

Robust Textual Inference via Graph Matching

Aria D. Haghighi

Dept. of Computer Science
Stanford University
Stanford, CA
aria42@stanford.edu

Andrew Y. Ng

Dept. of Computer Science
Stanford University
Stanford, CA
ang@cs.stanford.edu

Christopher D. Manning

Dept. of Computer Science
Stanford University
Stanford, CA
manning@cs.stanford.edu

Abstract

We present a system for deciding whether a given sentence can be inferred from text. Each sentence is represented as a directed graph (extracted from a dependency parser) in which the nodes represent words or phrases, and the links represent syntactic and semantic relationships. We develop a learned graph matching approach to approximate entailment using the amount of the sentence's semantic content which is contained in the text. We present results on the Recognizing Textual Entailment dataset (Dagan et al., 2005), and show that our approach outperforms Bag-Of-Words and TF-IDF models. In addition, we explore common sources of errors in our approach and how to remedy them.

imate translations which structurally differ from our reference translation.

One sub-task underlying these applications is the ability to *recognize* semantic entailment; whether one piece of text follows from another. In contrast to recent work which has successfully utilized logic-based abductive approaches to inference (Moldovan et al., 2003; Raina et al., 2005b), we adopt a graph-based representation of sentences, and use graph matching approach to measure the semantic overlap of text. Graph matching techniques have proven to be a useful approach for tractable approximate matching in other domains including computer vision. In the domain of language, graphs provide a natural way to express the dependencies between words and phrases in a sentence. Furthermore, graph matching also has the advantage of providing a framework for structural matching of phrases that would be difficult to resolve at the level of individual words.

1 Introduction

A fundamental stumbling block for several NLP applications is the lack of robust and accurate semantic inference. For instance, question answering systems must be able to recognize, or infer, an answer which may be expressed differently from the query. Information extraction systems must also be able to recognize the variability of equivalent linguistic expressions. Document summarization systems must generate succinct sentences which express the same content as the original document. In Machine Translation evaluation, we must be able to recognize legit-

2 Task Definition and Data

We describe our approach in the context of the 2005 Recognizing Textual Entailment (RTE) Challenge (Dagan et al., 2005), but note that our approach easily extends to other related inference tasks. The system presented here was one component of our research group's 2005 RTE submission (Raina et al., 2005a) which was the top-ranking system according to one of the two evaluation metrics.

In the 2005 RTE domain, we are given a set of pairs, each consisting of two parts: 1) the *text*, a

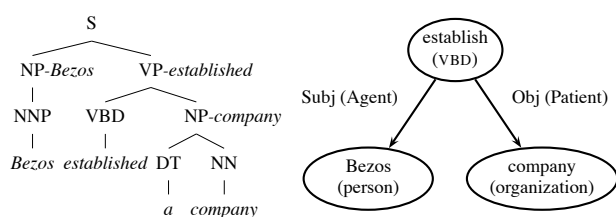


Figure 1: An example parse tree and the corresponding dependency graph. Each phrase of the parse tree is annotated with its head word, and the parenthetical edge labels in the dependency graph correspond to semantic roles.

small passage,¹ and the *hypothesis*, a single sentence. Our task is to decide if the hypothesis is “entailed” by the text. Here, “entails” does not mean strict *logical* implication, but roughly means that a competent speaker with basic world-knowledge would be happy to conclude the hypothesis given the text. This criterion has an aspect of relevance logic as opposed to material implication: while various additional background information may be needed for the hypothesis to follow, the text must substantially support the hypothesis.

Despite the informality of the criterion and the fact that the available world knowledge is left unspecified, human judges show extremely good agreement on this task – 3 human judges independent of the organizers calculated agreement rates with the released data set ranging from 91–96% (Dagan et al., 2005). We believe that this in part reflects that the task is fairly natural to human beings. For a flavor of the nature (and difficulty) of the task, see Table 1.

We give results on the data provided for the RTE task which consists of 567 development pairs and 800 test pairs. In both sets the pairs are divided into 7 tasks – each containing roughly the same number of entailed and not-entailed instances – which were used as both motivation and means for obtaining and constructing the data items. We will use the following toy example to illustrate our representation and matching technique:

Text: In 1994, Amazon.com was founded by Jeff Bezos.

