

Bridging Common Sense Knowledge Bases with Analogy by Graph Similarity

Yen-Ling Kuo and Jane Yung-jen Hsu

Department of Computer Science and Information Engineering
National Taiwan University
yjhsu@csie.ntu.edu.tw

Abstract

Present-day programs are brittle as computers are notoriously lacking in common sense. While significant progress has been made in building large common sense knowledge bases, they are intrinsically incomplete and inconsistent. This paper presents a novel approach to bridging the gaps between multiple knowledge bases, making it possible to answer queries based on knowledge collected from multiple sources without a common ontology. New assertions are found by computing graph similarity with principle component analysis to draw analogies across multiple knowledge bases. Experiments are designed to find new assertions for a Chinese commonsense knowledge base using the OMCS ConceptNet and similarly for WordNet. The assertions are voted by online users to verify that 75.77% / 77.59% for Chinese ConceptNet / WordNet respectively are good, despite the low overlap in coverage among the knowledge bases.

Keywords: commonsense knowledge base, graph-based analogy, social games, knowledge integration

Introduction

Programs are brittle as computers today are notoriously lacking in common sense. While it is natural for humans to cope with unfamiliar situations based on a huge store of background knowledge, computers often fail when given incomplete or inconsistent information. To break the software brittleness bottleneck, AI researchers proposed to equip computers with common sense knowledge about the world, so the general knowledge may take place whenever the domain-specific knowledge fails (Newell and Ernst 1965).

Codifying millions of pieces of knowledge that comprise human knowledge into machine usable forms has proved to be time-consuming and expensive. Started in 1984, the Cyc project aims to build a large knowledge base of real world facts and common sense (Lenat, Prakash, and Shepherd 1986). A team of knowledge engineers carefully crafted knowledge into CycL, a rigorous logic-based language to ensure its correctness (Lenat 1995). After 25 years, the OpenCyc 2.0 ontology contains hundreds of thousands of terms with millions of assertions relating the terms to each

other. In contrast, the Open Mind Common Sense (OMCS) project at MIT (Singh et al. 2002) took the Web 2.0 approach and appeal for contributions from online users. Over a million sentences in multiple languages have been collected and are encoded as semantic networks. While significant progress has been made in building large common sense knowledge bases, they are intrinsically incomplete and inconsistent.

In addition to the challenge of building up knowledge bases, our research tries to bridge the gaps between multiple knowledge bases, making it possible to answer queries based on data collected from multiple sources without a common ontology. This paper presents a novel approach to knowledge integration using analogy and glossary mapping. New assertions are discovered by computing graph similarity with principle component analysis to draw analogies across multiple knowledge bases.

This paper starts by reviewing related work on knowledge base integration. Following an overview of common sense representation, collection, and reasoning in ConceptNet, the proposed analogy-based inference procedure is introduced. We then present the experimental set-up to find new assertions for a Chinese commonsense knowledge base using the OMCS ConceptNet and similarly for WordNet. To evaluate the experimental results, the inferred assertions are voted by online users to verify their precision. The paper concludes with a comparison of the proposed procedure with blending.

Related Work

In this section, we review some approaches to challenges in integrating multiple knowledge bases.

Ontology mapping and merging In order to reuse concepts from multiple knowledge bases, it is often necessary to merge existing ontologies into a single ontology (Pinto and Martins 2001). *Ontology mapping* refers to the process of combining distributed and heterogeneous ontologies based on linguistic or structure similarity. A number of tools have been built to tackle this difficult problem, including PROMPT (Noy and Musen 2000), FCA-Merge (Stumme and Maedche 2001), and CHIMAERA (McGuinness et al. 2000), etc. Multiple knowledge bases can then be integrated

and reused according to the merged ontology.

Blending In (Havasi et al. 2009), *blending* was proposed as a technique to integrate common sense into other systems. In particular, blending of common sense knowledge with domain-specific knowledge can be done by finding an analogical closure across multiple, previously separated sources of data. Two sparse matrices are combined linearly into a single, larger matrix. Reasoning with blended knowledge bases containing overlapping information can produce inferences that would not be produced from either input alone. However, the result of blending may suffer from noises introduced by linear combination when some knowledge sources may be unreliable.

