# Acquiring Word-Meaning Mappings
# for Natural Language Interfaces

**Cynthia A. Thompson**                                   CINDI@CS.UTAH.EDU
*School of Computing, University of Utah*
*Salt Lake City, UT 84112-3320*

**Raymond J. Mooney**                                   MOONEY@CS.UTEXAS.EDU
*Department of Computer Sciences, University of Texas*
*Austin, TX 78712-1188*

## Abstract

This paper focuses on a system, WOLFIE (WOrd Learning From Interpreted Examples), that acquires a semantic lexicon from a corpus of sentences paired with semantic representations. The lexicon learned consists of phrases paired with meaning representations. WOLFIE is part of an integrated system that learns to parse representations such as logical database queries.

Experimental results are presented demonstrating WOLFIE's ability to learn useful lexicons for a database interface in four different natural languages. The usefulness of the lexicons learned by WOLFIE are compared to those acquired by a similar system developed by Siskind (1996), with results favorable to WOLFIE. A second set of experiments demonstrates WOLFIE's ability to scale to larger and more difficult, albeit artificially generated, corpora.

In natural language acquisition, it is difficult to gather the annotated data needed for supervised learning; however, unannotated data is fairly plentiful. Active learning methods (Cohn, Atlas, & Ladner, 1994) attempt to select for annotation and training only the most informative examples, and therefore are potentially very useful in natural language applications. However, most results to date for active learning have only considered standard classification tasks. To reduce annotation effort while maintaining accuracy, we apply active learning to semantic lexicons. We show that active learning can significantly reduce the number of annotated examples required to achieve a given level of performance.

## 1. Introduction & Overview

A long-standing goal for the field of artificial intelligence is to enable computer understanding of human languages. Much progress has been made in reaching this goal, but much also remains to be done. Before artificial intelligence systems can meet this goal, they first need the ability to *parse* sentences, or transform them into a representation that is more easily manipulated by computers. Several knowledge sources are required for parsing, such as a grammar, lexicon, and parsing mechanism.

Natural language processing (NLP) researchers have traditionally attempted to build these knowledge sources by hand, often resulting in brittle, inefficient systems that take many hundreds of hours to build. Our goal here is to overcome this "knowledge acquisition bottleneck" by applying methods from machine learning. We develop and apply methods from *empirical* or *corpus-based* NLP to learn semantic lexicons, and from *active learning* to reduce the annotation effort required to learn them.

The semantic lexicon, or the mapping from words to meanings, is one NLP component that is typically challenging and time consuming to construct and update by hand. This paper describes a system, WOLFIE (WOrd Learning From Interpreted Examples), that acquires a semantic lexicon of word/meaning pairs from a corpus of sentences paired with semantic representations. The goal is to automate lexicon construction for an integrated NLP system that acquires both semantic lexicons and parsers for natural-language interfaces from a single training set of annotated sentences.

Although a few others (Riloff & Jones, 1999; Siskind, 1996; Hastings, 1996; Brent, 1991) have presented systems for learning information about lexical semantics, the developed system is unique in combining several features. First, it interacts with a system, CHILL (Zelle & Mooney, 1996; Zelle, 1995), that learns to parse sentences into semantic representations. Second, it uses a fairly straightforward batch, greedy learning algorithm that is fast and accurate. Third, it is easily extendible to new representation formalisms. Fourth, it requires no prior knowledge although it can exploit an initial lexicon if provided. Finally, it simplifies the mapping problem by making a *compositionality* assumption that states that the meaning of a sentence is composed from the meanings of the individual words and phrases in that sentence, in addition, perhaps to some "connecting" information specific to the representation at hand. This assumption is similar to the linking rules of Jackendoff (1990), and has been used in previous work on grammar and language acquisition (e.g., Haas and Jayaraman (1997), Siskind (1996)).

We test WOLFIE's ability to acquire a semantic lexicon for a natural language interface to a geographical database using a corpus of queries collected from human subjects and annotated with their logical form. In this test, WOLFIE is integrated with CHILL, which learns parsers but requires a semantic lexicon (previously built manually). The results demonstrate that the final acquired parser performs nearly as accurately at answering novel questions when using a learned lexicon as when using a hand-built lexicon. WOLFIE is also compared to an alternative lexicon acquisition system developed by Siskind (1996), demonstrating superior performance on this task. Finally, the corpus is translated into Spanish, Japanese, and Turkish, and experiments are conducted demonstrating an ability to learn successful lexicons and parsers for a variety of languages.

A second set of experiments demonstrates WOLFIE's ability to scale to larger and more difficult, albeit artificially generated, corpora. Overall, the results demonstrate a robust ability to acquire accurate lexicons directly usable for semantic parsing. With such an integrated system, the task of building a semantic parser for a new domain is simplified. A single representative corpus of sentence/representation pairs allows the acquisition of both a semantic lexicon and parser that generalizes well to novel sentences.

While building an annotated corpus is arguably less work than building an entire NLP system, it is still not a simple task. Redundancies and errors may occur in the training data. A goal should be to also minimize the annotation effort, yet still achieve a reasonable level of generalization performance. In the case of natural language, there is frequently a large amount of unannotated text available. We would like to automatically, but intelligently, choose which of the available sentences to annotate.

We do this here using a technique called *active learning*. Active learning is an emerging research area in machine learning that features systems that automatically select the most informative examples for annotation and training (Cohn et al., 1994). The primary goal of

active learning is to reduce the number of examples that the system is trained on, thereby reducing the example annotation cost, while maintaining the accuracy of the acquired information. To demonstrate the usefulness of our active learning techniques, we compared the accuracy of parsers and lexicons learned using examples chosen by active learning for lexicon acquisition, to those learned using randomly chosen examples, finding that active learning saved significant annotation cost over training on randomly chosen examples. This savings is demonstrated in the geography query domain.

In summary, this paper demonstrates a machine learning technique for inducing semantic lexicons; and by combining this with previous research, we show that an entire natural language interface can be acquired from one training corpus. Further, we demonstrate the application of active learning techniques to minimize the number of sentences to annotate as training input for the integrated learning system.

The remainder of the paper is organized as follows. Section 2 gives more background information on CHILL and introduces Siskind's lexicon acquisition system, which we will compare to WOLFIE in Section 5. Sections 3 and 4 formally define the learning problem and describe the WOLFIE algorithm in detail. In Section 5 we present and discuss experiments evaluating WOLFIE's performance in learning lexicons in a database query domain and for an artificial corpus. Next, Section 6 describes and evaluates our use of active learning techniques for WOLFIE. Sections 7 and 8 discuss related research and future directions, respectively. Finally, Section 9 summarizes our research and results.

## 2. Background

In this section we give an overview of CHILL, the system we are building on. We also describe Jeff Siskind's lexicon acquisition system.

### 2.1 CHILL

The output produced by WOLFIE can be used to assist a larger language acquisition system; in particular, it is currently used as part of the input to a parser acquisition system called CHILL (Constructive Heuristics Induction for Language Learning). CHILL uses *inductive logic programming* (Muggleton, 1992; Lavrač & Džeroski, 1994) to learn a deterministic shift-reduce parser (Tomita, 1986) written in Prolog. The input to CHILL is a corpus of sentences paired with semantic representations, the same input required by WOLFIE. The parser learned is capable of mapping the sentences into their correct representations, as well as generalizing well to novel sentences. In this paper, we limit our discussion to CHILL's ability to acquire parsers that map natural-language questions directly into Prolog queries that can be executed to produce an answer (Zelle & Mooney, 1996). Following are two sample queries for a database on U.S. geography, paired with their corresponding Prolog query:

What is the capital of the state with the biggest population?
`answer(C, (capital(S,C), largest(P, (state(S), population(S,P)))))`.

What state is Texarkana located in?
`answer(S, (state(S), eq(C,cityid(texarkana,_)), loc(C,S)))`.
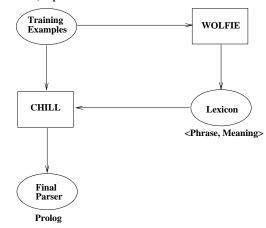
Figure 1: The Integrated System

CHILL treats parser induction as the problem of learning rules to control the actions of a shift-reduce parser. During parsing, the current context is maintained in a stack and a buffer containing the remaining input. When parsing is complete, the stack contains the representation of the input sentence. There are three types of operators the parser uses to construct logical queries. One is the introduction onto the stack of a predicate needed in the sentence representation due to a phrase's appearance at the front of the input buffer. These operators require a semantic lexicon as background knowledge. By using WOLFIE, the lexicon is provided automatically. For details on this and the other two parsing operators, see Zelle and Mooney (1996). Figure 1 illustrates the complete system.

CHILL solves the parser acquisition problem by learning rules to control the step by step actions of an initial, overly-general parsing shell. While the initial training examples are sentence/representation pairs, the examples given to the inductive logic programming component are positive and negative examples of states of the parser in which a particular operator should or should not be applied. These examples are automatically constructed by determining what sequence of operator applications leads to the correct parse. However, the overall learning task for which training data is provided is not a classification task.

## 2.2 Jeff Siskind's Lexicon Learning Research

The most closely related previous research into automated lexicon acquisition is that of Siskind (1996). As we will be comparing our system to his in Section 5, we describe the main features of his research in this section. His goal is to gain insight into the cognitive processes used by children in learning the lexicon, while we focus on learning a lexicon that will be useful for parsing. The form of the input required by his system is slightly different than for ours, in that he requires the representation of a sentence to be variable free, e.g., "John walked to school." would be represented as `go(john, to(school))`. His claim is that sentences do not contain unfilled argument positions. While this may be true for declarative sentences, questions may contain unfilled argument positions. As demonstrated

4

```
([capital], capital(_,_)),         ([state], state(_)),
([biggest], largest(_,_)),         ([in], loc(_,_)),
([highest,point], high_point(_,_)), ([long], len(_,_)),
([through], traverse(_,_)),         ([capital], capital(_)),
([has], loc(_,_))
```

Figure 2: Sample Semantic Lexicon

by the queries in the previous section, we do allow variables in the sentence representations presented to CHILL and WOLFIE.

His system takes an incremental, version-space approach to acquiring a lexicon. Learning proceeds in two stages. The first stage learns *which* symbols in the representation are to be used in the final "conceptual expression" that represents the meaning of a word. The second stage learns *how* these symbols are put together to form the final representation. For example, when learning the meaning of the word **raise**, the algorithm may learn the set {CAUSE, GO, UP} during the first stage and put them together to form the expression CAUSE($x$, GO($y$, UP)) during the second stage.

Siskind (1996) shows the effectiveness of his approach on a series of artificial corpora. While the system handles noise, ambiguity, referential uncertainty, and very large corpora, the usefulness of lexicons learned is only compared to the "correct," artificial lexicon. No demonstration of the system's usefulness in performing real natural language tasks has been performed until now.