Hypothesis: Bezos established a company.

¹Usually a single sentence, but occasionally longer.

3 Semantic Representation

3.1 The Need for Dependencies

Perhaps the most common representation of text for assessing content is “Bag-Of-Words” or “Bag-of-N-Grams” (Papineni et al., 2002). However, such representations lose syntactic information which can be essential to determining entailment. Consider a Question Answer system searching for an answer to *When was Israel established?* A representation which did not utilize syntax would probably enthusiastically return an answer from (the 2005 RTE text): *The National Institute for Psychobiology in Israel was established in 1979.*

In this example, it’s important to try to match relationships as well as words. In particular, any answer to the question should preserve the dependency between *Israel* and *established*. However, in the proposed answer, the expected dependency is missing although all the words are present.

Our approach is to view sentences as graphs between words and phrases, where dependency relationships, as in (Lin and Pantel, 2001), are characterized by the path between vertices.

Given this representation, we judge entailment by measuring not only how many of the *hypothesis* vertices are matched to the *text* but also how well the relationships between vertices in the hypothesis are preserved in their textual counterparts. For the remainder of the section we outline how we produce graphs from text, and in the next section we introduce our graph matching model.

3.2 From Text To Graphs

Starting with raw English text, we use a version of the parser described in (Klein and Manning, 2003), to obtain a parse tree. Then, we derive a dependency tree representation of the sentence using a slightly modified version of Collins’ head propagation rules (Collins, 1999), which make main verbs not auxiliaries the head of sentences. Edges in the dependency graph are labeled by a set of hand-created *tgrep* expressions. These labels represent “surface” syntax relationships such as *subj* for subject and *amod* for adjective modifier, similar to the relations in *Minipar* (Lin and Pantel, 2001). The dependency graph is the basis for our graphical representation, but it is enhanced in the following ways:

Task	Text	Hypothesis	Entailed
Question Answer (QA)	Prince Charles was previously married to Princess Diana, who died in a car crash in Paris in August 1997.	Prince Charles and Princess Diana got married in August 1997.	False
Machine Translation (MT)	Sultan Al-Shawi, a.k.a the Attorney, said during a funeral held for the victims, "They were all children of Iraq killed during the savage bombing."	The Attorney, said at the funeral, "They were all Iraqis killed during the brutal shelling."	True
Comparable Documents (CD)	Napster, which started as an unauthorized song-swapping Web site, has transformed into a legal service offering music downloads for a monthly fee.	Napster illegally offers music downloads.	False
Paraphrase Recognition (PP)	Kerry hit Bush hard on his conduct on the war in Iraq.	Kerry shot Bush.	False
Information Retrieval (IR)	The country's largest private employer, Wal-Mart Stores Inc., is being sued by a number of its female employees who claim they were kept out of jobs in management because they are women.	Wal-Mart sued for sexual discrimination.	True

Table 1: Some Textual Entailment examples. The last three demonstrate some of the harder instances.

1. Collapse Collocations and Named-Entities: We "collapse" dependency nodes which represent named entities (e.g., *Jeff Bezos* in Figure fig-example) and also collocations listed in WordNet, including verbs and their adjacent particles (e.g., *blow_off* in *He blew off his work*).
2. Dependency Folding: As in (Lin and Pantel, 2001), we found it useful to fold certain dependencies (such as modifying prepositions) so that modifiers became labels connecting the modifier's governor and dependent directly. For instance, in the text graph in Figure 2, we have changed *in* from a word into a relation between its head verb and the head of its NP complement.
3. Semantic Role Labeling: We also augment the graph representation with Probank-style semantic roles via the system described in (Toutanova et al., 2005). Each predicate adds an arc labeled with the appropriate semantic role to the head of the argument phrase. This helps to create links between words which share a deep semantic relation not evident in the surface syntax. Additionally, modifying phrases are labeled with their semantic types (e.g., *in 1991* is linked by a *Temporal* edge in the text graph of Figure 2), which should be useful in Question Answering tasks.
4. Coreference Links: Using a co-rereference resolution tagger, **coref** links are added through-

out the graph. These links allowed connecting the referent entity to the vertices of the referring vertex. In the case of multiple sentence texts, it is our only "link" in the graph between entities in the two sentences.

