called “Robots and Us¹,” which will continue for 3 more years. LEARNER2 has collected more than 100,000 raw entries from museum visitors of all ages, collecting *meronymy* (*part-of*), *typical purpose*, *similarity*, and other semantic relations about everyday objects. LEARNER2 uses a template-based, fill-in-the-blank approach. This approach focuses the collection effort on specific types of knowledge, which is an extension introduced over LEARNER2’s predecessor, LEARNER. For example, to learn about parts of a “*car*,” LEARNER2 partially instantiates a template to form a fill-in-the-blank knowledge acquisition question:

“a *car* has a piece or a part called a(n) _____”

To exclude malformed entries, the collected knowledge is automatically postprocessed, removing all entries not found in a large lexicon (which removed approximately 25% of the 100,000 raw entries). Spelling mistakes are also discarded to avoid introducing errors by automatically correcting them. The postprocessed knowledge is available as the Learner2-v1.1 dataset². To simplify evaluation, we focus on the meronymy statements (there were a total of 24,747 such statements). LEARNER2 used a seed set of 326 objects (selected from WordNet’s tree of “instrumentation or device”). Users were allowed to introduce other objects as well. The seed objects were semi-automatically selected to exclude very rare objects; the resulting set contains objects such as *axe*, *briefcase*, and *compass*. Since the collection effort focused mainly

on the seed set of 326 objects, we restrict our analysis to them. The resulting set analyzed in this paper contains a total of 6,658 entries, specifying 2,088 distinct statements.

Phenomena Identified in the Collected Statements

In this section, we introduce and motivate the issues present in the data: the ineffective coverage, the presence in the collected knowledge of statements which would need to be identified and discarded, and the presence of statements which are neither clearly acceptable nor clearly discardable but may be one or the other upon further qualification. Later in the paper, we present a quantitative analysis of these issues.

Coverage

Systems that collect knowledge from volunteers typically collect what can be called “spontaneous” contributions, that is statements about whatever topic or object comes to the contributor’s mind. As a result, there can be high redundancy in typical items and also spotty coverage in more unusual ones. This was the case in the collection we analyze here. Some statements are entered dozens of times at the expense of other acceptable statements, which are never entered. To illustrate, the 5 most frequently contributed (0.24% of all distinct) statements attracted a total of 533 (8.0% of all collected) entries:

<i>part-of</i> (handle, hammer)	136	<i>part-of</i> (blade, knife)	99
<i>part-of</i> (wheel, car)	121	<i>part-of</i> (wing, airplane)	75
<i>part-of</i> (engine, car)	102		

At the same time, some useful statements such as *part-of*(radiator, car), *part-of*(crankshaft, car), and *part-of*(aileron, airplane) were never entered.

These observations raised the issue of whether to stop collecting redundant contributions and if so how many times should a statement be collected before the utility of additional identical contributions becomes negligible. Another important issue is whether and how to steer contributors to contribute new statements when the collection contains a sizeable amount of what could be considered the most common or typical statements. In the “coverage” subsection of the next section, we will show a detailed analysis based on data from the LEARNER2 corpus regarding these issues.

Categories of acceptability of the collected knowledge: the good, the bad, and the needing qualification

Another important set of phenomena that we observed in the LEARNER2 data is a wide variety of quality or acceptability of the knowledge. There are statements arising from contributors occasionally disregarding the collection instructions, such as *part-of*(chicken, knife) and

¹ <http://www.smm.org/robots/>

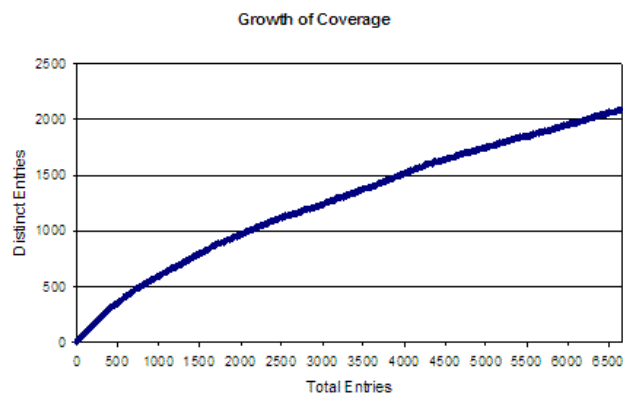
² The live system and the collected data are available at <http://learner.isi.edu>

part-of(truck, pot), that should clearly be discarded. We noticed that these are a very small portion of the collection.