## 3. The Lexicon Acquisition Problem

Although in the end our goal is to acquire an entire natural language interface, we currently divide the task into two parts, the lexicon acquisition component and the parser acquisition component. In this section, we discuss the problem of acquiring semantic lexicons that assist parsing and the acquisition of parsers. The training input consists of natural language sentences paired with their meaning representations. From these pairs we extract a lexicon consisting of phrases paired with their meaning representations. Some training pairs were given in the previous section, and a sample lexicon is shown in Figure 2.

To present the learning problem more formally, some definitions are needed. While in the following we use the terms "string" and "substring," these extend straight-forwardly to natural language sentences and phrases, respectively. We also refer to labeled trees, making the assumption that the semantic meanings of interest can be represented as such. Most common representations can be recast as labeled trees or forests, and our formalism extends easily to the latter.

**Definition:** Let $\Sigma_V$, $\Sigma_E$ be finite alphabets of vertex labels and edge labels, respectively. Let $V$ be a set of vertices (nodes), $E \subseteq V \times \Sigma_E \times V$ a set of labeled edges, and $l$ a total function $l : V \to \Sigma_V$. $l$ is called a "vertex labeling." The triple $G = (V, E, l)$ is a *labeled graph*.

**Definition:** A *labeled tree* is a connected, acyclic labeled graph.

**String $s_1$: "The girl ate the pasta with the cheese."**

**Tree $t_1$:**          **$t_1$ with its node and edge labels:**



**An Interpretation from $s_1$ to $t_1$:**

$f_1$("girl") = 2
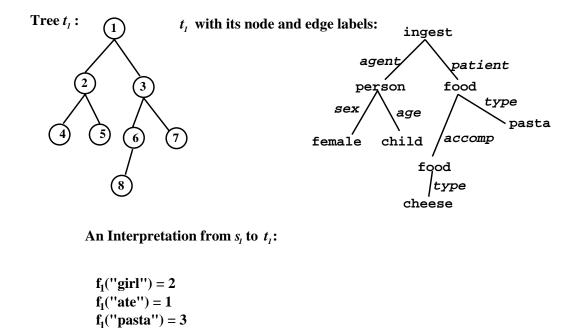$f_1$("ate") = 1
$f_1$("pasta") = 3
$f_1$("the cheese") = 6

Figure 3: Labeled Trees and Interpretations

Figure 3 shows the tree $t_1$ (with vertices 1-8) on the left, with associated node and edge labels on the right. The function $l$ is:
{(1, ingest), (2, person), (3, food), (4, female), (5, child), (6, food), (7, pasta), (8, cheese)}.
$t_1$ is a semantic representation of the sentence $s_1$: "The girl ate the pasta with the cheese." Using a conceptual dependency (Schank, 1975) representation in Prolog list form, the meaning is:

```
[ingest, agent:[person, sex:female, age:child],
      patient:[food, type:pasta, accomp:[food,type:cheese]]].
```

**Definition:** An *interpretation* $f$ from a string $s$ to a labeled, rooted tree $t$ is a one-to-one function mapping a disjoint subset of the substrings of $s$ into $nodes(t)$ such that $root(t) \in range(f)$.

The interpretation provides information about what parts of the meaning of a sentence originate from which of its phrases. In Figure 3, we show an interpretation, $f_1$ of $s_1$ to $t_1$. Note that "with" is not in the domain of $f_1$. Because we disallow overlapping substrings in the domain, both "cheese" and "the cheese" could not map to nodes in $t_1$.
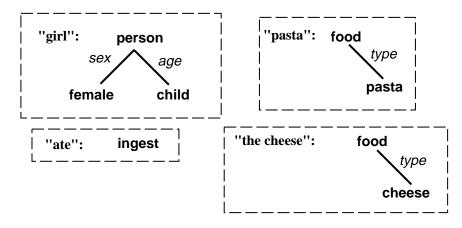
Figure 4: Meanings

**Definition:** Given an interpretation $f$ of string $s$ to tree $t$, and an element $p$ of $domain(f)$, the *meaning* of $p$ relative to $s$, $f$, $t$ is the connected subgraph of $t$ whose nodes include $f(p)$ and all its descendents *except* any other nodes of $range(f)$ and their descendents. This subgraph is also a tree, and let us assign $f(p)$ as its root.

Meanings in this sense concern the "lowest level" of phrasal meanings, occurring at the terminal nodes of a semantic grammar, namely the entries in the semantic lexicon. The grammar can then be used to construct the meanings of longer phrases and entire sentences. Figure 4 shows the meanings for each phrase in the domain of the interpretation function shown in Figure 3. We show only the labels on the nodes and arcs for readability.

**Definition:** Given a set of strings $S = \{s_1, \ldots, s_n\}$, a set of labeled trees $T = \{t_1, \ldots, t_n\}$, and a set of interpretation functions $F = \{f_1, \ldots, f_n\}$, where $f_i$ maps $s_i$ to $t_i$, let the *language* $\mathcal{L}_S = \{p_1, \ldots, p_k\}$ of $S$ be the union of all substrings[1] in the domain of $F$. For each $p_j \in \mathcal{L}_S$, the *meaning set* of $p_j$, denoted $M_{S,T,F}(p_j)$,[2] is the set of all meanings of $p_j$ relative to $s_i, t_i, f_i$ for $1 \leq i \leq n$. We consider two meanings to be the same if they are isomorphic trees taking labels into account.

For example, given sentence $s_2$: "The man ate the cheese," with associated labeled tree $t_2$ pictured in Figure 5, and $f_2$ defined as: $f_2(\text{"ate"}) = 1$, $f_2(\text{"man"}) = 2$, $f_2(\text{"the cheese"}) = 3$; the meaning set for "the cheese" with respect to $S = \{s_1, s_2\}$, $T = \{t_1, t_2\}$, $F = \{f_1, f_2\}$, is $\{[\texttt{food, type:cheese}]\}$, just one meaning though $f_1$ and $f_2$ map "the cheese" to different nodes in the two trees, because the subgraphs denoting the meaning of "the cheese" for the two functions are isomorphic.

---

1. We consider two substrings to be the same string if they contain the same characters in the same order, irrespective of their positions within the larger string in which they occur.
2. We omit the subscripts on $M$ when the sets $S$, $T$, and $F$ that the meanings are relative to are obvious from context.

**String $s_2$: "The man ate the cheese."**

**Tree $t_2$:**                                        $t_2$ **with its node and edge labels:**
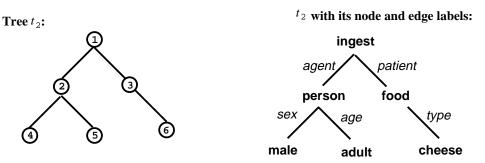


Figure 5: A Second Tree

**Definition:** Given a set of strings $S = \{s_1, \ldots, s_n\}$, a set of labeled trees $T = \{t_1, \ldots, t_n\}$, and a set of interpretation functions $F = \{f_1, \ldots, f_n\}$, the *covering lexicon* expressed by $S$, $T$, $F$ is

$$\{(p, m) : p \in \mathcal{L}_S, m \in M(p)\}.$$

The covering lexicon $L$ expressed by $S = \{s_1, s_2\}$, $T = \{t_1, t_2\}$, and $F = \{f_1, f_2\}$ is:

{ ("girl", [person, sex:female, age:child]),
("man", [person, sex:male, age:adult]), ("ate", [ingest]),
("pasta", [food,type:pasta]), ("the cheese", [food,type:cheese]) }.

The idea of a covering lexicon is that it should be sufficient to reconstruct the labeled *nodes* of the trees. Arc labels may or may not be included, since the idea is that some of them are due to syntax, which the parser will provide. The hope is that semantic arc links will be included in the lexicons learned. Because we only include in the covering lexicon phrases (substrings) that are in the domains of the $f_i$'s, words with the empty tree as meaning are not included in the covering lexicon. Note also that we will in general use "phrase" to mean substrings of sentences, whether they consist of one word, or more than one.

Finally, we are ready to define the learning problem at hand.

**The Lexicon Acquisition Problem:**
**Given:** a set of strings $S = \{s_1, \ldots, s_n\}$ and a set of labeled trees $T = \{t_1, \ldots, t_n\}$,
**Find:** a set of interpretation functions, $F = \{f_1, \ldots, f_n\}$, such that the cardinality of the covering lexicon expressed by $S$, $T$, $F$ is minimized. If such a set is found, we say we have found a *minimal* set of interpretations (or a *minimal covering lexicon*). □

Less formally, a learner is presented with a set of sentences ($S$) paired with their meanings ($T$); the goal of learning is to find the smallest lexicon consistent with this data. This lexicon is the paired listing of all phrases in $domain(F)$ (where $F$ is the set of interpretation functions found) with each of the elements in their meaning sets. The motivation for finding a lexicon of minimal size is the usual bias towards simplicity of representation and generalization beyond the training data. While this definition allows for phrases of any

length, we will usually want to limit the length of phrases to be considered for inclusion in the domain of the interpretation functions, for efficiency purposes.

Note that the covering lexicon given with the previous example is not a *minimal* covering lexicon. For the two sentences given, we could find minimal, though rather degenerate, lexicons such as:

```
{ ("girl", [[ingest, agent:[person, sex:female, age:child],
        patient:[food, type:pasta, accomp:[food,type:cheese]]]]),
  ("man", [ingest, agent:[person, sex:male, age:adult],
        patient:[food,type:cheese]]) }
```

This type of lexicon becomes less likely as the size of the corpus grows.

This definition of the lexicon acquisition problem differs from that given by other authors, including Riloff and Jones (1999), Siskind (1996), Manning (1993), Brent (1991) and others. For example, Riloff and Jones (1999) define semantic lexicons as a grouping of words into semantic categories. Manning (1993) and Brent (1991) describe work on learning selectional restrictions. Our definition focuses on deriving a word's meaning from representations of the sentences in which it appears.

Our definition of the problem makes some assumptions about the training input. First, by making $f$ a function instead of a relation, the definition assumes that the meaning for each phrase in a sentence appears once in the representation of that sentence, the *single-use* assumption. Second, by making $f$ one-to-one, it assumes *exclusivity*, that each node in a sentence's representation is due to only one phrase in the sentence. Third, it assumes that a phrase's meaning is a connected subgraph of a sentence's representation, not a more distributed representation, the *connectedness* assumption. While the first assumption may not hold for some representation languages, it does not present a problem in the domains we have considered. The second and third assumptions are perhaps less problematic with respect to general language use.

Our definition also assumes *compositionality*: that the meaning of a sentence is derived from the meanings of the phrases it contains, but is not derived from external sources such as noise. In other words, all the nodes of a sentence's representation are included within the meaning of some word or phrase in that sentence. Since we allow multi-word phrases in the lexicon (e.g., ([**kick the bucket**], die(_))), this assumption seems fairly unproblematic. In addition, future versions could allow a loosening of this restriction to allow for information about the context of surrounding sentences or of the situation to be included in the meaning of a sentence, and for noisy input.