For the remainder of the paper, we will refer to the text as T and hypothesis as H , and will speak of them in graph terminology. In addition we will use H_V and H_E to denote the vertices and edges, respectively, of H .

4 Entailment by Graph Matching

We take the view that a hypothesis is entailed from the text when the cost of matching the hypothesis graph to the text graph is low. For the remainder of this section, we outline a general model for assigning a match cost to graphs.

For hypothesis graph H , and text graph T , a *matching* M is a mapping from the vertices of H to those of T . For vertex v in H , we will use $M(v)$ to denote its "match" in T . As is common in statistical machine translation, we allow nodes in H to map to fictitious NULL vertices in T if necessary. Suppose the cost of matching M is $\text{Cost}(M)$. If \mathcal{M} is the set of such matchings, we define the cost of matching H to T to be

$$\text{MatchCost}(H, T) = \min_{M \in \mathcal{M}} \text{Cost}(M) \quad (1)$$

Suppose we have a model, $\text{VertexSub}(v, M(v))$, which gives us a cost in $[0, 1]$, for substituting vertex v in H for $M(v)$ in T . One natural cost model

is to use the normalized cost for each of the vertex substitutions in M :

$$\text{VertexCost}(M) = \frac{1}{Z} \sum_{v \in H_V} w(v) \text{VertexSub}(v, M(v)) \quad (2)$$

Here, $w(v)$ represents the weight or relative importance for vertex v , and $Z = \sum_{v \in H_V} w(v)$ is a normalization constant. In our implementation, the weight of each vertex was based on the part-of-speech tag of the word or the type of named entity, if applicable. However, there are several other possibilities including using TF-IDF weights for words and phrases.

Notice that when $\text{Cost}(M)$ takes the form of (2), computing $\text{MatchCost}(H, T)$ is equivalent to finding the minimal cost bipartite graph-matching, which can be efficiently computed using linear programming.

We would like our cost-model to incorporate some measure of how relationships in H are preserved in T under M . Ideally, a matching should preserve all local relationships; i.e, if $v \rightarrow v' \in H_E$, then $M(v) \rightarrow M(v') \in T_E$. When this condition holds for all edges in H , H is isomorphic to a subgraph of T .

What we would like is an *approximate* notion of isomorphism, where we penalize the distortion of each edge relation in H . Consider an edge $e = (v, v') \in H_E$, and let $\phi_M(e)$ be the path from $M(v)$ to $M(v')$ in T .

Again, suppose we have a model, $\text{PathSub}(e, \phi_M(e))$ for assessing the “cost” of substituting a direct relation $e \in H_E$ for its counterpart, $\phi_M(e)$, under the matching. This leads to a formulation similar to (2), where we consider the normalized cost of substituting each edge relation in H with a path in T :

$$\text{RelationCost}(M) = \frac{1}{Z} \sum_{e \in H_E} w(e) \text{PathSub}(e, \phi_M(e)) \quad (3)$$

where $Z = \sum_{e \in H_E} w(e)$ is a normalization constant. As in the vertex case, we have weights for each hypothesis edge, $w(e)$, based upon the edge’s label; typically subject and object relations are more important to match than others. Our final matching cost is given by a convex mixture of

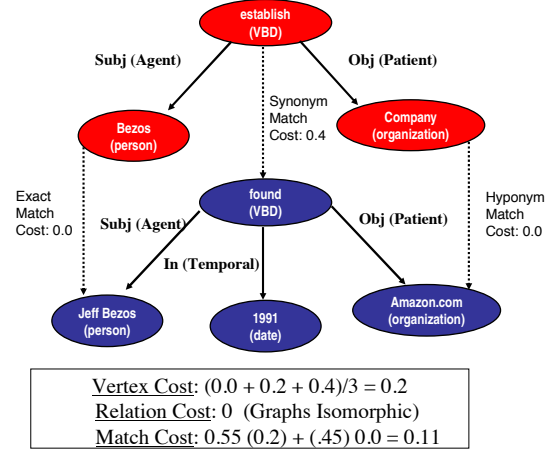


Figure 2: Example graph matching ($\alpha = 0.55$) for example pair. Dashed lines represent optimal matching.

the vertex and relational match costs: $\text{Cost}(M) = \alpha \text{VertexCost}(M) + (1 - \alpha) \text{RelationCost}(M)$.

