

Lexical Discovery with an Enriched Semantic Network

Doug Beeferman
School of Computer Science
Carnegie Mellon University

March 12, 1998

Abstract

The study of lexical semantics has produced a systematic analysis of relationships between content words that has greatly benefited both lexical search tools and natural language processing systems. We describe research toward a common algorithmic core for these two applications. We first introduce a database system called FreeNet that facilitates the description and exploration finite binary relations. We then describe the design and implementation of Lexical FreeNet, a semantic network that mixes WordNet-derived semantic relations with data-derived and phonetically-derived relations. We discuss how Lexical FreeNet has aided in lexical discovery, the pursuit of linguistic and factual knowledge by the computer-aided exploration of lexical relations, and relate this to ongoing research in how the system can be employed in natural language processing algorithms for segmentation and summarization.

Submission Type: Paper for Coling-ACL '98 Workshop, "Usage of WordNet in Natural Language Processing Systems"

Topic Areas: Lexical semantics, semantic networks, text processing, WordNet

Author of Record: Doug Beeferman

Under consideration for other conferences (specify)? None

Lexical Discovery with an Enriched Semantic Network

Abstract

The study of lexical semantics has produced a systematic analysis of relationships between content words that has greatly benefited both lexical search tools and natural language processing systems. We describe research toward a common algorithmic core for these two applications. We first introduce a database system called FreeNet that facilitates the description and exploration finite binary relations. We then describe the design and implementation of Lexical FreeNet, a semantic network that mixes WordNet-derived semantic relations with data-derived and phonetically-derived relations. We discuss how Lexical FreeNet has aided in lexical discovery, the pursuit of linguistic and factual knowledge by the computer-aided exploration of lexical relations, and relate this to ongoing research in how the system can be employed in natural language processing algorithms for segmentation and summarization.

1 Motivation

Numerous statistical natural language processing systems treat content words as information-carrying units acting in combination to express the meaning of a discourse. Most information retrieval systems use vector-space models of meaning, in which a document is represented in the high-dimensional space defined by the vocabulary. Large vocabulary speech recognizers, summarizers, and topic segmenters, too, are highly lexical, focusing on the meanings of individual words in favor of parsing and trying to infer the meanings of whole sentences, function words and all.

The study of lexical semantics has produced a systematic analysis of relationships between content words that has greatly benefited these systems. WordNet (Miller, 1990a), a semantic network originally designed to be a tool for lexical discovery (Miller, 1985), now plays a role in a number of NLP algorithms (Rosenzweig, 1998). As an electronic dictionary it is an exceptional resource, but the semantic connections it exposes to the human user are at least as valuable to a well-informed program, be it a word sense disambiguator or a search engine.

This paper introduces Lexical FreeNet, a database system leveraging WordNet and other knowledge and data sources to provide a common algorithmic core for lexical-semantic approaches to knowledge discovery by people and text processing by machines. We view these two applications as intimately related, with the semantic network serving as a common data structure.

Semantic networks

Semantic networks are widely used and occupy a firm niche in artificial intelligence research. Such a network views a knowledge space (lexical knowledge, for instance) as a graph of nodes that represent concepts (lexemes, for instance), with directed edges between nodes that are related in some way (generalization, for instance).

Fundamentally, a semantic network is a set of relations over the same vocabulary, or a multigraph with typed edges. Many networks familiar to us from such whimsical applications as computing the “Kevin Bacon number” of an actor or the “Erdős number” of a researcher employ only one relation. The genericity of the multigraph has far broader implications for linguistic knowledge, however. Our goal shall be to tease out the salient graph-theoretic properties that might make these networks more useful to clients such as a segmentation program, summarizer, search engine, or human being.

Lexical Discovery

By *lexical discovery* we mean the pursuit of linguistic and factual knowledge by the computer-aided exploration of lexical relations. A semantic network allied with the proper user interface can be a useful tool in its own right. By organizing words semantically rather than alphabetically, WordNet provides a means by which users can navigate its vocabulary logically, establishing nontrivial connections between concepts and not simply character sequences. Exploring the WordNet *hyponym* tree starting at the word **mammal**, for instance, reveals to us that aardvarks are mammals; exploring WordNet’s *meronymy* relation at the word **mammal** reveals to us that mammals have **hair**. From these two explorations we can accurately conclude that aardvarks have hair.