Analogy Instead of merging the knowledge bases explicitly, knowledge base integration can be achieved using analogy. As a fundamental function of human cognition, analogy is essential for implementing intelligent behaviors in AI systems. In structure-mapping theory (Gentner 1983), an analogy is defined as a set of correspondences that align objects in the source domain (i.e. base domain) with objects in the target domain. Implementations of structure-mapping theory include (Krishnamurthy 2009; Sowa and Majumdar 2003; Turney 2008) and the popular structure-mapping engine (SME) (Falkenhainer, Forbus, and Gentner 1989), which uses propositional logic representation and manually constructed LISP input.

For knowledge bases represented as semantic networks, CrossBridge is an efficient algorithm for finding analogies (Krishnamurthy 2009). The graphical representation lessens the effort of knowledge base transformation in the concept mapping process. Dimensionality reduction is used to reduce the computation required in finding analogies. This paper proposes a solution based on CrossBridge, which will be explained in more detail, and its results will be evaluated against blending.

Common Sense Knowledge Base

Common sense reasoning may call for different knowledge representations to satisfy specific problems and requirements (Singh 2002). The two most prominent common sense knowledge bases, Cyc (Lenat 1995) and ConceptNet (Liu and Singh 2004), made different choices of knowledge representation. Cyc chooses a formal logical framework that is appropriate for representing precise and unambiguous facts carefully encoded by knowledge engineers. On the other hand, ConceptNet represents user generated sentences as a large semantic network, which is better suited for contextual reasoning and computing concept similarity.

ConceptNet knowledge representation

In this research, the semantic network representation of ConceptNet is adopted due to its flexibility and ease in knowledge base integration. The OMCS ConceptNet represents all sentences in the corpus as a directed graph (Havasi, Speer, and Alonso 2007). The nodes of this graph are *concepts*, and its labeled edges are *relations* of common sense knowledge

connecting two concepts. There are over 20 relation types defined in ConceptNet. For example,

- UsedFor (*a*, *b*), e.g. [Spoon] is used for [eating].
- IsA (*a*, *b*), e.g. [Dog] is an [animal].

ConceptNet knowledge collection

The current ConceptNet corpora contain over one million statements in English, as well as about a quarter million statements in Chinese and Portuguese, respectively. The English and Portuguese corpora were collected from over 15,000 contributors at the OMCS website within a span of just about 10 years¹. In addition, about 20% of the English sentences were collected via Verbosity, a human computation game (von Ahn, Kedia, and Blum 2006). Unfortunately, user contributions in the other languages never took off. With innovations in community-based social games, human computation games benefit from rich interactions inherent in a community. The up-to-date knowledge in the Chinese ConceptNet was successfully collected and verified via question-answering between players within a year (Kuo et al. 2009). However, the knowledge collected from these sources differs on their depth and diversity, which also makes ConceptNet itself incomplete.

Common sense reasoning in ConceptNet

It is straightforward to equip a variety of applications with common sense by querying the OMCS ConceptNet using APIs. For example, one may ask if a specific assertion is present in the corpus or its frequency.

AnalogySpace (Speer, Havasi, and Lieberman 2008) generalizes the reasoning method called *cumulative analogy* (Chklovski 2003) so that it is robust enough in large and noisy semantic network. The assertions are divided into *concepts* and *features*, i.e. descriptions of concepts such as “UsedFor eating” or “dog IsA”. The knowledge in ConceptNet is represented as a sparse matrix, and its most prominent features can be identified by using singular value decomposition (SVD). Concept similarity is defined in terms of their shared features.

Concept Mapping by Analogy

Notations

Before describing the proposed method for knowledge base integration using concept mapping by analogy, we first define the terminology used in this paper. A *domain* is a knowledge base which is transformed into the semantic network representation with *concepts/relations* as its nodes/edges. The *features* we used to describe a concept are consisted of its neighbors and the relations it involves. The *graph similarity* between graphs refers to the similarity of their structures.