This definition also allows training input in which:

1. Words and phrases have multiple meanings. That is, ambiguity (polysemy) can occur in the lexicon.

2. Several phrases map to the same meaning. That is, synonymy may occur in the lexicon.

3. Some words in a sentence do not map to any meanings, leaving them unused in the assignment of words to meanings.[3]

---

3. These words may, however, serve as cues to a parser on how to assemble sentence meanings from word meanings.

4. Phrases of contiguous words map to part of a sentence's meaning representation.

Of particular note is lexical ambiguity (1 above). Note that we could have also derived an ambiguous lexicon such as:

> { ("girl", [person, sex:female, age:child]), ("ate", [ingest]),
> ("ate", [ingest, agent:[person,sex:male,age:adult]]),
> ("pasta", [food,type:pasta]), ("the cheese", [food,type:cheese]) }.

from our sample corpus. In this lexicon, "ate" is an ambiguous word. The earlier example minimizes ambiguity resulting in an alternative, more intuitively pleasing lexicon. We will exploit a bias of minimal ambiguity in our learning algorithm.

Also note that our definition could extend to training input in which sentences themselves are ambiguous (paired with more than one meaning). Since we specified $S$ as a set, not a bag, the current definition does not allow duplicates. However, this is easily accommodated and allows ambiguous syntax as well as ambiguous semantics. In fact, the training data that we consider in Section 5 does have some ambiguous sentences.

Our definition of the lexicon acquisition problem does not fit cleanly into the traditional definition of learning for classification. Each training example contains a sentence and its semantic parse, and we are trying to extract semantic information about some of the phrases in that sentence. So each example potentially contains information about multiple target concepts (phrases), and we are trying to pick out the relevant "features," or nodes of the representation, corresponding to the correct meaning of each phrase. Of course, our assumptions of single-use, exclusivity, connectedness, and compositionality impose additional constraints. In addition to this "multiple examples in one" learning scenario, we do not have access to negative examples, nor can we derive any implicit negatives, because of the possibility of ambiguous and synonymous phrases.

In some ways the problem is related to clustering, which is also capable of learning multiple, potentially non-disjoint categories. However, it is not clear how a clustering system could be made to learn the phrase/meaning mappings needed for parsing. Finally, current systems that learn multiple concepts commonly use examples for other concepts as negative examples of the concept currently being learned. The implicit assumption made by doing this is that concepts are disjoint, an unwarranted assumption in the presence of synonymy.

## 4. The WOLFIE Algorithm and an Example

### 4.1 Solving the Lexicon Acquisition Problem

A first attempt to solve the Lexicon Acquisition Problem might be to examine all interpretation functions across the corpus, then choose the one(s) with minimal lexicon size. The number of possible interpretation functions for a given input pair is dependent on both the size of the sentence and its representation. In a sentence with $w$ words, there are $\Theta(w^2)$ possible phrases, not a particular challenge. However, the number of possible interpretation functions grows extremely quickly with the size of the input. For a sentence with $p$ phrases

and an associated tree with $n$ nodes, the number of possible interpretation functions is:

$$c!(n-1)! \sum_{i=1}^{c} \frac{1}{(i-1)!(n-i)!(c-i)!}.$$

where $c$ is $min(p,n)$. When $n = p$, the largest term of this equation is $c! = p!$, which grows at least exponentially with $p$, so in general the number of interpretation functions is too large to allow enumeration. Therefore, finding a lexicon by examining all interpretations across the corpus, then choosing the one(s) with minimal cardinality, is clearly not tractable.

Instead of finding all interpretations, one could find a set of candidate meanings for each phrase, from which the final meaning(s) for that phrase could be chosen in a way that minimizes lexicon size. One way to find candidate meanings is to *fracture* (Siskind, 1992) the meanings of sentences in which a phrase appears, which in our formalism corresponds to finding all possible connected subgraphs of the meaning. Fracturing would also lead to an exponential blowup in the number of candidate meanings for a phrase: A lower bound on the number of connected subgraphs for a full binary tree with $n$ nodes is obtained by noting that any subset of the $(n+1)/2$ leaves may be deleted and still maintain connectivity of the remaining tree. Thus, counting all of the ways that leaves can be deleted gives us a lower bound of $2^{(n+1)/2}$ fractures.[4] While in practice we have not encountered very large trees, this could limit the scalability of an approach that uses fractures.

Another option is to force CHILL to essentially induce a lexicon on its own. In this model, we would provide to CHILL an ambiguous lexicon in which each phrase is paired with every fracture of every sentence in which it appears. CHILL would then have to decide which set of fractures leads to the correct parse for each training sentence, and would only include those in a final learned parser-lexicon combination. Thus the search would become exponential. Furthermore, even with small representations, it would likely lead to a system with poor generalization ability.

Another difficulty with fracturing is that though it may generate candidate lexicon entries, the lexicon is no longer guaranteed to be a *covering* lexicon as defined previously. For example, Figure 6 shows an example of a sentence/representation pair and two lexicons, one that does cover the pair and one that does not. Either of these lexicons might be generated by a scheme that uses fracturing alone.

If we could efficiently find some good candidates, a standard induction algorithm could then attempt to use them as a source of training examples for each phrase. However, any attempt to use the list of candidate meanings of one phrase as negative examples for another phrase would be flawed. The learner could not know in advance which phrases are possibly synonymous, and thus which phrase lists to use as negative examples of other phrase meanings. Also, many representation components would be present in the lists of more than one phrase. This is a source of conflicting evidence for a learner, even without the presence of synonymy. Since only positive examples are available, one might think of using most specific conjunctive learning, or finding the intersection of all the representations for each phrase, as proposed by Anderson (1977). However, the meanings of a polysemous phrase are disjunctive, and this intersection would be empty.

---

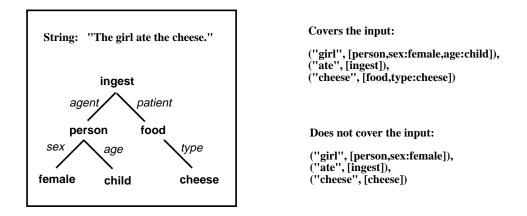4. Thanks to net-citizen Dan Hirshberg for help with this analysis

Figure 6: Covering Example

## 4.2 Our solution: WOLFIE

The above analysis leads us to believe that the Lexicon Acquisition Problem is computationally intractable. Therefore, we can not perform an efficient search for the best lexicon. Nor can we use a standard induction algorithm. Instead, our approach is to generate a set of candidate lexicon entries, from which the final learned lexicon is derived by greedily choosing the "best" lexicon item at each point, in the hopes of finding a final (minimal) covering lexicon. Our method "carve outs" one subgraph at a time from a sentence's representation and assigns it to the meaning of a phrase in that sentence. We do not actually learn interpretation functions, so do not guarantee that we will find a covering lexicon. Even if we were to search for interpretation functions, using a greedy search would also not guarantee covering the input, and of course it also does not guarantee that a minimal lexicon is found. However, we will later present experimental results demonstrating that our greedy approach performs well.

The WOLFIE algorithm, outlined in Figure 7, is an implemented attempt to find an approximate solution to the Lexicon Acquisition Problem. The algorithm's first step is to derive an initial set of candidate meanings for each phrase.[5] The second step is to choose final lexicon entries from this candidate set, one at a time, updating the candidate set as we do so, taking into account our assumptions of single-use, connectedness, and exclusivity. The basic scheme for choosing entries from the candidate set is to maximize the prediction of meanings given phrases. We attempt to find the "maximally common" meaning for each phrase while still allowing coverage of most of the input.

Let us explain the algorithm in further detail by way of an example, using Spanish instead of English to illustrate the difficulty somewhat more clearly. Consider the following corpus:

1. ¿ Cuál es el capital del estado con la población más grande?
   answer(C, (capital(S,C), largest(P, (state(S), population(S,P)))))).

2. ¿ Cuál es la punta más alta del estado con la area más grande?
   answer(P, (high_point(S,P), largest(A, (state(S), area(S,A)))))).

---

5. The current implementation is limited to one and two word phrases, but easily extended to longer phrases with a linear increase in complexity.

**For** each phrase, $p$ (of at most two words):
    1.1) Collect the training examples in which $p$ appears
    1.2) Calculate LICS from (sampled) pairs of these examples' representations
    1.3) For each $l$ in the LICS, add $(p, l)$ to the set of candidate lexicon entries
**Until** the input representations are covered, or no candidate lexicon entries remain do:
    2.1) Add the best (phrase, meaning) pair from the candidate entries to the lexicon
    2.2) Update candidate meanings of phrases in the same sentences as the phrase just learned
**Return** the lexicon of learned (phrase, meaning) pairs.

Figure 7: WOLFIE Algorithm Overview

3. ¿ En que estado se encuentra Texarkana?
```
answer(S, (state(S), eq(C,cityid(texarkana,_)), loc(C,S))).
```

4. ¿ Qué capital es la más grande?
```
answer(A, largest(A, capital(A))).
```

5. ¿ Qué es la area de los estados unitos?
```
answer(A, (area(C,A), eq(C,countryid(usa)))).
```

6. ¿ Cuál es la población de un estado que bordean a Utah?
```
answer(P, (population(S,P), state(S), next_to(S,M), eq(M,stateid(utah)))).
```

7. ¿ Qué es la punta más alta del estado con la capital Madison?
```
answer(C, (high_point(B,C), loc(C,B), state(B),
    capital(B,A), eq(A,cityid(madison,_)))).
```

Although not required, for simplification, let us assume sentences are stripped of phrases that we know have empty meanings (e.g., [**qué**], [**es**], [**con**], [**la**]) and that it is known that some phrases refer directly to given database constants (e.g., location names).

    Initial candidate meanings for a phrase are produced by computing the common substructures between sampled pairs of representations of sentences that contain it. We do not intersect all pairs for efficiency reasons, but randomly sample a subset, as in GOLEM (Muggleton & Feng, 1990). We derive common substructure by computing the Largest Isomorphic Connected Subgraphs (LICS) of two labeled trees, taking labels into account in the isomorphism. The Largest Common Subgraph problem is solvable in polynomial time if, as we assume, both inputs are trees (Garey & Johnson, 1979). The exact algorithm is complicated a bit by the variables and conjunction present in the representation above. Therefore, we use LICS with an addition similar to computing the Least General Generalization of first-order clauses (Plotkin, 1970), that is, the most specific clause subsuming two given clauses. Specifically, we find the LICS between two trees and then compute the Least General Generalization of the resulting subexpressions. For example, given the two trees:

```
answer(C, (capital(S,C), largest(P, (state(S), population(S,P))))).
```

```
answer(P, (high_point(S,P), largest(A, (state(S), area(S,A)))).,
```

the common meaning is `largest(_,state(_))`.[6]

The sets of initial candidate meanings for some of the phrases in the sample corpus are:

| Phrase | LICS | From Sentences |
|--------|------|----------------|
| [**capital**]: | `largest(_,_)` | 1,4 |
| | `capital(_,_)` | 1,7 |
| | `state(_)` | 1,7 |
| [**grande**]: | `largest(_,state(_))` | 1,2 |
| | `largest(_,_)` | 1,4; 2,4 |
| [**estado**]: | `largest(_,state(_))` | 1,2 |
| | `state(_)` | 1,3; 1,7; 2,3; 2,6; 2,7; 3,6; 6,7 |
| | `(population(S,_), state(S))` | 1,6 |
| | `capital(_,_)` | 1,7 |
| | `high_point(_,_)` | 2,7 |
| | `(state(S), loc(_,S))` | 3,7 |
| [**punta mas**]: | `high_point(_,_)` | 2,7 |
| | `state(_)` | 2,7 |
| [**encuentra**]: | `(state(S), loc(_,S))` | 3 |

Note that [**estado**] has six candidate meanings, and in some cases two meanings are generated by the same pair of representations of sentences in which words appear. For phrases appearing in only one sentence (e.g., [**encuentra**]), the entire sentence representation is used as an initial candidate meaning. Such candidates are typically generalized in step 2.2 of the algorithm to only the correct portion of the representation before they are added to the lexicon; we will see an example of this below.