Notice that minimizing $\text{Cost}(M)$ is computationally hard since if our PathSub model assigns zero cost only for preserving edges, then $\text{RelationCost}(M) = 0$ if and only if H is isomorphic to a subgraph of T . Since subgraph isomorphism is an NP-complete problem, we cannot hope to have an efficient exact procedure for minimizing the graph matching cost. As an approximation, we can efficiently find the matching M^* which minimizes $\text{VertexCost}(\cdot)$; we then perform local greedy hill-climbing search, beginning from M^* , to approximate the minimal matching. The allowed operations are changing the assignment of any hypothesis vertex to a text one, and, to avoid ridges, swapping two hypothesis assignments

5 Node and Edge Substitution Models

In the previous section we described our graph matching model in terms of our VertexSub model, which gives a cost for substituting one graph vertex for another, and PathSub , which gives a cost for substituting the path relationship between two paths in one graph for that in another. We now outline these models.

5.1 Vertex substitution cost model

Our $\text{VertexSub}(v, M(v))$ model is based upon a sliding scale, where progressively higher costs are

given based upon the following conditions:

- **Exact Match:** v and $M(v)$ are identical words/phrases.
- **Stem Match:** v and $M(v)$'s stems match or one is a derivational form of the other; e.g., matching *coaches* to *coach*.
- **Synonym Match:** v and $M(v)$ are synonyms according to *WordNet* (Fellbaum, 1998). In particular we use the top 3 senses of both words to determine synsets.
- **Hypernym Match:** v is a "kind of" $M(v)$, as determined by *WordNet*. Note that this feature is asymmetric.
- **WordNet Similarity:** v and $M(v)$ are similar according to *WordNet::Similarity* (Pedersen et al., 2004). In particular, we use the measure described in (Resnik, 1995). We found it useful to only use similarities above a fixed threshold to ensure precision.
- **LSA Match:** v and $M(v)$ are distributionally similar according to a freely available Latent Semantic Indexing package,² or for verbs similar according to *VerbOcean* (Chklovski and Pantel, 2004).
- **POS Match:** v and $M(v)$ have the same part of speech.
- **No Match:** $M(v)$ is NULL.

Although the above conditions often produce reasonable matchings between text and hypothesis, we found the recall of these lexical resources to be far from adequate. More robust lexical resources would almost certainly boost performance.

5.2 Path substitution cost model

Our $\text{PathSub}(v \rightarrow v', M(v) \rightarrow M(v'))$ model is also based upon a sliding scale cost based upon the following conditions:

- **Exact Match:** $M(v) \rightarrow M(v')$ is an en edge in T with the same label.
- **Partial Match:** $M(v) \rightarrow M(v')$ is an en edge in T , not necessarily with the same label.
- **Ancestor Match:** $M(v)$ is an ancestor of $M(v')$. We use an exponentially increasing cost for longer distance relationships.

- **Kinked Match:** $M(v)$ and $M(v')$ share a common parent or ancestor in T . We use an exponentially increasing cost based on the maximum of the node's distances to their least common ancestor in T .

These conditions capture many of the common ways in which relationships between entities are distorted in semantically related sentences. For instance, in our system, a partial match will occur whenever an edge type differs in detail, for instance use of the preposition *towards* in one case and *to* in the other. An ancestor match will occur whenever an indirect relation leads to the insertion of an intervening node in the dependency graph, such as matching *John is studying French farming* vs. *John is studying French farming practices*.