Problem Definition

- Input

¹<http://openmind.media.mit.edu/>

1. Two knowledge bases – the source domain S and the target domain T .
 2. A few mappings of concepts and relations between the two knowledge networks, indicating overlapping information between the two data sets.
 3. A query concept c in the target domain T .
- Output
 1. A list of features found from the source domain S that can be used to answer queries about concept c .

Inference Procedure

1. **Cross domain analogy mapping from T to S** Analogical processing plays an important role in searching for similar concepts across domains under the assumption that similar concepts usually participate in similar relations. Also, analogy reduces the search space from the network to its subset. Using the subgraphs constructed by c and its neighbors in T we can find a set of analogical mappings M in S by any graph isomorphism/similarity algorithm. Each mapping $m \in M$ consists of correspondences of concepts in both domains and a structural evaluation score.
2. **Check glossary mapping constraint** We have observed that every concept similar to c shares some features with c if they are in related domains (i.e. two domains having overlapping information). For every mapping $m \in M$ with structural evaluation score exceeding a specified threshold, it is checked against the glossary mapping constraint. If any of the neighbors of some concept mapping c appears in the input glossary mappings, m is marked as a candidate mapping.
3. **Transfer features from S to T** Any concept verified by analogy and glossary mapping is very likely a corresponding concept of c in S . Once we have the set C' of candidate mappings of c , we can transfer the common features F in subsets of C' to describe c . Each feature $f \in F$ is assigned a *score* which is the size of the set of concepts sharing F . Finally, the new assertion for answering the query of c consists of the concept c and the feature f .

Example

To give a comprehensive understanding of how this knowledge mapping and reusing process works, we illustrate an example of finding new assertions for concept “狗(dog)” in Chinese ConceptNet as shown in Figure 1.

First, we use “狗(dog)” and its neighbors to form a subgraph such as the left part of Figure 1. Then, we find the analogies “run-(CapableOf)-cat-(IsA)-animal” and “eat-(CapableOf)-lion-(IsA)-animal” in S (the right part of Figure 1) by structure mapping of *relation structure* “a-(CapableOf)-c-(IsA)-b” where a and b are neighbors of c . Since run and animal are in glossary mappings, we are confident that cat and lion are similar to “狗(dog)” based on some identical features. Using the common features of cat and lion, we create a new assertion “狗(dog) Has fur” for answering the query of “狗(dog)” in Chinese ConceptNet.

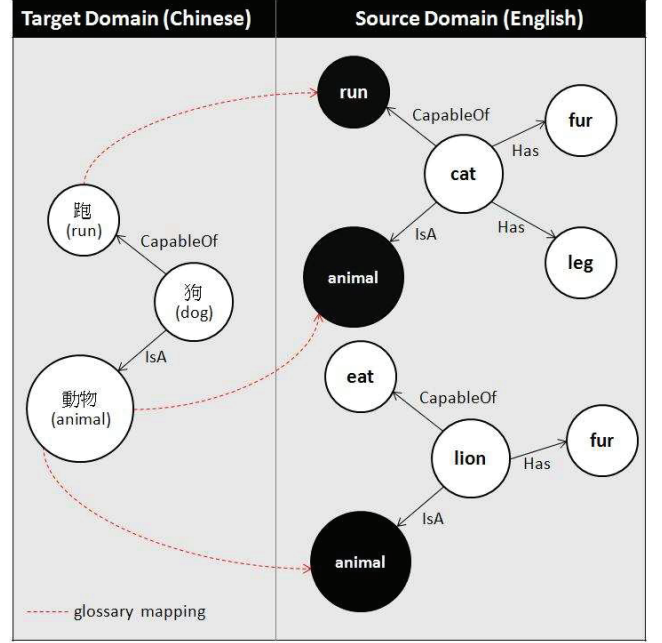


Figure 1: Example of concept mapping between Chinese and English ConceptNet. “狗(dog)” is mapped to “cat” and “lion” by analogy and glossary mapping constraint.

Experiments

We tested the proposed method by experimenting with real knowledge bases. Given that the English ConceptNet is the biggest corpus, we decided to leverage it as a base to find new assertions for the other corpora. It is also observed that the proposed method is less effective when the size of the knowledge base is too small, e.g. for French or Japanese. Therefore, we conducted the experiment to extend the Chinese ConceptNet.