Judging the quality of the collection in terms of its correctness or accuracy is a non-trivial task. This is the case with many kinds of knowledge and is not specific to *part-of* relations. Whether a given statement is indeed a *part-of* statement involves a number of subtleties. For example, (Winston et al, 1987) have discussed the types of the *part-of* relation, such as component/integral object, member/group, place/area, and others while Miller (1990) highlights instances of non-transitivity of the relation. The issues we observed had more to do with how the notion of the *part-of* relation and the terms in the relation need to be qualified to determine whether a given statement is acceptable.

Given the lack of a formal or intensional definition of correct *part-of* relations, we decided not to treat correctness as an all-or-nothing matter but rather as something that can be increased by additional context to the statement. For example, *part-of*(film, camera) was entered by several contributors and is not clearly wrong. Yet, the statement does not hold for digital cameras, or newly purchased, not yet loaded cameras, and so on. What should be counted as an object and therefore as its parts is also not always clear cut. For example, acceptability of *part-of*(elevator shaft, elevator), and *part-of*(sail, mast) depends on whether the elevator refers to just the elevator cab or to the whole elevator structure, and whether the mast refers to the structure with the sail and the rigging or just the bare structure. Other statements are questionable because the part was not tangible, as in *part-of*(hole, tube), *part-of*(flame, torch). Word senses can also play a role. For example, *part-of*(row, table); *part-of*(mouse, computer) drew disagreement in evaluation scores. Although collecting explicit information on senses in which words are used would be useful, such collection involves an entire set of research issues (e.g., Mihalcea and Chklovski, 2004) which have not been engaged by LEARNER2.

Given that our ultimate goal is to collect common knowledge about everyday objects, we would prefer to keep all of these statements in some form within the collection. This is a very challenging issue, and one that we discuss below in more detail. It is worth noting that such statements are often not included in manually engineered and highly curated resources such as WordNet. In construction of knowledge repositories by knowledge engineers, the knowledge encoded is typically prescriptive. That is, if a statement is often, but not necessarily true, it would likely not be included. For example, WordNet specifies that a dog is a mammal, but does not provide any indication that dogs are (often) pets. By contrast, the statements we collect tend to include statements which are only sometimes true, such as *part-of*(remote control, stereo) and *part-of*(rope, pulley). Harnessing the ability to collect such statements and perhaps qualifications of the context in which they hold may be a potential strength of the approach of collecting from volunteers.



Detailed Analysis of the Collected Statements

In this section, we analyze in detail how contributor statements are distributed and the impact of this distribution on coverage. We also suggest possible indicators of acceptability of knowledge and analyze their merits based on the data collected.

Coverage

Out of a total of 6,658 entries collected, only 2,088 are distinct; 68.6% of entries were spent on getting redundant knowledge, adding nothing to coverage. Furthermore, examining all entries contributed three or more times reveals that 4,416 entries (66.3% of all entries) yielded only 350 distinct entries (16.8% of all distinct entries). This suggests that contributor effort was inefficiently exploited and could be redirected to areas that have poorer coverage.

Furthermore, as shown in Figure 1, as the collection grows, the ratio of distinct to all statements contributed so far keeps decreasing. The diminishing returns seem to come from two sources. The first source is simple saturation of distinct answers. As the more frequent answers are collected the new ones become increasingly rare. The second source stems from the variability in the number of acceptable answers to a question. For example, even though in the collection studied all parts of a *hammer* and an *axe* have probably been collected, many parts of a *watch* have not yet been. Yet, the system currently keeps querying about objects without any preference for those about which knowledge is less complete. Hence, coverage suffers from contributor effort not being directed both at the question and at the answer level.