Lexical exploration need not be limited to one step at a time, however. Viewing a semantic network as a computational structure awaiting graph-theoretic queries gives us the freedom to demand services beyond mere lookup. “Does the aardvark have hair?”, or “What is the closest connection between aardvarks and hair?” or “How interchangeably can the words **aardvark** and **anteater** be used?” are all reasonable questions with answers staring us in the face. Of course, the idea of finding shortest paths in semantic networks (through so-called *activation-spreading* or *intersection search*) is as at least as old as Dijkstra himself. But these questions have typ-

ically been asked of very limited graphs, networks for domains far narrower than the lexical space of English, say. We feel that formalizing how WordNet can be employed for this broader sort of lexical discovery is a good start. But we also feel that it is necessary first to enrich the network with information that, as we shall see, cannot be easily gleaned from WordNet's current battery of relations. The very large electronic corpora and wide variety of linguistic resources that today's computing technology has enabled will in turn enable this.

Text processing

Lexical discovery need not be limited to people, either. Several sense tagging algorithms traverse WordNet in some fashion. Other semantic networks are explored by algorithms for text segmentation (Kozima and Furugori, 1993), summarization (Mani and Bloedorn, 1997) and information extraction (Niwa et al., 1997). These techniques base decisions not only on the knowledge that two observed words are related, but also on *how* they are related. This relationship information is harder to extract by purely statistical means from text corpora because of the inherent limitations in the quantity and conceptual coverage of available data. Knowledge resources such as WordNet offer coverage of content words that is guaranteed to be "comprehensive enough" conceptually by the standards of the lexicographers who create them.

But where conceptual coverage is a strength, topicality is a weakness of knowledge-engineered relations. No fixed resource can offer an up-to-date set proper names, for example, without significant, ongoing human effort. Also, many relations do not fall into categories that lexicographers view as important enough to code painstakingly into their dictionaries. For example, the **anteater** (but not the **aardvark**) is the mascot of the University of California at Irvine, but you wouldn't know it from WordNet, and justifiably so. Yet a statistical language model recognizes the existence of some relationship between the two concepts.

Enriching a knowledge-engineered semantic network with data-derived relations offers the chance to mix these complementary resources in natural language processing applications, overcoming data sparseness on one hand, and conceptual coverage on the other.

Outline

The remainder of this paper is organized as follows. We shall first describe in Section 2 the FreeNet database system for the expression and analysis of relational data. In Section 3 we'll describe the design and construction of an instance of this database called Lexical FreeNet. We'll conclude by providing examples of applications of Lexical FreeNet to lexical discovery, and then briefly discuss intended ap-

plications of the system to the semantic processing of text.

2 FreeNet

FreeNet, an acronym for *finite relation expression network*, is a system for describing and exploring finite binary relations. Here we mean relation in the mathematical sense, *i.e.* a set of ordered pairs. We concern ourselves with finite sets of pairs of tokens drawn from a finite set of tokens, or *vocabulary*.

2.1 Tokens and relations

A *token* in FreeNet is simply a normalized string of characters drawn from a finite vocabulary. The vocabulary might be a dictionary of English words, a set of movie titles, or a set of names of researchers. The system is assumed to implement normalization as a function from input strings to strings.

A *relation* in FreeNet is a finite set of ordered pairs of tokens, or *links*. Each relation has a name that, like a token, is simply a normalized string of characters drawn from a finite vocabulary (which we shall do better to call an *alphabet*, for reasons made clear below.)

Use of the FreeNet system can be seen to consist of three distinct processing phases: the *relation computation* stage, in which a set of relations is derived from some knowledge or data source and transduced to an explicit set of labeled ordered pairs; the *graph construction* stage, in which this set of labeled pairs is transduced to an efficient multigraph representation; and the *query* stage, in which a client can interact with the system to find paths in the multigraph that match a certain specification.