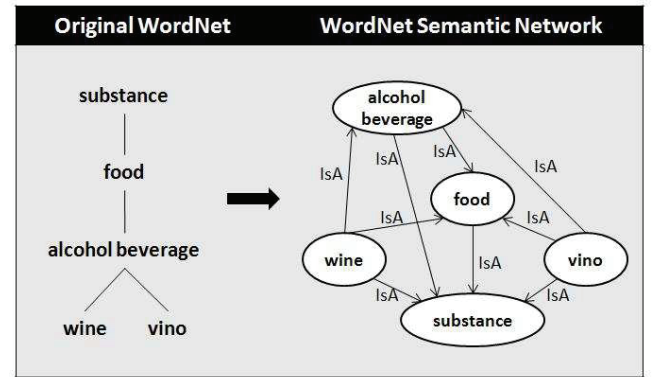


Figure 2: Transform WordNet to semantic network by adding IsA relation to all ancestor of “wine” and “vino.”

Dataset As a comparison, the experiment is conducted on another network built with about 3000 concepts from Word-

Net (Miller 1995). WordNet is a lexical database built by linguists and psychologists where each unique meaning of a word is grouped into *synset*. Synsets are connected to each other through 11 relation types, with hypernym being the most popular link. WordNet maintains a hierarchical structure as Figure 2, which is transformed into a semantic network by adding IsA relation to each concept and all of its hypernym. Table 1 shows the corresponding ConceptNet links for the other four relation types, which are also added to the semantic network.

Table 1: Mapping of WordNet and ConceptNet relation

WordNet relation	ConceptNet relation
Hypernym	IsA
Part Holonym	PartOf
Substance Meronym	MadeOf
Attribute	HasProperty
Entailment	Cause

Since WordNet only contains lexical information of words, we try to extend its semantic meaning with some of the assertions of ConceptNet in our experiment. Table 2 is a summary of the Chinese ConceptNet and WordNet data sets as compared with the English ConceptNet.

Table 2: Statistics of ConceptNet and WordNet

	English ConceptNet	Chinese ConceptNet	WordNet
# of concepts	274,477	64,467	155,287
# of assertions	815,275	205,526	–
# of relation types	20	20	11
% overlap in concepts	–	1.27%	4.79%
% overlap in relations	–	100%	45.45%

Finding new assertion by CrossBridge CrossBridge is adopted as the algorithm for finding analogies among semantic networks. In CrossBridge, the structure of a graph is described by the graph’s *structure vector*, which is a set of relation structures mapped between domains by structure mapping. Instead of searching for isomorphic subgraphs, it searches for subgraphs with similar structure vectors after applying principal component analysis to reduce the dimension of the structure vectors. We modified CrossBridge to handle multiple semantic networks, possibly in different domains, cultures, or languages.

To find analogies between Chinese and English ConceptNet, we treat Chinese ConceptNet as the target domain and English ConceptNet as the source domain. The glossary mapping table is created from the Langdao Chinese-English dictionary. Following the example illustrated in the previ-

ous section, our experiments find new assertions for 8,000 concepts in Chinese ConceptNet.

Similarly, we experimented with WordNet as the target domain and English ConceptNet as the source domain. Given that both knowledge bases are in English, the glossary mapping table is defined as the same words in both networks.

The example of new assertions created by our method are listed in Table 3. In addition to examining the new assertions produced for Chinese ConceptNet and WordNet, we have also observed some interesting properties of our method.