Towards classifying knowledge by acceptability

Given the considerations discussed above on how to judge quality and acceptability and lack of a working definition, we turned to evaluation by majority vote of human judges, a methodology previously selected by Berland and Charniak (1999) and Girju (2003) to evaluate automatic

text extraction techniques. While an imperfect indicator of acceptability, as has been pointed out by Berland and Charniak (1999), majority vote provides a practical way to assess it. In our analysis, we asked 3 subjects (judges) to rate collected statements on a scale ranging from 0 to 3 (“is not,” “probably is not,” “probably is,” and “is” a part-of relation). Statements were presented in a random order.

We consider two potential indicators of the acceptability of statements: redundancy and generation frequency.

We first examine whether the number of times a statement has been entered (its *redundancy*) is indicative of the opinion of the judges. To that end, we sampled the data from several redundancy categories: 1, 2, 3, “4 or more”. These categories were not shown to the judges. Our sample had 869 items in all: 250 items of redundancy 1, 250 of redundancy 2, 119 (all available) from category 3, and another 250 items of redundancies 4 or more. In the presented analysis, the sampled data from the “4 or more” category was additionally broken out into the categories of “exactly 4” and “5 or more”.

Table 2. Redundant contributions and majority vote

# times statement was contributed (redundancy)	# distinct statements in this category	% for which majority voted “is” or “is probably” part-of
1 or more (all statements)	2,088	70.5%
2 or more	735	89.8%
3 or more	469	93.8%
4 or more	350	95.9%
5 or more	271	97.5%
exactly 1	1,353	60.0%
exactly 2	266	82.8%
exactly 3	119	87.4%
exactly 4	79	90.6%

Table 2 presents the results. The number of times statements were contributed is shown as well as the proportion of statements rated as “is” or “is probably” part-of by the majority of the three judges. The bottom of the table shows the number of statements contributed a given number of times. When our sampling is weighted by true number of statements in each sampled subset of statements, 70.5% of all statements receive the majority vote of judges. Of statements contributed more than once, majority vote is received by 89.8%. The majority vote increases monotonically with the number of times a statement has been contributed, with 97.5% of statements with contributed five or more times receiving the majority vote. In our evaluation sample, all 52 statements with contribution frequency of 15 or more were accepted (although in 3 cases one judge dissented). All three judges accepted all 35 evaluated statements that were entered 23 times or more (the maximum times a statement was entered is 136).

The entries that we suggested earlier as ones that should clearly be discarded, such as *part-of*(chicken, knife) and *part-of*(truck, pot), primarily had the redundancy of 1. In the set evaluated by judges, there were 25 such statements. Also, all such statements received the lowest evaluation score from at least two of three judges, giving promise to future work on their identification.

A second potential indicator that we used in our analysis is *generation frequency*. It is based on the notion that more common or typical statements that are spontaneously brought up by many users are more likely to be acceptable. We define the generation frequency (*gf*) of a statement about a given part and an object as the frequency with which this part has been contributed out of a total number of times a statement has been made about any part of this object. For example, *part-of*(handle, hammer) was contributed 136 times out of a total of 203 statements about parts of a hammer. This yields the generation frequency of *part-of*(handle, hammer) to be $136/203=0.67$. We expected answers with higher generation frequencies to be rated more highly.

Table 3 shows the results. We show separately the results for statements contributed once, twice and so on into two sets: those with *gf* below 0.1 and those with *gf* of at least 0.1 (splitting the evaluation data into two sets of roughly equal size). Surprisingly, for redundancy greater than 1, items with lower generation frequency tend to be more acceptable to judges than items with the higher generation frequency. This finding suggests that collecting low-frequency items may not negatively impact the quality of the collection.