FreeNet consists of software to do the second and third phases. Implementation of a specific instance of FreeNet requires the user to write software to do the first phase, but support software exists for an optional *filtering* substage that constrains the input pair set in certain ways—eliminating pairs that contain stopwords, enforcing limits on the fanout of tokens, and enforcing strength thresholds, for instance.

The second phase, graph building, simply entails providing a set of triples (two tokens and a relation) to the system. The order in which the triples appear in the input does not matter, as it is the onus of FreeNet to reorder the links as necessary and store the graph efficiently.

The third phase, querying, is the chief novel contribution, and is described below.

2.2 Regular expressions

The power behind FreeNet lies in the client's ability to compose primitive relations to build more complex relations that it may use in its queries.

The primary mechanism for building complex relations is the *regular expression* over the alphabet of relation names. Just as a regular expression over

ASCII characters specifies a regular set of strings recursively in terms of other sets, so too can a regular expression over relation names specify a set of ordered pairs recursively in terms of other sets and various operators.

The following grammar specifies allowable regular expressions in FreeNet.

```

regexp =
    <rel>                (relation name)
  | (regexp)             (parenthesization)

    (Binary operators)
  | regexp regexp       (concatenation)
  | regexp | regexp     (union)
  | regexp & regexp     (conjunction)

    (Unary operators)
  | regexp*             (transitive closure)
  | regexp'             (inverse)
  | regexp-             (complement)
  | regexp%             (sibling)

```

These regexp-building operators are described below.

Concatenation

The concatenation operator is used to compose two relations directly. The expression $r_1 r_2$ denotes the set of pairs (a, b) such that for some token c , $(a, c) \in r_1$ and $(c, b) \in r_2$. For example, a network implementing a genealogy database might offer primitive **parent** and **brother** relations. In that case, the relation denoted by the regular expression **(parent brother)** is what we know of as the **uncle** relation.

Conjunction

Conjunction takes the intersection of two relations: plainly, the intersection of their respective pair sets. The expression $r_1 \& r_2$ denotes the set of pairs (a, b) such that $(a, b) \in r_1$ and $(a, b) \in r_2$.

Supposing that in a lexical semantic net we have the relations **required_by** and **requires**, then a symmetric **sympiotic_with** relation might be implemented as their conjunction.

Union

The union operator is used to join two relations. The expression $r_1 | r_2$ denotes the set of pairs (a, b) such that $(a, b) \in r_1$ or $(a, b) \in r_2$. In an Erdős-number like application, for example, two authors may be “related” if they have coauthored a paper *or* if one has cited the other.

Transitive closure

We commonly reason about the transitive closure of relations. The transitive closure operator implements homogeneous reachability—is there a path between the tokens using links only of a certain type? Namely, let r^1 denote the relation r and r^i

$i > 1$ denote the relation $(r \ r^{(i-1)})$. Then r^* denotes the union of all r^i as i ranges from 0 to infinity. (Note that since we assume finite relations, this set is always finite.) In the genealogy example, **parent*** would be what we consider the “ancestor” relation.

Inverse, Complement, and Sibling

A few more unary operators are minor conveniences in building relations. The inverse operator swaps every pair: r^- denotes the set of pairs (a, b) such that $(b, a) \in r$. Taking the union of a relation with its inverse produces a new relation that is guaranteed to be symmetric.

The complement operator produces a set containing all pairs *but* those in a certain relation. r' denotes the set of pairs (a, b) such that $(b, a) \notin r$. (The vocabulary is assumed to be fixed after the graph is built, and so the universe is well-defined.)

The sibling operator produces pairs that have in common their relation with a certain other token. $r\%$ denotes the set of pairs (a, b) such that $a \neq b$ and there exists a c such that $(a, c) \in r$ and $(b, c) \in r$. Thus **(parent-)%** relation is the genealogical sibling relation formed by applying the inverse operator and then the sibling operator to the “parent” relation. (Note that the sibling operator cannot quite be implemented in terms of the inverse and conjunction or concatenation operators, because the “sibling” of a relation excludes reflexivity.)

Notes

A simple structural induction can be used to prove that any relation built from these operators is also a relation. Additional operators to support set addition and subtraction with constant sets are also conceivable.