After deriving initial candidates, the greedy search begins. The heuristic used to evaluate candidates attempts to help assure that a small lexicon is learned. It is the weighted sum of two components, where $p$ is the phrase and $m$ its candidate meaning:

1. $P(m \mid p) \times P(p \mid m) \times P(m) = P(p) \times P(m \mid p)^2$

2. The generality of $m$

The first component is analogous the the cluster evaluation heuristic used by COBWEB (Fisher, 1987), which measures the utility of clusters based on attribute-value pairs and categories, instead of meanings and phrases. The probabilities are estimated from the training data and then updated as learning progresses to account for phrases and meanings already covered. We will see how this updating works as we continue through our example of the algorithm. The goal of this part of the heuristic is to maximize the probability of predicting the correct meaning for a randomly sampled phrase. The equality holds by Bayes Theorem. Looking at the right side, $P(m \mid p)^2$ is the expected probability that meaning $m$ is correctly guessed for a given phrase, $p$. This assumes a strategy of probability matching, in which a meaning $m$ is chosen for $p$ with probability $P(m \mid p)$ and correct with the same probability. The other term, $P(p)$, biases the component by how common the phrase is. Interpreting the left side of the equation, the first term biases towards lexicons with low ambiguity, the second towards low synonymy, and the third towards frequent meanings.

---

6. Since CHILL initializes the parse stack with the `answer` predicate, it is first stripped from the input given to WOLFIE.

The second component of the heuristic, *generality*, is computed as the negation of the number of nodes in the meaning's tree structure, and helps prefer smaller, more general meanings. Learning a meaning with fewer terms helps evenly distribute the labels in a sentence's representation among the meanings of the phrases in that sentence, and thus leads to a lexicon that is more likely to be correct. To see this, we note that some pairs of words tend to frequently co-occur (**grande** and **estado** in our example), and so their joint representation (meaning) is likely to be in the list of candidate meanings for both words. By preferring a more general meaning, we easily ignore these incorrect joint meanings. For example, in the candidate set above, if all else were equal, the generality portion of the heuristic would prefer `state(_)`, with generality value -1, over `largest(_,state(_))` and `(state(S),loc(_,S))`, each with generality value -2, as the meaning of **estado**.

In the example below and all experiments, we use a weight of 10 for the first component of the heuristic, and a weight of 1 for the second. The first component has smaller absolute values and is therefore given a higher weight. Results are not overly-sensitive to the weights and automatically setting them using cross-validation on the training set (Kohavi & John, 1995) had little effect on overall performance. To break ties, less ambiguous (those with currently fewer meanings) and shorter phrases are preferred. Below we illustrate the calculation of the heuristic measure for some of the above fourteen pairs, and its value for all. The calculation shows the sum of multiplying 10 by the COBWEB-based component of the heuristic and multiplying 1 by the generality component. The first component is simplified as follows:

$$P(p) \times P(m \mid p)^2 = \frac{\mid p \mid}{total\_phrases} \times \frac{\mid m \cap p \mid^2}{\mid p \mid^2} = \frac{\mid m \cap p \mid^2}{\mid p \mid},$$

where $\mid p \mid$ is the number of times phrase $p$ appears in the corpus, and $\mid m \cap p \mid$ is the number of times that meaning $m$ is paired with phrase $p$. We can ignore *total_phrases* since the number of phrases in the corpus is the same for each pair, and has no effect on the ranking.

([**capital**], `largest(_,_)`): $10(2^2/3) + 1(-1) = 12.33$,

([**capital**], `capital(_,_)`): $12.33$,

([**capital**], `state(_,_)`): $12.33$,

([**grande**], `largest(_,state(_))`): $10(2^2/3) + 1(-2) = 11.3$,

([**grande**], `largest(_,_)`): $29$,

([**estado**], `largest(_,state(_))`): $10(2^2/5) + 1(-2) = 6$,

([**estado**], `state(_)`): $10(5^2/5) + 1(-1) = 49$,

([**estado**], `(population(S,_), state(S))`): $6$,

([**estado**], `capital(_,_)`): $7$,

([**estado**], `high_point(_,_)`): $7$,

([**estado**], `(state(S), loc(_,S))`): $6$,

([**punta mas**], `high_point(_,_)`): $19$,

([**punta mas**], `state(_)`): $10(2^2/2) + 1(-1) = 19$,

([**encuentra**], `(state(S), loc(_,S))`): $10(1^2/1) + 1(-2) = 8$.

The highest scoring pair is ([**estado**], state(_)), so it is added to the lexicon.

Next is the candidate generalization phase (step 2.2). One of the key ideas of the algorithm is that each (phrase, meaning) choice can constrain the candidate meanings of phrases yet to be learned. Given the assumption that each portion of the representation is due to at most one phrase in the sentence (exclusivity), once part of a representation is covered, no other phrase in the sentence can be paired with that meaning (at least for that sentence). Therefore, in step 2.2 the candidate meanings for words in the same sentences as the word just learned are generalized to exclude the representation just learned. Thus, now that state(_) is covered in sentences 1, 2, 3, 6, and 7, the candidates for several other words in those sentences are generalized. For example, the meaning (state(S), loc(_,S)) for [**encuentra**], is generalized to loc(_,_), with a new heuristic value of $10(1^2/1)+1(-1) = 9$. Also, our single use assumption allows us to remove all candidate pairs containing [**estado**] from the pool, since the learned pair covers all occurrences of [**estado**] in that pool.

Note that the pairwise matchings to generate candidate items, together with this updating of the candidate set, enable multiple meanings to be learned for ambiguous phrases. For example, note that **capital** is ambiguous in this data set, though its ambiguity is an artifact of the way that the query language was designed, and one doesn't ordinarily think of it as an ambiguous word. However, this small example illustrates the capability to handle ambiguity.

Subsequently, the greedy search continues until the resulting lexicon covers the training corpus, or until no candidate phrase meanings remain. In rare cases, learning errors occur that leave some portions of representations uncovered. In our example, the following lexicon is learned:

([**estado**], state(_)),
([**grande**], largest(_)),
([**area**], area(_)),
([**punta**], high_point(_,_)),
([**poblacion**], population(_,_)),
([**capital**], capital(_,_)),
([**encuentra**], loc(_,_)),
([**alta**], loc(_,_)),
([**bordean**], next_to(_)),
([**capital**], capital(_)).

In the next section, we discuss the ability of Wolfie to learn lexicons that are useful for parsers and parser acquisition.

## 5. Evaluation of Wolfie

The following two sections discuss experiments testing Wolfie's success in learning lexicons for both real and artificial corpora, comparing it in several cases to a previously developed lexicon learning system.

## 5.1 A Database Query Application

This section describes our experimental results on a database query application. The first corpus discussed contains 250 questions about U.S. geography, paired with their Prolog query to extract the answer to the question from a database. This domain was originally chosen due to the availability of a hand-built natural language interface, *Geobase*, to a database containing about 800 facts. *Geobase* was supplied with Turbo Prolog 2.0 (Borland International, 1988), and designed specifically for this domain. The questions in the corpus were collected by asking undergraduate students to generate English questions for this database. To broaden the test, we had the same 250 sentences translated into Spanish, Turkish, and Japanese. The Japanese translations are in word-segmented Roman orthography. Translated questions were paired with the appropriate logical queries from the English corpus.

To evaluate the learned lexicons, we measured their utility as background knowledge for CHILL. This is performed by choosing a random set of 25 test examples and then creating lexicons and parsers using increasingly larger subsets of the remaining 225 examples (increasing by 50 examples each time). After training, the test examples are parsed using the learned parser. We then submit the resulting queries to the database, compare the answers to those generated by submitting the correct representation to the database, and record the percentage of correct (matching) answers. By using the difficult "gold standard" of retrieving a correct answer, we avoid measures of partial accuracy which we believe do not adequately measure final utility. We repeated this process for ten different random training and test sets and evaluated performance differences using a two-tailed, paired $t$-test with a significance level of $p \leq 0.05$.

We compared our system to an incremental (on-line) lexicon learner developed by Siskind (1996), and originally evaluated only on artificial data. To make a more equitable comparison to our batch algorithm, we ran his in a "simulated" batch mode, by repeatedly presenting the corpus 500 times, analogous to running 500 epochs to train a neural network. We also removed WOLFIE's ability to learn phrases of more than one word, since the current version of Siskind's system does not have this ability. We also made comparisons to the parsers learned by CHILL when using a hand-coded lexicon as background knowledge.

In this and similar applications, there are many terms, such as state and city names, whose meanings can be automatically extracted from the database. Therefore, all tests below were run with such names given to the learner as an initial lexicon; this is helpful but not required.

### 5.1.1 COMPARISONS USING ENGLISH

The first experiment was a comparison of the two systems on the original English corpus. In this experiment, unlike in the example of Section 4, we did not strip sentences of phrases known to have empty meanings; however, we still provided meanings for phrases referring to database constants. Figure 8 shows learning curves for CHILL when using the lexicons learned by WOLFIE (CHILL+Wolfie) and by Siskind's system (CHILL+Siskind). The uppermost curve (CHILL+handbuilt) shows CHILL's performance when given the hand-built lexicon. CHILL-testlex shows the performance when words that never appear in the training data are deleted from the hand-built lexicon (since a learning algorithm has no chance
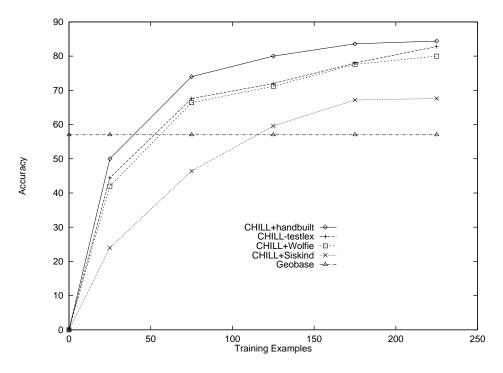
Figure 8: Accuracy on English Geography Corpus

of learning these). Finally, the horizontal line shows the performance of the *Geobase* benchmark.