Table 3. Generation frequency and majority vote

# times statement was contributed (redundancy)	Gen freq < 0.1		Gen freq ≥ 0.1	
	% receiving majority vote “is” or “is probably”	Num in sample	% receiving majority vote “is” or “is probably”	Num in sample
exactly 1	55.7%	158	67.4%	92
exactly 2	89.8%	127	75.6%	123
exactly 3	92.3%	52	83.6%	67
exactly 4	92.0%	25	89.3%	28
5 or more	100.0%	44	96.7%	153

To sum up, the number of times that a statement has been contributed is a strong indicator of majority vote and therefore acceptability to judges. However, high generation frequency, for statements contributed more than once, is not. The mixed assessment of statements contributed once suggests that more information is needed about the acceptability of these statements. The positive assessment of the statements contributed many times suggests that they require relatively little further assessment effort.

Human Agreement when Evaluating Statements

In our study, we have relied on three evaluators to rate acceptability of the collected knowledge. The evaluators rated the 869 items at a rate of around 10 items per minute, which suggests that a large volume of evaluations can be carried out relatively quickly.

At the same time, we observed some signs for the need of calibration of the evaluators. While the most permissive evaluator in our study rated 85.5% as “is” or “is probably” part-of, the least permissive one assigned one of those ratings to only 65.8% of the same statements. The overall inter-annotator agreement of the judges who received little instruction was 76.6%, while agreement on answers with redundancy 4 or more was 85.1%. This suggests that there may be significant individual differences between evaluators’ assessments, presenting a challenge to future validation efforts.

Analysis of Comparable Resources

An area that requires further work is the detailed comparisons of the content of our collections versus resources created through other approaches such as ontology engineering and text extraction. Here, we present some initial results which indicate that the approaches may be complementary and amenable to combination.

Extracting the part-of relation from text has been attempted by Berland and Charniak (1999), reporting 55% accuracy and citing issues such as lack of unequivocal syntactic indicators for the part-of relation in text. Girju et al., (2003), resorted to ontological knowledge and a large amount of manually annotated training data to improve extraction precision, reporting precision of 83% on an extracted set of 119 statements. For statements contributed 2 or more times, our accuracy is 89.8%, which surpasses the results from text extraction. Still, automatic extraction from very large corpora (e.g., Hearst, 1992; Etzioni et al, 2004; Riloff & Jones 1999; Schubert, 2002) may uncover valuable statements to augment or seed volunteer collection efforts.

WordNet also contains part-of relations, although an appropriate comparison is difficult to formulate because our collection does not differentiate among word senses, while WordNet often resorts to very rare senses. For example, one of the senses of a “pen” in WordNet is a “female swan,” which, as a “whole object” has a part “part” and, as a “bird” has parts such as “oyster” (a small muscle of a bird). Such part-of statements would be highly unlikely to be contributed by (or acceptable to) volunteers, because these statements are not what we may call “common knowledge.” Another example of confusion of the common and the rare cases in WordNet is that the senses of “cat” and “dog” which are closest to each other are actually those in which “cat” and “dog” are a kind of person. In an initial effort to compare the contents of both resources, we looked at the part-of relations explicitly stated in WordNet (not inherited, not composite, not

derived) of primary senses of the concepts studied. We found an overlap of statements LEARNER2 collected with WordNet to be 10-15%. We plan to conduct a more thorough analysis in the future. One important issue to consider is how volunteer contributors could complement WordNet with indications of “common” knowledge and typical or default word usage.

Conclusions and Future Work

In summary, these are the main findings supported by the representative collection we analyzed and the implications of these conclusions for future work:

- *Redundancy of a statement, achieved by repeated collection from multiple volunteers, is a useful indication of the quality of the statement.* In future work, it would be a good idea to continue to leverage repeated collection of statements to obtain better quality. Note that some collection systems, such as LEARNER, do not do this.
- *If little guidance is provided to the contributors some statements are collected with very high redundancy, which consumes contributor effort that may have higher utility in other areas.* This inefficiency suggests that limiting redundancy in repeated contributions would be useful. Note that the current systems collecting statements from volunteers, including OMCS, OMICS, LEARNER or LEARNER2, currently do not carry out such management of redundancy. In the next point, we discuss guidance and redirection of contributor effort in more detail.
- *As the collection grows through spontaneous contributions, there are diminishing returns in the coverage achieved.* This suggests that redirection of contributor effort in the collection process towards providing not yet collected or not sufficiently redundant knowledge would be beneficial. We briefly outline some methods which could be used to redirect contributor effort to broaden coverage: (a) *guide contributors away from known answers*, keeping a “taboo list” made up of the top most frequent answers to a question; (b) *collect knowledge about insufficiently covered objects*, using some saturation criterion to guide contributors towards objects about which new answers continue to be contributed (for example, a hammer has fewer parts than a car, and the collection process should reflect this); (c) *allow entry of several answers per question*. This allows the contributor to engage with the question more deeply, recalling additional, less salient knowledge; (d) *prompt contributors with possible answers*. This method uses collected knowledge to suggest other possible answers or to generate similar answers (as was done in LEARNER but not in LEARNER2 or other systems); (e) a method synergistic with text extraction approaches is also possible, where knowledge collected from volunteers can be used to aid extraction of text fragments from text corpora (e.g. the Web or an

encyclopedia) and such text fragments may also be used to prompt contributors.

- *Typical statements generated by many users (those with high generation frequency) may not necessarily be of better quality.* This suggests that concentrating the contributions on typical statements may not necessarily augment the quality of the collection.
- *Human agreement on quality of a statement varies across statements. Evaluators also vary in their degree of acceptance of statements.* This suggests the need for careful design of any future validation and refinement (qualification) mechanisms. Evaluators may need to receive more instruction on the criteria to be used in evaluating knowledge. The quality of the evaluators can be tracked by offering for evaluation “gold standard” items which would allow the collection system to “evaluate the evaluators.”

A key factor in determining the criteria and thresholds of quality of the statements may be the intended usage or application of the knowledge. For example, an intelligent user interface application may require the list of the most typical parts, most agreed upon parts, while an application for reference resolution is better served by a more inclusive, even if less reliable, list of statements. As more applications of such collections are developed, more will become known about the criteria and thresholds which are most useful to these applications.

Another important area of future work is augmenting the statements in a collection with additional information also collected from volunteers. For example, the statements could be annotated with indicators of disagreements and flagged for further refinement and qualification. Qualification of the knowledge may take a number of different forms, which are perhaps rooted in deep knowledge representation issues. Qualifying the part-of relation could take the form of further qualifying the relation itself. For example, *part-of*(idea, textbook) could be refined to *intangible-part-of*(idea, textbook). Senses of the terms in the relation could also be explicitly specified, especially in the cases where non-default senses are used or the default sense is unclear, as for “head” in *part-of*(head, beer) or for “table” in *part-of*(row, table). Another type of qualification could be specifying whether the relation holds for all or only for some instances of the objects, as for *part-of*(remote control, radio), *part-of*(airbag, car).

Knowledge collection from volunteer contributors is a novel and promising approach to creating broad coverage knowledge repositories. The design of the collection tools, the collection process, and the nature of the contributions can be greatly improved through the kind of empirical analysis that we have presented in this paper.

Acknowledgments

We gratefully acknowledge funding for this work by DARPA under contract no. NBCHD030010.