2.3 Queries

Queries in FreeNet specify a set of paths in the graph that are to be searched for. A *path specification* is a sequence of tokens or token variables with interleaved relation regexps. More precisely, every query is of the form $(W \langle \text{regexp} \rangle)^* W$, where W is either a *constant token* or a *variable* w_i , and $\langle \text{regexp} \rangle$ is a regular expression over relations, as defined above.

FreeNet returns a shortest path (or all paths) in the multigraph that match the query, binding the variables in the query to concrete tokens. The output includes the names of all of the primitive relation links traversed.

Queries in the public, limited implementation of FreeNet can take one of only four forms, each parameterized by one or two tokens; but these demonstrate what are expected to be common queries. Below, the “ANY” regexp is the union of all available (or selected) primitive relations. The comma (“,”) represents the universal relation, linking all pairs

of tokens; the comma relation can thus be used in FreeNet queries to implement conjunction of clauses.

- *Shortest path*: This query takes two arguments s and t , and outputs the result of the query “ s ANY* t ”. This finds a shortest path, using any of the selected relations, between the source and the target.
- *Fanout*: This query takes a single argument s and outputs the result of “ s ANY w_1 ”. This simply shows all words related in some way to the source.
- *Intersection search*: This query takes two arguments s and t and outputs the result of “ s ANY w_1 , t ANY w_1 ”. This is useful for finding what two tokens “have in common” in terms of primitive relationships with other tokens. The two relations involved in such a path need not be identical.
- *Coercion*: This query takes two arguments s and t , two relations $re1$ and $re2$, and outputs the result of “ s $re1$ w_1 $re2$ w_2 $re1$ t ”. This is useful for a wide variety of constraint-solving, such as, in the lexical semantic net case, pun and rhyme generation.

2.4 Implementation issues

A FreeNet multigraph is stored sparsely for efficient offline (disk) access as a list of variable-length adjacency lists. Each element in an adjacency list is a single 32-bit word that describes an arc by combining its destination token ID and relation ID; the source token ID for an arc is implicit in its row. An index of offsets into the list is precomputed and stored together with hash tables for the token and relation namespaces. At no point in query processing is more than a single line of the list (equivalently, a set of links emanating from the same source node) in memory at once.

Graph construction

A number of optimizations in the layout of the multigraph on disk are essential if arbitrary searches over large multigraphs are to be efficient. Of particular concern is disk seek time, because traversing the graph entails accessing different rows of the adjacency list representation in rapid succession. One simple preprocessing step is to sort each row of the representation by the word identifier’s row location, so that all of the nodes emanating from a fixed source can be accessed with a unidirectional sweep.

A trickier concern is the ordering of the rows themselves. We desire to order the rows so that related words tend to appear near each other so that seek time between them is minimized. We can formalize this problem by asking for an ordering that minimizes the average offset difference between a

randomly chosen edge in the multigraph. This problem is at least as computationally hard as the well-studied, NP-complete *bandwidth* problem in graph theory (Papadimitriou, 1976), which is to find a linear ordering of the vertices of a given graph such that the maximum difference in the ordering between any two adjacent vertices is minimal. We are studying approximation algorithms (Blum et al., to appear) that allow this preprocessing step to be carried out efficiently during database construction.

Querying

Supporting arbitrary FreeNet queries, with the full battery of regular expression operators, is a non-trivial data structures problem, because it is prohibitively expensive to add new links with the occurrence of a new regexp. Instead, the graph is static. Each relation in the “alphabet” of relations is converted to an ASCII character, and stock regexp processing software is used to convert each regexp in a query to a state machine. A query is converted to a single state machine by concatenating its constituent regexp state machines, interleaving “constraint points” that enforce the identity of multiple bindings of the same variable. A dynamic set of state IDs and backtrace IDs is associated with each token to support breadth-first search.

The limited query types above are implemented without all this machinery, by simply performing breadth-first-search on the graph, maintaining a single backtrace ID for each node, and allowing or prohibiting certain relations as specified by the user. Coercion is implemented as a hard-coded path constraint.

3 Lexical FreeNet

Lexical FreeNet is an instance of FreeNet supporting a range of lexical semantic applications. It achieves this by mixing statistically-derived and knowledge-derived relations.