The results show that a lexicon learned by WOLFIE led to parsers that were almost as accurate as those generated using a hand-built lexicon. The best accuracy is achieved by the hand-built lexicon, followed by the hand-built lexicon with words only in the test set removed, followed by WOLFIE, followed by Siskind's system. All the systems do as well or better than *Geobase* by the time they reach 125 training examples. The differences between WOLFIE and Siskind's system are statistically significant at all training example sizes. These results show that WOLFIE can learn lexicons that support the learning of successful parsers, and that are better from this perspective than those learned by a competing system. Also, comparing to the CHILL-testlex curve, we see that most of the drop in accuracy from a hand-built lexicon is due to words in the test set that the system has not seen during training. In fact, none of the differences between CHILL+Wolfie and CHILL-testlex are statistically significant.

One of the implicit hypotheses of our problem definition is that coverage of the training data implies a good lexicon. The results show a coverage of 100% of the 225 training examples for WOLFIE versus 94.4% for Siskind. In addition, the lexicons learned by Siskind's system were more ambiguous and larger than those learned by WOLFIE. WOLFIE's lexicons had an average of 1.1 meanings per word, and an average size of 56.5 words (after 225 training examples) versus 1.7 meanings per word and 154.8 entries in Siskind's lexicons. For comparison, the hand-built lexicon had 1.2 meanings per word and 88 entries. These differences undoubtedly contribute to the final performance differences.
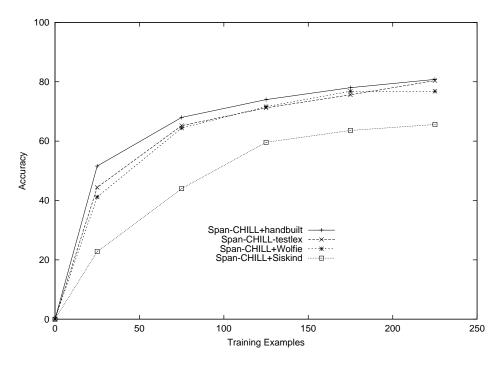
18

Figure 9: Accuracy on Spanish

### 5.1.2 Performance for Other Natural Languages

Next, we examined the performance of the two systems on the Spanish version of the corpus. Figure 9 shows the results. The differences between WOLFIE and Siskind's system are statistically significant at all training set sizes. We also again show the performance with hand-built lexicons, both with and without phrases present only in the testing set. The performance compared to the hand-built lexicon with test-set phrases removed, is still competitive, with the difference being significant only at 225 examples.

Figure 10 shows the accuracy of learned parsers with WOLFIE's learned lexicons for all four languages. This time we did not strip sentences of phrases known to have empty meanings, since this background knowledge may not be available for all languages. The performance differences among the four languages are quite small, demonstrating that our methods are not language dependent.

### 5.1.3 A Larger Corpus

Next, we present results on a larger, more diverse corpus from the geography domain, where the additional sentences were collected from computer science undergraduates in an introductory AI course. The set of questions in the smaller corpus was collected from students in a German class, with no special instructions on the complexity of queries desired. The AI students tended to ask more complex and diverse queries: their task was to give five interesting questions and the associated logical form for a homework assignment. They were requested to give at least one sentence whose representation included a predicate containing embedded predicates, for example `largest(S, state(S))`, and we asked for
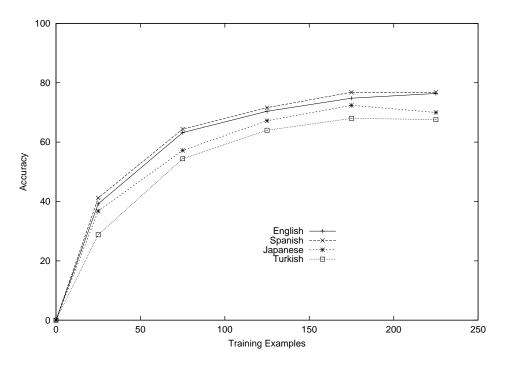
Figure 10: Accuracy on All Four Languages

variety in their sentences. There were 221 new sentences, for a total of 471 (including the original 250 sentences).

For these experiments, we split the data into 425 training sentences and 46 test sentences, for 10 random splits, then trained WOLFIE and then CHILL as before. Our goal was to see whether WOLFIE was still effective for this more difficult corpus, since there were approximately 40 novel words in the new sentences. Therefore, we tested against the performance of CHILL with an extended hand-built lexicon. For this test, we stripped sentences of phrases known to have empty meanings. We did not use phrases of more than one word, since these do not seem to make a significant difference in this domain. For these results, we include only the results for CHILL using hand-built lexicons that do not include phrases only appearing in the test set.

Figure 11 shows the resulting learning curves. The differences between CHILL and WOLFIE are statistically significant at 175, 225, 325, and 425 examples. The more mixed results here indicate both the difficulty of the domain and the more variable vocabulary. However, the improvement of machine learning methods over the *Geobase* hand-built interface is much more dramatic for this corpus.

### 5.1.4 LICS VERSUS FRACTURING

One component of the algorithm not yet evaluated explicitly is the candidate generation method. As mentioned in Section 3, we could use fractures of representations of sentences in which a phrase appears to generate the candidate meanings for that phrase, instead of LICS. We used this approach and compared it to using the standard method of using the largest isomorphic connected subgraphs of sampled pairs of representations as candidate meanings.
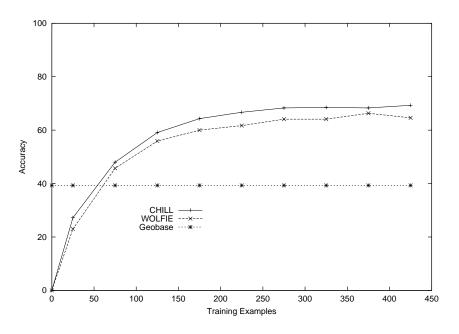
20

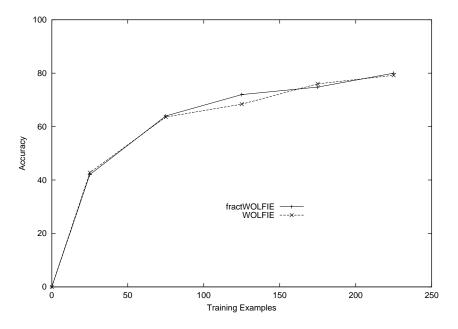Figure 11: Accuracy on the Larger Geography Corpus



Figure 12: Fracturing vs. LICS: Accuracy

To attempt a more fair comparison, we also sampled representations for fracturing, using the same number of source representations as the number of pairs sampled for LICS.

The accuracy of CHILL when using the resulting learned lexicons as background knowledge are shown in Figure 12. Using fracturing (`fractWOLFIE`) shows little or no advantage; none of the differences between the two systems are statistically significant.

Figure 13: Fracturing vs. LICS: Number of Candidates

In addition, the number of initial candidate lexicon entries from which to choose is much larger for fracturing than our traditional method, as shown in Figure 13. This is true even though we sampled the same number of representations as pairs for LICS, because there are a larger number of fractures for an arbitrary representation than the number of LICS for an arbitrary pair. Finally, WOLFIE's learning time when using fracturing is greater than that when using LICS, as shown in Figure 14, where the CPU time is shown in seconds.

In summary, these differences show the utility of LICS as a method for generating candidates: a more thorough method does not result in better performance, and also results in longer learning times. In a domain with larger representations, the differences in learning time would likely be even more dramatic.

## 5.2 Artificial Data

The previous section showed that WOLFIE successfully learns lexicons for a real-world corpus. However, this demonstrates success on only a relatively small corpus and with one representation formalism. We now show that our algorithm scales up well with more lexicon items to learn, more ambiguity, and more synonymy. These factors are difficult to control when using real data as input. Also, there are no large corpora available that are annotated with semantic parses. We therefore present experimental results on an artificial corpus. In this corpus, both the sentences and their representations are completely artificial, and the sentence representation is a variable-free representation, as suggested by the work of Jackendoff (1990) and others. We previously (Thompson, 1995) presented results demonstrating learning representations of a different form, that of a case-role representation (Fillmore, 1968) augmented with Conceptual Dependency (Schank, 1975) information.
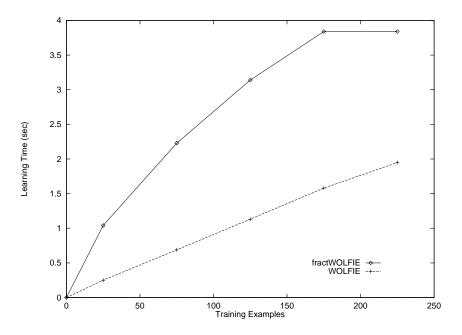
22

Figure 14: Fracturing vs. LICS: Learning Time

For each corpus presentence here, a random lexicon mapping words to simulated meanings was first constructed.[7] This *original* lexicon was then used to generate a corpus of random utterances each paired with a meaning representation. After using this corpus as input to WOLFIE, the learned lexicon was compared to the original lexicon, and *weighted precision* and *weighted recall* of the learned lexicon were measured. Precision measures the percentage of the lexicon entries (i.e., (word, meaning) pairs) that the system learns that are correct. Recall measures the percentage of the lexicon entries in the hand-built lexicon that are correctly learned by the system.

$$precision = \frac{\text{\# correct pairs}}{\text{\# pairs learned}}$$

$$recall = \frac{\text{\# correct pairs}}{\text{\# pairs in hand-built lexicon}}$$

To get weighted precision and recall measures, we then weight the results for each pair by the word's frequency in the entire corpus (not just the training corpus). This models how likely we are to have learned the correct meaning for an arbitrarily chosen word in the corpus.

We generated several lexicons and associated corpora, varying the ambiguity rate (number of meanings per word) and synonymy rate (number of words per meaning). Meaning representations were generated using a set of "conceptual symbols" that combined to form the meaning for each word. The number of conceptual symbols used in each lexicon will be noted when we describe each corpus below. In each lexicon, 47.5% of the senses were variable-free to simulate noun-like meanings, and 47.5% contained from one to three variables to denote open argument positions to simulate verb-like meanings. The remainder of

---

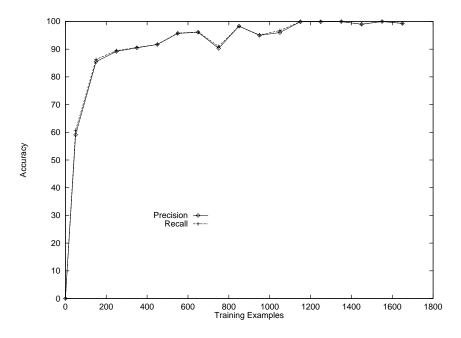7. Thanks to Jeff Siskind for the initial corpus generation software, which we enhanced for these tests.