Table 3: Example of new assertions

Chinese ConceptNet	WordNet
打電動(<i>play video games</i>) HasProperty fun	<i>pickle</i> AtLocation grocery store
打電動(<i>play video games</i>) IsA activity	<i>hotdog</i> HasProperty edible
child Desires 打電動(<i>play video games</i>)	<i>piccolo</i> AtLocation band
person CapableOf 心情不好(<i>have a bad mood</i>)	<i>wind instrument</i> UsedFor play music
上網(<i>surf the net</i>) HasFirst-Subevent check email	<i>woodwind</i> UsedFor make music
上網(<i>surf the net</i>) UsedFor acquire knowledge	<i>flute</i> UsedFor play orchestra
study Causes 上網(<i>surf the net</i>)	<i>mandolin</i> UsedFor Jazz
computer UsedFor 上PTT(<i>log on PTT</i>) ¹	<i>alcoholic drink</i> AtLocation bar
dolphin CapableOf 吃魚(<i>eat fish</i>)	<i>pail</i> UsedFor contain
打b(<i>use BBS</i>) ² HasSubevent typing	<i>vessel</i> UsedFor carry water

1. Our method succeeded in finding new assertions given concepts in Chinese ConceptNet without glossary mapping nor good machine translation result. For example, “上網 (surf the net)” is translated to “Internet” by Google Translate. Unfortunately, features of “Internet” produce problematic descriptions about “上網 (surf the net).” The proposed method maps “上網 (surf the net)” to “learn” and “work” via analogy and glossary mapping constraint and created the new assertion “study Causes 上網 (surf the net)” based on common features of “learn” and “work”. The other assertions listed in the left column of Table 3 are discovered for concepts “上PTT (log on PTT)”, “打b (use BBS)”, and “心情不好 (have a bad mood)” etc, which cannot be found in the dictionary nor machine translation.
2. The new assertions found for concepts in the same WordNet hierarchy are similar to each other. For example, wind instrument, woodwind, and flute are related to play music

¹PTT is the largest bulletin board system in Taiwan.

²BBS is the abbreviation of bulletin board system.

or make music (see the right column in Table 3). This result corresponds to our intuition that the concepts in the same hierarchy shares common features, and therefore we can transfer the features from one concept to describe another concept. From the example of WordNet, our method also reflects the structure of a knowledge base and has the ability to adapt to any knowledge base with different structure from ConceptNet.

3. We can easily explain why an assertion is produced for a concept by reversing our process. As the example illustrated in *Concept Mapping by Analogy* section, we are able to generate explanations like “狗 (dog) HasA fur because lion/cat HasA fur and 狗 (dog) is analogous to lion/cat in some properties.” Such explanations may help convince users of the resulting assertions.

Evaluation

To ensure that the knowledge base is extended correctly using the proposed method, we need to find out if the new assertions discovered are good common sense. For quantitative evaluation of our results, all new assertions found are rated by online players of the Virtual Pet Game (Kuo et al. 2009). In our experiments, there are three sources of the new assertions: 1,655 assertions from analogy between Chinese and English ConceptNet with a score ≥ 2 ; 241 assertions from analogy between WordNet and English ConceptNet with a score ≥ 2 ; and 1,560 assertions from blending of Chinese and English ConceptNet.

All new assertions are shuffled with the original Chinese ConceptNet to be voted. Each assertion is rated as either good or bad by 3 randomly selected players, and it is treated as a good assertion if two or more players rated the assertion as good. Otherwise, it is considered as a bad assertion. Table4 summarizes the percentage of good assertions from these sources. The performance of the proposed method for Chinese and English ConceptNet, the proposed method for WordNet and English ConceptNet, and Blending of Chinese and English ConceptNet are 75.77%, 77.59%, and 41.03% respectively.

Table 4: Result of users’ rating on new assertions produced by different method

Method	Analogy between Chinese and English ConceptNet	Analogy between WordNet and English ConceptNet	Blending of Chinese and English ConceptNet
% of good assertions	75.77%	77.59%	41.03%

Analogy v.s. Blending

The proposed method outperforms Blending in extending Chinese ConceptNet. The linear combination in Blending brings all features from one knowledge base to their similar concepts in another knowledge base, which in turn introduces many noises to the combined KB. Noises from

such over-generalization may be problematic even if the two knowledge bases share the same structure, e.g. Chinese and English ConceptNet. For example, many concepts about animals (e.g. dog, cat, penguin, etc) are linked incorrectly to the feature “AtLocation forest” in the result from Blending. On the other hand, we suffer less from such errors in the analogy-based approach because new features were only generated from the common features of a set of concepts verified by analogy and glossary mapping constraints. These common features can be transferred better because they are not simply possible features but the prominent features with regard to the specified query concept in the source domain.