Tokens

The tokens in Lexical FreeNet are the words that appear in at least one of the program’s various data sources. This includes over 130,000 words from the CMU Pronouncing Dictionary v1.6d (CMU, 1997), 160,000 words and multiple-word phrases from WordNet 1.6, and 60,000 words from the broadcast news transcripts used to train the trigger relation. The intersection between these three sources is significant, of course, and in total there are slightly under 200,000 distinct tokens, including phrases.

Relations

Lexical FreeNet includes seven *semantic* relations, two *phonetic* relations, and one *orthographic* relation. These relations connect the token set with about seven million links, costing 30 MB of disk

space. A summary of the relations is shown in Figure 2. Below we use a bidirectional arrow (\Leftrightarrow) to indicate a symmetric relation, and a unidirectional arrow (\Rightarrow) to indicate an asymmetric relation.

“Synonym of” ($\stackrel{\text{SYN}}{\Leftrightarrow}$)

This relation is computed by taking, for each synonym set (or *synset*) in all WordNet 1.6 word categories, the cross-product of the synonym set with itself, excluding reflexive links. That is to say, we include all pairs of lexemes in each synset except the links from a lexeme to itself. Thus we mix different lexeme senses into the same soup, conflating, for example, the two noun and verb senses of **BIKE** in **bike** $\stackrel{\text{SYN}}{\Leftrightarrow}$ **bicycle** and **bike** $\stackrel{\text{SYN}}{\Leftrightarrow}$ **pedal**.

“Triggers” ($\stackrel{\text{TRG}}{\Leftrightarrow}$)

Trigger pairs are ordered word pairs that co-occur significantly in data; that is, they are pairs that appear near each other in text more frequently than would be expected if the words were unrelated. Given a large corpus of text data, we built the asymmetric trigger relation by finding the pairs in the cross-product of the vocabulary that have the highest average *mutual information*, as in (Rosenfeld, 1994; Beeferman et al., 1997a). Mutual information is one measure of whether an observed co-occurrence of two vocabulary words is not due to chance. Word pairs with high mutual information are likely to be semantically related in some way.

We chose 160 million words of Broadcast News data (LDC, 1997) for this computation, and defined co-occurrence as “occurring within 500 words”, approximately the average document length. We selected the top 350,000 trigger pairs from the ranking to use in the relation, putting the size of the relation on par with the synonym relation.¹ Some of the top trigger pairs discovered by this procedure are shown in Table 1. In our implementation we limit the number of trigger links emanating from a token to the top 50, and prune away links that include any member of a hand-coded stopword set that includes function words.

“Specializes” ($\stackrel{\text{SPC}}{\Rightarrow}$) and “Generalizes” ($\stackrel{\text{GEN}}{\Rightarrow}$)

The specialization relation captures the lexical inheritance system underlying WordNet nouns (Miller, 1990b) and verbs (Fellbaum, 1990). It is computed by taking, for each pair of WordNet synsets that appear as parent and child in the WordNet *hyponym* trees, the cross-product of the pair. For example, **shoe** $\stackrel{\text{SPC}}{\Rightarrow}$ **footwear**.

The generalization relation is simply the inverse of specialization relation, or SPC-. For example: **tree** $\stackrel{\text{GEN}}{\Rightarrow}$ **cypress**.

¹We used the Trigger Toolkit, available at <http://www.cs.cmu.edu/aberger/software.html>, for this computation

<i>s</i>	<i>t</i>
Los Angeles	United States
White House	President Clinton
New York	health care
...	...
campaign	Bush
Haitian	Aristide
films	film
fed	rates
court	evidence
care	insurance

Figure 1: The top six trigger pairs (*s*, *t*), ranked by mutual information, in the Lexical FreeNet trigger relation, and the 500th through 505th-ranked pairs. The highest-ranked pairs tend to be distance-one bigram phrases, while the remainder co-occur at greater distances.