Figure 15: Baseline Artificial Corpus

the words (the remaining 5%) had the empty meaning to simulate function words. In addition, the functors in each meaning could have a depth of up to two and an arity of up to two. An example noun-like meaning is `f23(f2(f14))`, and a verb-meaning `f10(A,f15(B))`; the conceptual symbols in this example are `f23`, `f2`, `f14`, `f10`, and `f15`. By using these multi-level meaning representations we demonstrate the learning of more complex representations than those in the geography database domain: none of the hand-built meanings for phrases in that lexicon had functors embedded in arguments. We used a grammar to generate utterances and their meanings from each original lexicon, with terminal categories selected using a distribution based on Zipf's Law (Zipf, 1949). Under Zipf's Law, the occurrence frequency of a word is inversely proportional to its ranking by occurrence.

We started with a baseline corpus generated from a lexicon of 100 words using 25 conceptual symbols and no ambiguity or synonymy; 1949 sentence/meaning pairs were generated. We split this into five training sets of 1700 sentences each. Figure 15 shows the weighted precision and recall curves for this initial test. This demonstrates good scalability to a slightly larger corpus and lexicon than that of the U.S. geography query domain.

A second corpus was generated from a second lexicon, also of 100 words using 25 conceptual symbols, but increasing the ambiguity to 1.25 meanings per word. This time, 1937 pairs were generated and the corpus split into five sets of 1700 training examples each. Weighted precision at 1650 examples drops to 65.4% from the previous level of 99.3%, and weighted recall to 58.9% from 99.3%. The full learning curve is shown in Figure 16. A quick comparison to Siskind's performance on this corpus confirmed that his system achieved comparable performance, showing that with current methods, this is close to the best performance that we are able to obtain on this more difficult corpus. One possible explanation for the smaller performance difference between the two systems on this corpus versus the geography domain is that in this domain, the correct meaning for a word is not
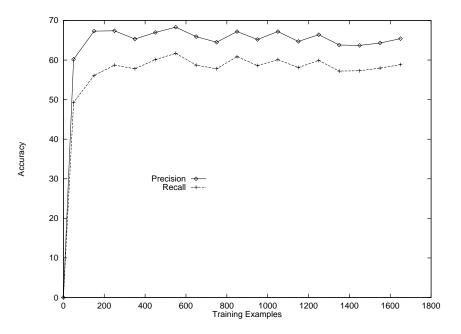
24

Figure 16: A More Ambiguous Artificial Corpus

| Ambiguity Level | Number of Examples |
|:---:|:---:|
| 1.0 | 150 |
| 1.25 | 450 |
| 2.0 | 1450 |

Table 1: Number of Examples to Reach 75% Precision

necessarily the most "general," in terms of number of nodes, of all its candidate meanings. Therefore, the generality portion of the heuristic may negatively influence the performance of WOLFIE in this domain.

Finally, we show the change in performance with increasing ambiguity and increasing synonymy. Figure 17 shows the weighted precision and recall at 1050 examples at increasing levels of ambiguity, holding the synonymy level constant. Figure 18 shows the results at increasing levels of synonymy, holding ambiguity constant. Increasing the level of synonymy does not effect the results as much as increasing the level of ambiguity, which is as we expected. Having multiple competing meanings for a word decreases the amount of evidence available for each meaning, making the learning task more difficult. On the other hand, increasing the level of synonymy does not have the potential to mislead the learner.

The number of training examples required to reach a certain level of accuracy is also informative. In Table 1, we show the point at which a standard precision of 75% was first reached for each level of ambiguity. Note, however, that we only measured accuracy after each set of 100 training examples, so the numbers in the table are approximate.

We performed a second test of scalability on two corpora generated from lexicons an order of magnitude larger than those in the above tests. In these tests, we use a lexicon
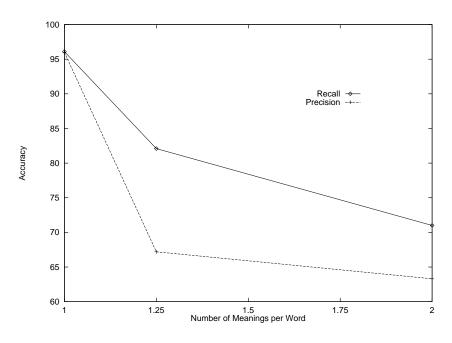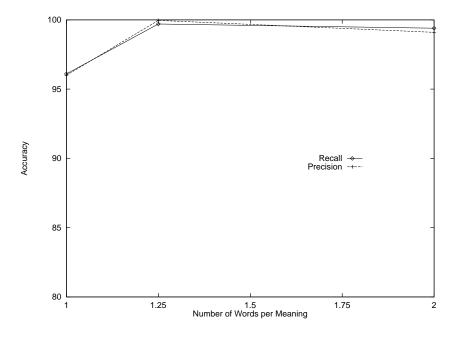
Figure 17: Increasing the Level of Ambiguity



Figure 18: Increasing the Level of Synonymy

containing 1000 words and using 250 conceptual symbols. We give results for a corpus with no ambiguity and a corpus generated from a lexicon with ambiguity and synonymy similar to that found in the WordNet database (Beckwith, Fellbaum, Gross, & Miller, 1991); the ambiguity there is approximately 1.68 meanings per word and the synonymy 1.3 words per meaning. These corpora contained 9904 (no ambiguity) and 9948 examples, respectively, and we split the data into five sets of 9000 training examples each. For the easier large corpus, the maximum average of weighted precision and recall was 85.6%, at 8100 training examples, while for the harder corpus, the maximum average was 63.1% at 8600 training examples.

## 6. Active Learning

As indicated in the previous sections, we have built an integrated system for language acquisition that is flexible and useful. However, a major difficulty remains: the construction of training corpora. Though annotating sentences is still arguably less work than building an entire system by hand, the annotation task is also time-consuming and error-prone. Further, the training pairs often contain redundant information. We would like to minimize the amount of annotation required while still maintaining good generalization accuracy.

To do this, we turned to methods in *active learning*. Active learning is a research area in machine learning that features systems that automatically select the most informative examples for annotation and training (Angluin, 1988; Seung, Opper, & Sompolinsky, 1992), rather than relying on a benevolent teacher or random sampling. The primary goal of active learning is to reduce the number of examples that the system is trained on, while maintaining the accuracy of the acquired information. Active learning systems may construct their own examples, request certain types of examples, or determine which of a set of unsupervised examples are most usefully labeled. The last approach, *selective sampling* (Cohn et al., 1994), is particularly attractive in natural-language learning, since there is an abundance of text, and we would like to annotate only the most informative sentences. For many language learning tasks, annotation is particularly time-consuming since it requires specifying a complex output rather than just a category label, so reducing the number of training examples required can greatly increase the utility of learning.

In this section, we explore the use of active learning, specifically selective sampling, for lexicon acquisition, and demonstrate that with active learning, fewer examples are required to achieve the same accuracy obtained by training on randomly chosen examples.

The basic algorithm for selective sampling is relatively simple. Learning begins with a small pool of annotated examples and a large pool of unannotated examples, and the learner attempts to choose the most informative additional examples for annotation. Existing work in the area has emphasized two approaches, *certainty-based* methods (Lewis & Catlett, 1994), and *committee-based* methods (McCallum & Nigam, 1998; Freund, Seung, Shamir, & Tishby, 1997; Liere & Tadepalli, 1997; Dagan & Engelson, 1995; Cohn et al., 1994); we focus here on the former.

In the certainty-based paradigm, a system is trained on a small number of annotated examples to learn an initial classifier. Next, the system examines unannotated examples, and attaches certainties to the predicted annotation of those examples. The $k$ examples with the lowest certainties are then presented to the user for annotation and retraining. Many

Apply the learner to $n$ bootstrap examples, creating a classifier.
Until no examples remain or the annotator is unwilling to label more examples, do:
    Use most recently learned classifier to annotate each unlabeled instance.
    Find the $k$ instances with the lowest annotation certainty.
    Annotate these instances.
    Train the learner on the bootstrap examples and all examples annotated so far.

Figure 19: Selective Sampling Algorithm

methods for attaching certainties have been used, but they typically attempt to estimate the probability that a classifier consistent with the prior training data will classify a new example correctly.

Figure 19 presents abstract pseudocode for certainty-based selective sampling. In an ideal situation, the batch size, $k$, would be set to one to make the most intelligent decisions in future choices, but for efficiency reasons in retraining batch learning algorithms, it is frequently set higher. Results on a number of classification tasks have demonstrated that this general approach is effective in reducing the need for labeled examples (see citations above).

Applying certainty-based sample selection to WOLFIE requires determining the certainty of a complete annotation of a potential new training example, despite the fact that individual learned lexical entries and parsing operators perform only part of the overall annotation task. Therefore, our general approach is to compute certainties for pieces of an example, in our case, phrases, and combine these to obtain an overall certainty for an example. Since lexicon entries contain no explicit uncertainty parameters, we used WOLFIE's heuristic measure to estimate uncertainty.

To choose the sentences to be annotated in each round, we first bootstrapped an initial lexicon from a small corpus, keeping track of the heuristic values of the learned items. Then, for each unannotated sentence, we took an average of the heuristic values of the lexicon entries learned for phrases in that sentence, giving a value of zero to unknown words, and eliminating from consideration any words that we assume are known in advance, such as database constants. The sentences with the lowest values were chosen for annotation, added to the bootstrap corpus, and a new lexicon learned. Our technique is summarized in Figure 20.

To evaluate our technique, we compared active learning to learning from randomly selected examples, again measuring the effectiveness of learned lexicons as background knowledge for CHILL. We again used the (smaller) U.S. Geography corpus, as in the original WOLFIE tests, using the lexicons as background knowledge during parser acquisition (and using the *same* examples for parser acquisition).

For each trial in the following experiments, we first randomly divide the data into a training and test set. Then, $n = 25$ bootstrap examples are randomly selected from the training examples and in each step of active learning, the least certain $k = 10$ examples of the remaining training examples are selected and added to the training set. The result

Learn a lexicon with the examples annotated so far
1) For each phrase in an unannotated sentence:
    If it has entries in the learned lexicon,
        then its certainty is the average of the heuristic values of those entries
    Else, if it is a one-word phrase
        then its certainty is zero
2) To rank sentences use:
    $$\frac{\text{Total certainty of phrases from step 1}}{\text{\# of phrases counted in step 1}}$$
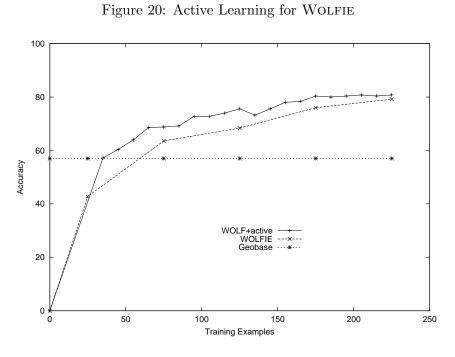
Figure 20: Active Learning for WOLFIE



Figure 21: Using Lexicon Certainty for Active Learning

of learning on this set is evaluated after each step. The accuracy of the resulting learned parsers was compared to the accuracy of those learned using randomly chosen examples to learn lexicons and parsers, as in Section 5. In other words, we can think of the $k$ examples in each round as being chosen randomly.