Chinese ConceptNet v.s. WordNet

The proposed method discovered over 75% of good assertions in Chinese ConceptNet and WordNet despite the percentage of overlapping concepts are very low in both networks (1.27% and 4.79% respectively.) It is not difficult for the typical users or knowledge engineers to filter out the bad assertions from our results.

Unlike the experiment in Chinese ConceptNet, only a portion of WordNet (about 3000 concepts) was used in the process of finding new assertions. The positive result implies that we can choose the concepts and their neighbors in building the extended network. The quality of new assertions is the same as when the entire knowledge base is considered. Therefore, our method can be applied to any large knowledge base by dividing them into smaller networks.

Precision of top k assertions

The quality of new assertions can be further analyzed by charting the average user ratings as sorted by their scores. Figures 3 shows the percentage of good assertions for assertions with top k scores of the proposed method vs. Blending respectively in the experiment to extend Chinese ConceptNet.

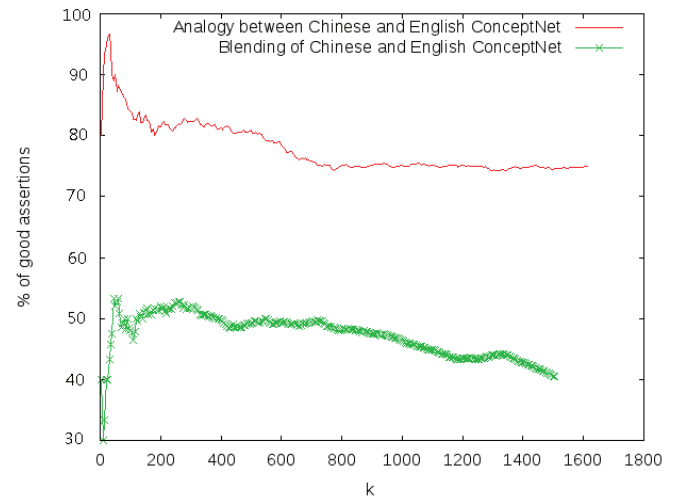


Figure 3: The % of good assertions in top k new assertions identified by Analogy vs. Blending

The score of our method is defined in the *Concept Mapping by Analogy* section and the score of Blending is defined as the cosine similarity of the specified feature and concept. One can find that the quality of our method exceeds 90% when $k \leq 50$ and maintains at least 75% for other values of k , whereas the performance of Blending was about 40% to 50% regardless of k and with the tendency to go down for larger k . Hence, we can trade off the quality and quantity of new assertions by specifying the best k . This property is helpful when we need the flexibility for different usages of the integrated knowledge bases. For example, if accuracy is important for a specific purpose, we should consider assertions with the highest scores.

Conclusion

This paper proposed a novel approach to bridging knowledge bases with analogy by graph similarity. By finding analogy and glossary mapping between two knowledge bases, we are able to leverage multiple knowledge bases to create new assertions for answering queries even if they do not share the same ontology or language. Experiments have been conducted to extend an incomplete knowledge base, such as the Chinese ConceptNet collected via the Virtual Pet Game, based on the English ConceptNet. Despite an extremely low overlap (1.27%) in concept coverage, the proposed method succeeded in finding good new assertions. Evaluation by online voting showed that the top k assertions found have a high precision of over 90%, with an average of over 75% when $k = 1600$. It significantly outperforms the precision of 41% for blending. The assertions found by analogy may be used to create new questions/answers for the Virtual Pet game to guide further data collection. Thereby, it is possible to fill the gaps between multiple knowledge bases automatically.

Acknowledgments

The authors would like to thank the anonymous reviewers for their valuable comments. This research was supported by the National Science Council (NSC 98-2622-E-002-001-CC2) and the “Environment and Behavior Perception for Intelligent Living”, a joint project by the Institute for Information Industry funded by the Ministry of Economy Affairs of Taiwan.