“Part of” ($\stackrel{\text{PAR}}{\Rightarrow}$) and “Comprises” ($\stackrel{\text{COM}}{\Rightarrow}$)

The $\stackrel{\text{PAR}}{\Rightarrow}$ relation captures meronymy, another inheritance system which can informally be thought of as a “part of” tree over nouns. It is computed by taking, for each pair of WordNet synsets that are related in WordNet by the meronym relation, the cross-product of the pair. For example, **shoe** $\stackrel{\text{SPC}}{\Rightarrow}$ **footwear**. The “comprises” relation is simply its inverse, PAR-, as in **tree** $\stackrel{\text{COM}}{\Rightarrow}$ **cypress**.

“Antonym of” ($\stackrel{\text{ANT}}{\Leftrightarrow}$)

The antonym relation uses the antonym relation defined in WordNet for nouns, verbs, adjectives, and adverbs. It is computed by taking, for each pair of WordNet synsets that are related in WordNet by the antonym relation, the cross-product of the pair. For example, **clear** $\stackrel{\text{SPC}}{\Leftrightarrow}$ **opaque**.

“Phonetically similar to” ($\stackrel{\text{SIM}}{\Leftrightarrow}$) and “Rhymes with” ($\stackrel{\text{RHY}}{\Leftrightarrow}$)

To allow users to cross the dimensions of sound and meaning in their queries, two phonetic relations are added to the mix in Lexical FreeNet. These relations, while amusing for shortest path queries, are not expected to contribute to the text processing applications discussed later in this paper. Both relations leverage the phonetic and lexical stress transcriptions in the CMU Pronouncing Dictionary.

The $\stackrel{\text{SIM}}{\Leftrightarrow}$ relation is computed by adding every pair of words in the vocabulary that have pronunciations which differ in *edit distance* by at most some number of edits. Edit distance is computed using a dynamic programming algorithm as the minimum number of substitutions, insertions, and dele-

(a)

Database breakdown by relation		
Relation	Symbol	Number of links
Triggers	TRG	354800
Synonymous	SYN	249156
Generalizes	GEN	261275
Specializes	SPC	261281
Comprises	COM	23650
Part of	PAR	23650
Antonym of	ANT	18982
Rhymes	RHY	4533536
Sounds like	SIM	1483360
Anagram	ANA	91072
197598 distinct tokens		

(b)

	TRG	SYN	GEN	SPC	COM	PAR	ANT	RHY	SIM	ANA
TRG	354800									
SYN	1873	249156								
GEN	972	3390	261275							
SPC	1164	3390	1628	261281						
COM	330	1805	150	263	23650					
PAR	391	1805	263	150	310	23650				
ANT	469	38	28	28	5	5	18982			
RHY	1279	1576	1042	1042	30	30	1576	4533536		
SIM	10129	1518	150	150	15	15	92	587015	1483360	
ANA	363	1338	18	18	25	25	0	7871	9946	91072

Figure 2: Statistics on the relations in Lexical FreeNet. (a) The number of links in each relation. (b) Relation crossover counts. Each cell reports the number of word pairs that exist in both relations. One of the 5 pairs counted in the cell at (ANT, COM), for example, is (DAY, NIGHT).

tions (unweighted, and blind to nearness in substitution) to the first word’s phonetic sequence required to reach the second word’s phonetic sequence. In our current implementation we limit the relation to pairs with edit distance at most 1, e.g. **cancel** $\overset{\text{SIM}}{\longleftrightarrow}$ **candle**.

The $\overset{\text{RHY}}{\longleftrightarrow}$ relation is computed by adding each pair of words that have pronunciations such that their phonetic suffixes including and following the primary stressed syllables match, e.g. **Reno** $\overset{\text{RHY}}{\longleftrightarrow}$ **Casino**.

“Anagram of” ($\overset{\text{ANA}}{\longleftrightarrow}$)

The final relation, $\overset{\text{ANA}}{\longleftrightarrow}$, is almost, but not quite, completely useless, symmetrically linking lexemes that use the same distribution of letters, as in **Geraldine** $\overset{\text{ANA}}{\longleftrightarrow}$ **realigned**. This is perhaps best described as a “wormhole” in lexical space.

Extensions

A portion of the wealth of WordNet was discarded in Lexical FreeNet—the verb entailment relation, for instance. Adjectives are somewhat slighted by the system, as their WordNet description in terms of

bipolar attributes (Gross and Miller, 1990) is largely ignored.