References

- Belasco, A., Curtis, J., Kahlert, R., Klein, C., Mayans, C., Reagan, P. 2002. Representing Knowledge Gaps Effectively. In *Practical Aspects of Knowledge Management, (PAKM)*, Vienna, Austria, December 2-3.
- Berland, M. and Charniak, E.. 1999. Finding Parts in Very Large Corpora. In *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics (ACL-99)*.
- Chklovski, T. 2003a. *Using Analogy to Acquire Commonsense Knowledge from Human Contributors*, PhD thesis. MIT Artificial Intelligence Lab technical report AITR-2003-002
- Chklovski, T. 2003b. LEARNER: A System for Acquiring Commonsense Knowledge by Analogy. In *Proceedings of Second International Conference on Knowledge Capture (K-CAP 2003)*.
- Chklovski, T. and Pantel, P. 2004. Path Analysis for Refining Verb Relations. In *Proceedings of KDD Workshop on Link Analysis and Group Detection (LinkKDD-04)*. Seattle, WA.
- Chklovski, T. 2005. Designing Interfaces for Guided Collection of Knowledge about Everyday Objects from Volunteers. In *Proceedings of the 2005 Conference on Intelligent User Interfaces (IUI-05)*. San Diego, CA.
- Etzioni, O., Cafarella, M., Downey, D., et al. 2004. Methods for Domain-Independent Information Extraction from the Web: An Experimental Comparison. In *Proceedings of the Nineteenth National Conference on Artificial Intelligence (AAAI-2004)*. San Jose, CA.
- Friedland, N., Allen, P., Matthews, G., Witbrock, M. et al. 2004. Project Halo: Towards a Digital Aristotle. *AI Magazine*, 25(4): Winter 2004, 29-48
- Girju, R., Badulescu, A., and Moldovan, D. 2003. Learning Semantic Constraints for the Automatic Discovery of Part-Whole Relations. Proc. of the *Human Language Technology Conference (HLT)*, Edmonton, Canada.
- Gangemi, A., Guarino, N., Masolo, C., Oltramari, A. 2003. Sweetening WORDNET with DOLCE. *AI Magazine* 24(3): 13-24.
- Gupta, R., and Kochenderfer, M. 2004. Common sense data acquisition for indoor mobile robots. In *Nineteenth National Conference on Artificial Intelligence (AAAI-04)*.
- Hearst, M. 1992. Automatic acquisition of hyponyms from large text corpora. In *Proceedings of the 14th International Conference on Computational Linguistics*.
- Lam, C. and Stork, D. 2003. Evaluating classifiers by means of test data with noisy labels. *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI-03)*, Acapulco, Mexico, pp. 513–518

- Lenat, D. 1995. CYC: A large-scale investment in knowledge infrastructure. *Communications of the ACM*, 38 (11)
- Lieberman, H., Liu, H., Singh, P., and Barry, B. 2004. Beating common sense into interactive applications. *AI Magazine*, Winter 2004, 25(4):63-76. AAAI Press.
- Mihalcea, R. and Chklovski, T. 2004. Building Sense Tagged Corpora with Volunteer Contributions over the Web, book chapter in *Current Issues in Linguistic Theory: Recent Advances in Natural Language Processing*, Nicolas Nicolov and Ruslan Mitkov (eds), John Benjamins Publishers.
- Mihalcea, R. 2004. WordNet bibliography, available online at <http://enr.smu.edu/~rada/wn>
- Miller, G. 1990. Nouns in WordNet: A Lexical Inheritance System. *International Journal of Lexicography*, 3(4), 245-264
- Richardson, M., Domingos, P. 2003. Building large knowledge bases by mass collaboration. In *International Conference on Knowledge Capture (K-CAP03)*
- Riloff, E. and Jones, R. 1999. Learning Dictionaries for Information Extraction by Multi-Level Bootstrapping. In *Proceedings of the Sixteenth National Conference on Artificial Intelligence (AAAI-99)*, pp. 474-479.
- Schubert, L. 2002. Can we derive general world knowledge from texts? In *Proceedings of the Human Language Technology Conference (HLT 2002)*, San Diego, CA, pp. 94-97
- Singh, P., Lin, T., Mueller, E., Lim, G., Perkins, T., Zhu, W. 2002. Open Mind Common Sense: Knowledge acquisition from the general public. In Meersman, R. and Tari, Z. (Eds.), LNCS: Vol. 2519. *On the Move to Meaningful Internet Systems: DOA/CoopIS/ODBASE* (pp. 1223-1237). Springer-Verlag.
- Winston, M. E., Chaffin, R. and Herrmann, D. 1987. A taxonomy of part-whole relations. *Cognitive Science*. 11:417-444