References

Chklovski, T. 2003. Learner: a system for acquiring commonsense knowledge by analogy. In *K-CAP '03: Proceedings of the 2nd International Conference on Knowledge Capture*.

Falkenhainer, B.; Forbus, K. D.; and Gentner, D. 1989. The structure-mapping engine: Algorithm and examples. *Artificial Intelligence* 20:1–63.

Gentner, D. 1983. Structure-mapping: A theoretical framework for analogy. *Cognitive Science* 7(2):155–170.

Havasi, C.; Pustejovsky, J.; Speer, R.; and Lieberman, H. 2009. Digital intuition: Applying common sense using dimensionality reduction. *IEEE Intelligent Systems* 24(4):24–35.

Havasi, C.; Speer, R.; and Alonso, J. 2007. Conceptnet 3: A flexible, multilingual semantic network for common sense knowledge. In *Recent Advances in Natural Language Processing*.

Krishnamurthy, J. 2009. Finding analogies in semantic networks using the singular value decomposition. Master's thesis, Massachusetts Institute of Technology.

Kuo, Y. L.; Lee, J. C.; Chiang, K. Y.; Wang, R.; Shen, E.; Chan, C. W.; and Hsu, J. Y.-j. 2009. Community-based game design: experiments on social games for commonsense data collection. In *Proceedings of the ACM SIGKDD Workshop on Human Computation*.

Lenat, D. B.; Prakash, M.; and Shepherd, M. 1986. Cyc: Using common sense knowledge to overcome brittleness and knowledge acquisition bottlenecks. *AI Magazine* 6(4):65–85.

Lenat, D. B. 1995. Cyc: A large-scale investment in knowledge infrastructure. *Communications of the ACM* 38(11):33–38.

Liu, H., and Singh, P. 2004. Conceptnet: A practical commonsense reasoning toolkit. *BT Technology Journal* 22(4):211–226.

McGuinness, D. L.; Fikes, R.; Rice, J.; and Wilder, S. 2000. The chimaera ontology environment. In *Proceedings of the Seventeenth National Conference on Artificial Intelligence and Twelfth Conference on Innovative Applications of Artificial Intelligence*.

Miller, G. A. 1995. Wordnet: A lexical database for english. *Communications of the ACM* 38(11):39–41.

Newell, A., and Ernst, G. 1965. The search for generality. In *Proceedings of IFIP Congress*.

Noy, N. F., and Musen, M. A. 2000. Prompt: Algorithm and tool for automated ontology merging and alignment. In *Proceedings of the Seventeenth National Conference on Artificial Intelligence and Twelfth Conference on Innovative Applications of Artificial Intelligence*.

Pinto, H. S., and Martins, J. P. 2001. A methodology for ontology integration. In *Proceedings of the 1st international conference on Knowledge capture*.

Singh, P.; Lin, T.; Mueller, E. T.; Lim, G.; Perkins, T.; and Zhu, W. L. 2002. Open mind common sense: Knowledge acquisition from the general public. In *On the Move to Meaningful Internet Systems, 2002 - DOA/Coop IS/ODBASE 2002 Confederated International Conferences DOA, CoopIS and ODBASE 2002*.

Singh, P. 2002. The public acquisition of commonsense knowledge. In *Proceedings of AAAI Spring Symposium*.

Sowa, J., and Majumdar, A. 2003. Analogical reasoning. In *Conceptual Structures for Knowledge Creation and Communication: Proceedings of ICCS 2003*.

Speer, R.; Havasi, C.; and Lieberman, H. 2008. Analogyspace: Reducing the dimensionality of common sense knowledge. In *Proceedings of AAAI-2008*.

Stumme, G., and Maedche, A. 2001. Fca-merge: Bottom-up merging of ontologies. In *Proceedings of the 7th International Conference on Artificial Intelligence*.

Turney, P. D. 2008. The latent relation mapping engine: Algorithm and experiments. *Journal of Artificial Intelligence Research* 33:615–655.

von Ahn, L.; Kedia, M.; and Blum, M. 2006. Verbosity: A game for collecting common-sense knowledge. In *ACM Conference on Human Factors in Computing Systems (CHI Notes)*, 75–78.