5.3 Learning Weights

Is it possible to learn weights for the relative importance of the conditions in the *VertexSub* and *PathSub* models? Consider the case where match costs are given only by equation (2) and vertices are weighted uniformly ($w(v) = 1$). Suppose that $\Phi(v, M(v))$ is a vector of features³ indicating the cost according to each of the conditions listed for matching v to $M(v)$. Also let w be weights for each element of $\Phi(v, M(v))$. First we can model the substitution cost for a given matching as:

$$\text{VertexSub}(v, M(v)) = \frac{\exp(w^T \Phi(v, M(v)))}{1 + \exp(w^T \Phi(v, M(v)))}$$

Letting $s(\cdot)$ be the 1-sigmoid function used in the right hand side of the equation above, our final matching cost as a function of w is given by

$$c(H, T; w) = \min_{M \in \mathcal{M}} \frac{1}{|H_V|} \sum_{v \in H} s(w^T \Phi(v, M(v))) \quad (4)$$

Suppose we have a set of text/hypothesis pairs, $\{(T^{(1)}, H^{(1)}), \dots, (T^{(n)}, H^{(n)})\}$, with labels $y^{(i)}$ which are 1 if $H^{(i)}$ is entailed by $T^{(i)}$ and 0 otherwise. Then we would like to choose w to minimize costs for entailed examples and maximize it for non-entailed pairs:

³In the case of our "match" conditions, these features will be binary.

²Available at <http://infomap.stanford.edu>

$$\ell(w) = \sum_{i:y^{(i)}=1} \log c(H^{(i)}, T^{(i)}; w) + \sum_{i:y^{(i)}=0} \log(1 - c(H^{(i)}, T^{(i)}; w))$$

Unfortunately, $\ell(w)$ is not a convex function. Notice that the cost of each matching, M , implicitly depends on the current setting of the weights w . It can be shown that since each $c(H, T; w)$ involves minimizing $M \in \mathcal{M}$, which depends on w , it is not convex. Therefore, we can't hope to globally optimize our cost functions over w and must settle for an approximation.

One approach is to use coordinate ascent over M and w . Suppose that we begin with arbitrary weights and given these weights choose $M^{(i)}$ to minimize each $c(H^{(i)}, T^{(i)}; w)$. Then we use a relaxed form of the cost function where we use the matchings found in the last step:

$$\hat{c}(H^{(i)}, T^{(i)}; w) = \frac{1}{|H_V|} \sum_{v \in H} s(w^T \Phi(v, M^{(i)}(v)))$$

Then we maximize w with respect to $\ell(w)$ with each $c(\cdot)$ replaced with the cost-function $\hat{c}(\cdot)$. This step involves only logistic regression. We repeat this procedure until our weights converge.

To test the effectiveness of the above procedure we compared performance against baseline settings using a random split on the development set. Picking each weight uniformly at random resulted in 53% accuracy. Setting all weights identically to an arbitrary value gave 54%. The procedure above, where the weights are initialized to the same value, resulted in an accuracy of 57%. However, we believe there is still room for improvement since carefully-hand chosen weights results in comparable performance to the learned weights on the final test set. We believe this setting of learning under matchings is a rather general one and could be beneficial to other domains such as Machine Translation. In the future, we hope to find better approximation techniques for this problem.

6 Checks

One systematic source of error coming from our basic approach is the implicit assumption of upwards

monotonicity of entailment; i.e., if T entails H then adding *more* words to T should also give us a sentence which entails H . This assumption, also made by other recent abductive approaches (Moldovan et al., 2003), does not hold for several classes of examples. Our formalism does not at present provide a general solution to this issue, but we include special case handling of the most common types of cases, which we outline below.⁴ These checks are done after graph matching and assume we have stored the minimal cost matching.

Negation Check

Text: Clinton's book is not a bestseller

Hypothesis: Clinton's book is a bestseller

To catch such examples, we check that each hypothesis verb is not matched to a text word which is negated (unless the verb pairs are antonyms) and vice versa. In this instance, the *is* in H , denoted by is_H , is matched to is_T which has a negation modifier, not_T , absent for is_H . So the negation check fails.

Factive Check

Text: Clonaid claims to have cloned 13 babies worldwide.

Hypothesis: Clonaid has cloned 13 babies.