Figure 21 shows the accuracy on unseen data of parsers learned using the lexicons learned by WOLFIE when examples are chosen randomly and actively. There is an annotation savings of around 50 examples by using active learning: the maximum accuracy is reached after 175 examples, versus 225 with random examples. The advantage of using active learning is clear from the beginning, though the differences between the two curves are only statistically significant at 175 training examples. Since we are learning both lexicons and parsers, but only choosing examples based on WOLFIE's certainty measures, the boost could

be improved even further if CHILL had a say in the examples chosen. See Thompson, Califf, and Mooney (1999) for a description of active learning for CHILL.

## 7. Related Work

### 7.1 Lexicon Acquisition

Work on automated lexicon and language acquisition dates back to Siklossy (1972), who demonstrated a system that learned transformation patterns from logic back to natural language. More recently, Pedersen and Chen (1995) describe a method for acquiring syntactic and semantic features of an unknown word, assuming access to an initial concept hierarchy, but they give no experimental results. Many systems (Fukumoto & Tsujii, 1995; Haruno, 1995; Johnston, Boguraev, & Pustejovsky, 1995; Webster & Marcus, 1995) focus only on acquisition of verbs or nouns, rather than all types of words. Also, these either do not experimentally evaluate their systems, or do not show the usefulness of the learned lexicons for a specific application.

Several authors (e.g., (Rooth, Riezler, Prescher, Carroll, & Beil, 1999; Collins, 1997; Ribas, 1994; Manning, 1993; Resnik, 1993; Brent, 1991)) discuss the acquisition of sub-categorization information for verbs, which is different from the information required for mapping to semantic representation, but could be useful as a source of information to further constrain the search. Li (1998) further expands on this by inducing clustering information. In a related vein, several authors (Collins & Singer, 1999; Riloff & Jones, 1999; Roark & Charniak, 1998; Schneider, 1998) induce semantic categories or labels for words, and in the latter case, additional relational information. The result is typically applied as a semantic lexicon for information extraction or entity tagging.

Several systems (Knight, 1996; Hastings, 1996; Russell, 1993) learn new words from context, assuming that a large initial lexicon and parsing system are already available. Tishby and Gorin (1994) learn associations between words and actions (as meanings of those words). They test their system on a corpus of sentences paired with representations but do not demonstrate the integration of learning a semantic parser using the learned lexicon. Similarly, Oates, Eyler-Walker, and Cohen (1999) discuss the acquisition of lexical hierarchies and their associated meaning as defined by the sensory environment of a robot.

Of course, the closest work is that of Jeff Siskind, which we described in Section 2 and whose system we ran comparisons to in Section 5. His approach is somewhat more general in that it handles noise and referential uncertainty (multiple possible meanings for a sentence), while ours is specialized for applications where a single meaning is available. The experimental results in Section 5 demonstrate the advantage of our method for such an application. His system does not currently handle multiple-word phrases. Also, his system operates in an incremental or on-line fashion, discarding each sentence as it processes it, while ours is batch. While he argues for psychological plausibility, we do not. In addition, his search for word meanings is most analogous to a version space search, while ours is a greedy search. Finally, his system does not compute statistical correlations between words and their possible meanings, while ours does. Another example of pursuing language learning from a cognitive perspective is found in Colunga and Gasser (1998), who use neural network modeling techniques for learning spatial concepts.

Another related goal is pursued by De Marcken (1994), who uses a flat list of tokens for semantic representations, but does not segment sentences into words. He uses a variant of expectation-maximization (Dempster, Laird, & Rubin, 1977), together with a form of parsing and dictionary matching techniques, to segment the sentences and associate the segments with their most likely meaning. On the Childes corpus, the algorithm achieves very high precision, but recall is not provided.

Our work also has ties to the work on automatic construction of translation lexicons (Smadja, McKeown, & Hatzivassiloglou, 1996; Melamed, 1995; Wu & Xia, 1995; Kumano & Hirakawa, 1994; Catizone, Russell, & Warwick, 1993; Gale & Church, 1991; Brown & et al., 1990). While most of these methods also compute association scores between pairs (in their case, word/word pairs) and use a greedy algorithm to choose the best translation(s) for each word, they do not take advantage of the constraints between pairs. One exception is Melamed (2000); however, his approach does not allow for phrases in the lexicon or for synonymy within one text segment, while ours does. Also, Yamazaki, Pazzani, and Merz (1995) learn both translation rules and semantic hierarchies from parsed parallel sentences in Japanese and English. Of course, the main difference between this body of work and this paper is that in our case we map words to *semantic structures*, not to other words.

## 7.2 Active Learning

With respect to additional active learning techniques, Cohn et al. (1994) were among the first to discuss certainty-based active learning methods in detail. They focus on a neural network approach to actively searching a version-space of concepts. Only a few of the researchers applying machine learning to natural language processing have utilized active learning (Hwa, 2001; Schohn & Cohn, 2000; Tong & Koller, 2000; Thompson et al., 1999; Argamon-Engelson & Dagan, 1999; Liere & Tadepalli, 1997; Lewis & Catlett, 1994), and the majority of these have addressed classification tasks such as part of speech tagging and text categorization. For example, Liere and Tadepalli (1997) apply active learning with committees to the problem of text categorization. They show improvements with active learning similar to those that we obtain, but use a committee of Winnow-based learners on a traditional classification task. Argamon-Engelson and Dagan (1999) also apply committee-based learning to part-of-speech tagging. In their work, a committee of hidden Markov models is used to select examples for annotation. Lewis and Catlett (1994) use *heterogeneous* certainty-based methods, in which a simple classifier is used to select examples that are then annotated and presented to a more powerful classifier.

However, many language learning tasks require annotating natural language text with a complex output, such as a parse tree, semantic representation, or filled template. The application of active learning to tasks requiring such complex outputs has not been well studied, the exceptions being Hwa (2001), Soderland (1999), Thompson et al. (1999). The latter two include work on active learning applied to information extraction, and Thompson et al. (1999) includes work on active learning for semantic parsing. Hwa (2001) describes an interesting method for evaluating a statistical parser's uncertainty, when applied for syntactic parsing.

## 8. Future Work

Although WOLFIE's current greedy search method has performed quite well, a better search heuristic or alternative search strategy could result in improvements. A second avenue of exploration is to apply our approach to learning to parse into more common SQL database queries. Such corpora should be easily constructible by recording queries submitted to existing SQL applications along with their English forms, or translating existing lists of SQL queries into English (presumably an easier direction to translate). The fact that the same training data can be used to learn both a semantic lexicon and a parser also helps limit the overall burden of constructing a complete natural language interface.

With respect to active learning, experiments on additional corpora are needed to test the ability of our approach to reduce annotation costs in a variety of domains. It would also be interesting to explore active learning for other natural language processing problems such as syntactic parsing, word-sense disambiguation, and machine translation.

Our current results have involved a certainty-based approach; however, proponents of committee-based approaches have convincing arguments for their theoretical advantages. Our initial attempts at adapting committee-based approaches to our systems were not very successful; however, additional research on this topic is indicated. One critical problem is obtaining diverse committees that properly sample the version space (Cohn et al., 1994).

## 9. Conclusions

Acquiring a semantic lexicon from a corpus of sentences labeled with representations of their meaning is an important problem that has not been widely studied. We present both a formalism of the learning problem and a greedy algorithm to find an approximate solution to it. WOLFIE demonstrates that a fairly simple greedy symbolic learning algorithm performs well on this task and obtains performance superior to a previous lexicon acquisition system on a corpus of geography queries. Our results also demonstrate that our methods extend to a variety of natural languages besides English, and that they scale fairly well to larger, more difficult corpora.

Most experiments in corpus-based natural language have presented results on some subtask of natural language, and there are few results on whether the learned subsystems can be successfully integrated to build a complete NLP system. The experiments presented in this paper demonstrated how two learning systems, WOLFIE and CHILL, were successfully integrated to learn a complete NLP system for parsing database queries into executable logical form given only a single corpus of annotated queries.

Active learning is a new area of machine learning that has been almost exclusively applied to classification tasks. We have demonstrated its successful application to more complex natural language mappings from phrases to semantic meanings, supporting the acquisition of lexicons and parsers. The wealth of unannotated natural language data, along with the difficulty of annotating such data, make selective sampling a potentially invaluable technique for natural language learning. Our results on realistic corpora indicate that example annotations savings as high as 22% can be achieved by employing active sample selection using only simple certainty measures for predictions on unannotated data. Improved sample selection methods and applications to other important language problems

hold the promise of continued progress in using machine learning to construct effective natural language processing systems.

## 10. Acknowledgements

## References

Anderson, J. R. (1977). Induction of augmented transition networks. *Cognitive Science*, *1*, 125–157.

Angluin, D. (1988). Queries and concept learning. *Machine Learning*, *2*, 319–342.

Argamon-Engelson, S., & Dagan, I. (1999). Committee-based sample selection for probabilistic classifiers. *Journal of Artificial Intelligence Research*, *11*, 335–360.

Beckwith, R., Fellbaum, C., Gross, D., & Miller, G. (1991). WordNet: A lexical database organized on psycholinguistic principles. In Zernik, U. (Ed.), *Lexical Acquisition: Exploiting On-Line Resources to Build a Lexicon*, pp. 211–232. Lawrence Erlbaum, Hillsdale, NJ.

Borland International (1988). *Turbo Prolog 2.0 Reference Guide*. Borland International, Scotts Valley, CA.

Brent, M. (1991). Automatic acquisition of subcategorization frames from untagged text. In *Proceedings of the 29th Annual Meeting of the Association for Computational Linguistics (ACL-91)*, pp. 209–214.

Brown, P., & et al. (1990). A statistical approach to machine translation. *Computational Linguistics*, *16*(2), 79–85.

Catizone, R., Russell, G., & Warwick, S. (1993). Deriving translation data from bilingual texts. In *Proceedings of the First International Lexical Acquisition Workshop*.

Cohn, D., Atlas, L., & Ladner, R. (1994). Improving generalization with active learning. *Machine Learning*, *15*(2), 201–221.

Collins, M., & Singer, Y. (1999). Unsupervised models for named entity classification. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing and Very Large Corpora (EMNLP/VLC-99)* University of Maryland.

Collins, M. J. (1997). Three generative, lexicalised models for statistical parsing. In *Proceedings of the 35th Annual Meeting of the Association for Computational Linguistics (ACL-97)*, pp. 16–23.

Colunga, E., & Gasser, M. (1998). Linguistic relativity and word acquisition: a computational approach. In *Proceedings of the Twenty First Annual Conference of the Cognitive Science Society*, pp. 244–249.

Dagan, I., & Engelson, S. P. (1995). Committee-based sampling for training probabilistic classifiers. In *Proceedings of the Twelfth International Conference on Machine Learning (ICML-95)*, pp. 150–157 San Francisco, CA. Morgan Kaufman.

De Marcken, C. (1994). The acquisition of a lexicon from paired phoneme sequences and semantic representations. In *Lecture Notes in Computer Science*, Vol. 862, pp. 66–77.

Dempster, A., Laird, N., & Rubin, D. (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society B*, *39*, 1–38.

Fillmore, C. J. (1968). The case for case. In Bach, E., & Harms, R. T. (Eds.), *Universals in Linguistic Theory*. Holt, Reinhart and Winston, New York.

Fisher, D. H. (1987). Knowledge acquisition via incremental conceptual clustering. *Machine Learning*, *2*, 139–172.

Freund, Y., Seung, H. S., Shamir, E., & Tishby, N. (1997). Selective sampling using the query by committee algorithm. *Machine Learning*, *28*, 133–168.

Fukumoto, F., & Tsujii, J. (1995). Representation and acquisition of verbal polysemy. In *Papers from the 1995 AAAI Symposium on the Representation and Acquisition of Lexical Knowledge: Polysemy, Ambiguity, and Generativity*, pp. 39–44 Stanford, CA.

Gale, W., & Church, K. (1991). Identifying word correspondences in parallel texts. In *Proceedings of the Fourth DARPA Speech and Natural Language Workshop*.

Garey, M., & Johnson, D. (1979). *Computers and Intractability: A Guide to the Theory of NP-Completeness*. Freeman, New York, NY.

Haas, J., & Jayaraman, B. (1997). From context-free to definite-clause grammars: a type-theoretic approach. *Journal of Logic Programming*, *30*(1), 1–23.

Haruno, M. (1995). A case frame learning method for Japanese polysemous verbs. In *Papers from the 1995 AAAI Symposium on the Representation and Acquisition of Lexical Knowledge: Polysemy, Ambiguity, and Generativity*, pp. 45–50 Stanford, CA.

Hastings, P. (1996). Implications of an automatic lexical acquisition mechanism. In Wermter, S., Riloff, E., & Scheler, C. (Eds.), *Connectionist, Statistical, and Symbolic Approaches to Learning for natural language processing*. Springer-Verlag, Berlin.

Hwa, R. (2001). On minimizing training corpus for parser acquisition. In *Proceedings of the Fifth Computational Natural Language Learning Workshop*.

Jackendoff, R. (1990). *Semantic Structures*. The MIT Press, Cambridge, MA.

Johnston, M., Boguraev, B., & Pustejovsky, J. (1995). The acquisition and interpretation of complex nominals. In *Papers from the 1995 AAAI Symposium on the Representation and Acquisition of Lexical Knowledge: Polysemy, Ambiguity, and Generativity*, pp. 69–74 Stanford, CA.

Knight, K. (1996). Learning word meanings by instruction. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence (AAAI-96)*, pp. 447–454 Portland, Or.

Kohavi, R., & John, G. (1995). Automatic parameter selection by minimizing estimated error. In *Proceedings of the Twelfth International Conference on Machine Learning (ICML-95)*, pp. 304–312 Tahoe City, CA.

Kumano, A., & Hirakawa, H. (1994). Building an MT dictionary from parallel texts based on linguistic and statistical information. In *Proceedings of the Fifteenth International Conference on Computational Linguistics*, pp. 76–81.

Lavrač, N., & Džeroski, S. (1994). *Inductive Logic Programming: Techniques and Applications*. Ellis Horwood.

Lewis, D. D., & Catlett, J. (1994). Heterogeneous uncertainty sampling for supervised learning. In *Proceedings of the Eleventh International Conference on Machine Learning (ICML-94)*, pp. 148–156 San Francisco, CA. Morgan Kaufman.

Li, H. (1998). *A probabilistic approach to lexical semantic knowledge acquisition and structural disambiguation*. Ph.D. thesis, University of Tokyo.

Liere, R., & Tadepalli, P. (1997). Active learning with committees for text categorization. In *Proceedings of the Fourteenth National Conference on Artificial Intelligence (AAAI-97)*, pp. 591–596 Providence, RI.

Manning, C. D. (1993). Automatic acquisition of a large subcategorization dictionary from corpora. In *Proceedings of the 31st Annual Meeting of the Association for Computational Linguistics (ACL-93)*, pp. 235–242 Columbus, OH.

McCallum, A. K., & Nigam, K. (1998). Employing EM and pool-based active learning for text classification. In *Proceedings of the Fifteenth International Conference on Machine Learning (ICML-98)*, pp. 350–358 Madison, WI. Morgan Kaufman.

Melamed, I. D. (2000). Models of translational equivalence among words. *Computational Linguistics*, *26*(2), 221–249.

Melamed, I. (1995). Automatic evaluation and uniform filter cascades for inducing $n$-best translation lexicons. In *Proceedings of the Third Workshop on Very Large Corpora*.

Muggleton, S., & Feng, C. (1990). Efficient induction of logic programs. In *Proceedings of the First Conference on Algorithmic Learning Theory* Tokyo, Japan. Ohmsha.

Muggleton, S. H. (Ed.). (1992). *Inductive Logic Programming*. Academic Press, New York, NY.

Oates, T., Eyler-Walker, Z., & Cohen, P. (1999). Using syntax to learn semantics: an experiment in language acquisition with a mobile robot. Tech. rep. 99-35, University of Massachusetts Computer Science Department.

Pedersen, T., & Chen, W. (1995). Lexical acquisition via constraint solving. In *Papers from the 1995 AAAI Symposium on the Representation and Acquisition of Lexical Knowledge: Polysemy, Ambiguity, and Generativity*, pp. 118–122 Stanford, CA.

Plotkin, G. D. (1970). A note on inductive generalization. In Meltzer, B., & Michie, D. (Eds.), *Machine Intelligence (Vol. 5)*. Elsevier North-Holland, New York.

Resnik, P. (1993). *Selection and information: a class-based approach to lexical relationships.* Ph.D. thesis, University of Pennsylvania, CIS Department.

Ribas, F. (1994). An experiment on learning appropriate selectional restrictions from a parsed corpus. In *Proceedings of the Fifteenth International Conference on Computational Linguistics*, pp. 769–774.

Riloff, E., & 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. 1044–1049 Orlando, FL.

Roark, B., & Charniak, E. (1998). Noun-phrase co-occurrence statistics for semi-automatic semantic lexicon construction. In *Proceedings of the 36th Annual Meeting of the Association for Computational Linguistics and COLING-98 (ACL/COLING-98)*, pp. 1110–1116.

Rooth, M., Riezler, S., Prescher, D., Carroll, G., & Beil, F. (1999). Inducing a semantically annotated lexicon via EM-based clustering. In *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics*, pp. 104–111.

Russell, D. (1993). *Language Acquisition in a Unification-Based Grammar Processing System Using a Real World Knowledge Base.* Ph.D. thesis, University of Illinois, Urbana, IL.

Schank, R. C. (1975). *Conceptual Information Processing.* North-Holland, Oxford.

Schneider, R. (1998). A lexically-intensive algorithm for domain-specific knowledge acquisition. In *Proceedings of the Joint Conference on New Methods in Language Processing and Computational Natural Language Learning*, pp. 19–28.

Schohn, G., & Cohn, D. (2000). Less is more: Active learning with support vector machines. In *Proceedings of the Seventeenth International Conference on Machine Learning (ICML-2000)*, pp. 839–846 Stanford, CA.

Seung, H. S., Opper, M., & Sompolinsky, H. (1992). Query by committee. In *Proceedings of the ACM Workshop on Computational Learning Theory* Pittsburgh, PA.

Siklossy, L. (1972). Natural language learning by computer. In Simon, H. A., & Siklossy, L. (Eds.), *Representation and meaning: Experiments with Information Processsing Systems.* Prentice Hall, Englewood Cliffs, NJ.

Siskind, J. M. (1992). *Naive Physics, Event Perception, Lexical Semantics and Language Acquisition*. Ph.D. thesis, Department of Electrical Engineering and Computer Science, Massachusetts Institute of Technology, Cambridge, MA.

Siskind, J. M. (1996). A computational study of cross-situational techniques for learning word-to-meaning mappings. *Cognition, 61*(1), 39–91.

Smadja, F., McKeown, K. R., & Hatzivassiloglou, V. (1996). Translating collocations for bilingual lexicons: A statistical approach. *Computational Linguistics, 22*(1), 1–38.

Soderland, S. (1999). Learning information extraction rules for semi-structured and free text. *Machine Learning, 34*, 233–272.

Thompson, C. A., Califf, M. E., & Mooney, R. J. (1999). Active learning for natural language parsing and information extraction. In *Proceedings of the Sixteenth International Conference on Machine Learning (ICML-99)*, pp. 406–414 Bled, Slovenia.

Thompson, C. A. (1995). Acquisition of a lexicon from semantic representations of sentences. In *Proceedings of the 33rd Annual Meeting of the Association for Computational Linguistics (ACL-95)*, pp. 335–337 Cambridge, MA.

Tishby, N., & Gorin, A. (1994). Algebraic learning of statistical associations for language acquisition. *Computer Speech and Language, 8*, 51–78.

Tomita, M. (1986). *Efficient Parsing for Natural Language*. Kluwer Academic Publishers, Boston.

Tong, S., & Koller, D. (2000). Support vector machine active learning with applications to text classification. In *Proceedings of the Seventeenth International Conference on Machine Learning (ICML-2000)*, pp. 999–1006 Stanford, CA.

Webster, M., & Marcus, M. (1995). Automatic acquisition of the lexical semantics of verbs from sentence frames. In *Proceedings of the 27th Annual Meeting of the Association for Computational Linguistics (ACL-89)*, pp. 177–184.

Wu, D., & Xia, X. (1995). Large-scale automatic extraction of an English-Chinese translation lexicon. *Machine Translation, 9*(3-4), 285–313.

Yamazaki, T., Pazzani, M., & Merz, C. (1995). Learning hierarchies from ambiguous natural language data. In *Proceedings of the Twelfth International Conference on Machine Learning (ICML-95)*, pp. 575–583 San Francisco, CA. Morgan Kaufmann.

Zelle, J. M. (1995). *Using Inductive Logic Programming to Automate the Construction of Natural Language Parsers*. Ph.D. thesis, Department of Computer Sciences, University of Texas, Austin, TX. Also appears as Artificial Intelligence Laboratory Technical Report AI 96-249.

Zelle, J. M., & Mooney, R. J. (1996). Learning to parse database queries using inductive logic programming. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence (AAAI-96)*, pp. 1050–1055 Portland, OR.

Zipf, G. (1949). *Human behavior and the principle of least effort.* Addison-Wesley, New York, NY.