Non-factive verbs (*claim*, *think*, *charged*, etc.) in contrast to factive verbs (*know*, *regret*, etc.) have sentential complements which do not represent true propositions. We detect such cases, by checking that each verb in H that is matched in T does not have a non-factive verb for a parent.

Superlative Check

Text: The Osaka World Trade Center is the tallest building in Western Japan.

Hypothesis: The Osaka World Trade Center is the tallest building in Japan.

In general, superlative modifiers (*most*, *biggest*, etc.) invert the typical monotonicity of entailment and must be handled as special cases. For any noun n with a superlative modifier (part-of-speech JJS) in H , we must ensure that all modifier relations of $M(n)$ are preserved in H . In this example, *building_H* has a superlative modifier *tallest_H*, so we must ensure that each modifier relation of *Japan_T*, a noun

⁴All the examples are actual, or slightly altered, RTE examples.

Method	Accuracy	CWS
Random	50.0%	0.500
Bag-Of-Words	49.5%	0.548
TF-IDF	51.8%	0.560
GM-General	56.8%	0.614
GM-ByTask	56.7%	0.620

Table 2: Accuracy and confidence weighted score (CWS) for test set using various techniques.

dependent of $building_T$, has a $Western_T$ modifier not in H . So it fails the superlative check.

Additionally, during error analysis on the development set, we spotted the following cases where our VertexSub function erroneously labeled vertices as similar, and required special case consideration:

- **Antonym Check:** We consistently found that the `WordNet::Similarity` modules gave high-similarity to antonyms.⁵ We explicitly check whether a matching involved antonyms and reject unless one of the vertices had a negation modifier.
- **Numeric Mismatch:** Since numeric expressions typically have the same part-of-speech tag (CD), they were typically matched when exact matches could not be found. However, mismatching numerical tokens usually indicated that H was not entailed, and so pairs with a numerical mismatch were rejected.

7 Experiments and Results

For our experiments we used the development and test sets from the Recognizing Textual Entailment challenge (Dagan et al., 2005). We give results for our system as well as for the following systems:

- **Bag-Of-Words:** We tokenize the text and hypothesis and strip the function words, and stem the resulting words. The cost is given by the fraction of the hypothesis not matched in the text.
- **TF-IDF:** Similar to Bag-Of-Words except that there is a `tf.idf` weight associated with each hypothesis word so that more “important” words are higher weight for matching.

⁵This isn’t necessarily incorrect, but is simply not suitable for textual inference.

Task	GM-General		GM-ByTask	
	Accuracy	CWS	Accuracy	CWS
CD	72.0%	0.742	76.0%	0.771
IE	55.9%	0.583	55.8%	0.595
IR	52.2%	0.564	51.1%	0.572
MT	50.0%	0.497	43.3%	0.489
PP	58.0%	0.741	58.0%	0.746
QA	53.8%	0.537	55.4%	0.556
RC	52.1%	0.539	52.9%	0.523

Table 3: Accuracy and confidence weighted score (CWS) split by task on the RTE test set.

We also present results for two graph matching (GM) systems. The GM-General system fits a single global threshold from the development set. The GM-ByTask system fits a different threshold for each of the tasks.

Our results are summarized in Table 2. As the result indicates, the task is particularly hard; all RTE participants scored between 50% and 60% in terms of overall accuracy (Dagan et al., 2005). Nevertheless, both GM systems perform better than either Bag-Of-Words or TF-IDF. CWS refers to Confidence Weighted Score (also known as average precision). This measure is perhaps a more insightful measure, since it allows the inclusion of a ranking of answers by confidence and assesses whether you are correct on the pairs that you are most confident that you know the answer to. To assess CWS, our n answers are sorted in decreasing order by the confidence we return, and then for each i , we calculate a_i , our accuracy on our i most confident predictions. Then $CWS = \frac{1}{n} \sum_{i=1}^n a_i$.

We also present results on a per-task basis in Table 3. Interestingly, there is a large variation in performance depending on the task.

8 Conclusion

We have presented a learned graph matching approach to approximating textual entailment which outperforms models which only match at the word level, and is competitive with recent weighed abduction models (Moldovan et al., 2003). In addition, we explore problematic cases of nonmonotonicity in entailment, which are not naturally handled by either subgraph matching or the so-called “logic form”

Text	Hypothesis	True Ans.	Our Ans.	Conf	Comments
A Filipino hostage in Iraq was released.	A Filipino hostage was freed in Iraq.	True	True	0.84	Verb rewrite is handled. Phrasal ordering does not affect cost.
The government announced last week that it plans to raise oil prices.	Oil prices drop.	False	False	0.95	High cost given for substituting word for its antonym.
Shrek 2 rang up \$92 million.	Shrek 2 earned \$92 million.	True	False	0.59	Collocation “rang up” is not known to be similar to “earned”.
Sonia Gandhi can be defeated in the next elections in India by BJP.	Sonia Gandhi is defeated by BJP.	False	True	0.77	“can be” does not indicate the complement event occurs.
Fighters loyal to Moqtada al-Sadr shot down a U.S. helicopter Thursday in the holy city of Najaf.	Fighters loyal to Moqtada al-Sadr shot down Najaf.	False	True	0.67	Should recognize non-Location cannot be substituted for Location.
C and D Technologies announced that it has closed the acquisition of Datel, Inc.	Datel Acquired C and D technologies.	False	True	0.64	Failed to penalize switch in semantic role structure enough

Table 4: Analysis of results on some RTE examples along with out guesses and confidence probabilities

inference of (Moldovan et al., 2003) and have proposed a way to capture common cases of this phenomenon. We believe that the methods employed in this work show much potential for improving the state-of-the-art in computational semantic inference.

9 Acknowledgments

Many thanks to Rajat Raina, Christopher Cox, Kristina Toutanova, Jenny Finkel, Marie-Catherine de Marneffe, and Bill MacCartney for providing us with linguistic modules and useful discussions. This work was supported by the Advanced Research and Development Activity (ARDA)’s Advanced Question Answering for Intelligence (AQUAINT) program.

References

- Timothy Chklovski and Patrick Pantel. 2004. VerbOcean: Mining the web for fine-grained semantic verb relations. In *EMNLP*.
- Michael Collins. 1999. *Head-driven statistical models for natural language parsing*. Ph.D. thesis, University of Pennsylvania.
- Ido Dagan, Oren Glickman, and Bernardo Magnini. 2005. The PASCAL recognizing textual entailment challenge. In *Proceedings of the PASCAL Challenges Workshop Recognizing Textual Entailment*.
- C. Fellbaum. 1998. *WordNet: An Electronic Lexical Database*. MIT Press.
- Dan Klein and Christopher D. Manning. 2003. Accurate unlexicalized parsing. In *ACL*, pages 423–430.
- Dekang Lin and Patrick Pantel. 2001. DIRT - discovery of inference rules from text. In *Knowledge Discovery and Data Mining*, pages 323–328.
- Dan I. Moldovan, Christine Clark, Sanda M. Harabagiu, and Steven J. Maorano. 2003. Cogex: A logic prover for question answering. In *HLT-NAACL*.
- K. Papineni, S. Roukos, T. Ward, and W. Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *ACL*.
- Ted Pedersen, Siddharth Parwardhan, and Jason Michelizzi. 2004. Wordnet::similarity – measuring the relatedness of concepts. In *AAAI*.
- Rajat Raina, Aria Haghighi, Christopher Cox, Jenny Finkel, Jeff Michels, Kristina Toutanova, Bill MacCartney, Marie-Catherine de Marneffe, Christopher D. Manning, and Andrew Y. Ng. 2005a. Robust textual inference using diverse knowledge sources. In *Proceedings of the First PASCAL Challenges Workshop*. Southampton, UK.
- Rajat Raina, Andrew Y. Ng, and Christopher D. Manning. 2005b. Robust textual inference via learning and abductive reasoning. In *Proceedings of AAAI 2005*. AAAI Press.
- Philip Resnik. 1995. Using information content to evaluate semantic similarity in a taxonomy. In *IJCAI*, pages 448–453.
- Kristina Toutanova, Aria Haghighi, and Christopher Manning. 2005. Joint learning improves semantic role labeling. In *Association of Computational Linguistics (ACL)*.