Other possible semantic relations include the more specialized knowledge-engineered links that appear in typically narrow-coverage semantic nets, such as “acts on”, “uses”, “stronger than”, and the like.

Data-driven approaches to relation induction that dig deeper than the collocation extraction of the trigger computation may prove useful and interesting. One approach (Richardson, 1997; Richardson et al., 1993) bootstraps a parser to induce many unconventional semantic relations from dictionary data. A link grammar (Sleator and Temperley, 1991) applied to data can conceivably be used to extract some interesting relations that live at the syntax/semantics interface.

For fun, where *fun* need not be equated with uselessness, other *orthographic* relations can be trivially computed, such as a host of edit-distance-based relations. This can facilitate popular games like finding “word ladders” of nearly identically-spelled words.

4 Applications

In this section we discuss applications of Lexical FreeNet, first to human users, and then to natural language processing applications.

4.1 Lexical discovery

A World Wide Web interface to Lexical FreeNet, depicted in Figure 3, is available and has become a popular online resource since its release in late January, 1998.² The program allows the user to issue one of the four template queries to the database described in Section 2.3. One of these query templates (“Fanout”) requires only a single source token as input, and this has become a popular lookup tool, providing some of the functionality of a thesaurus and rhyming dictionary. The other query functions require source *and* target tokens to be specified. Each token can itself contain spaces in the case of phrasal inputs, which are normalized to the underscore character in processing. The four basic queries allow the user to specify a subset of the ten primitives relations to permit in the output paths by clicking a series of checkboxes. Upon submission, the state of the checkboxes sets the ANY relation to be the union of checked relations.

An additional “Spell check” query mode allows the user to find database tokens that have similar (or exact) spelling to a given input token, where similarity is measured by an orthographic edit distance.

Upon submission, the system finds and displays the path or paths resulting from the query with arrow glyphs representing the various relations. Queries typically finish within an acceptable time

²See <http://www.link.cs.cmu.edu/lexfn/>

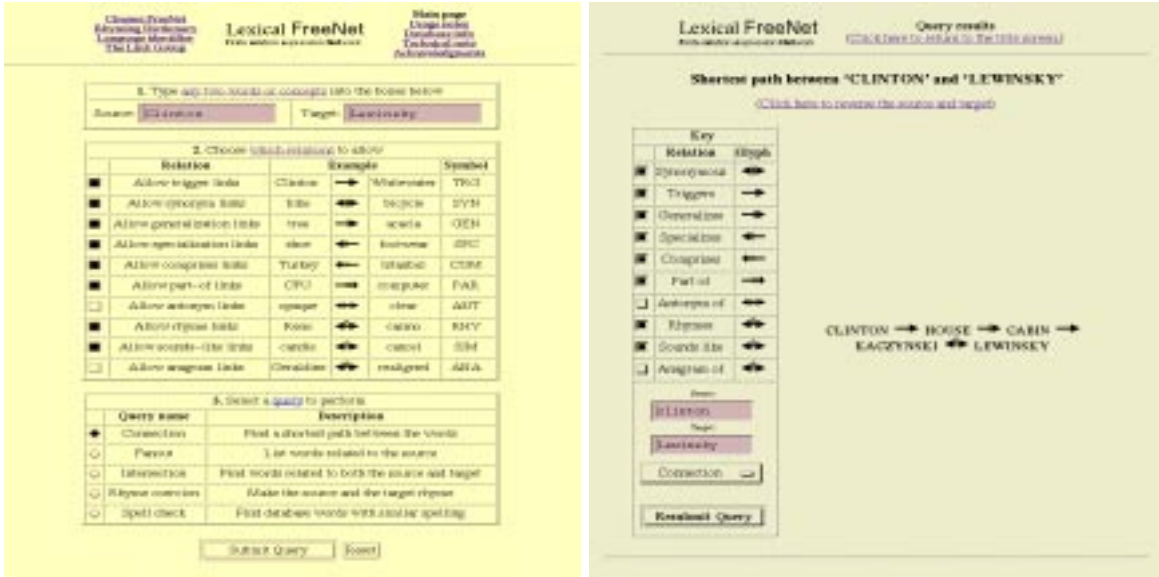


Figure 3: The front page of the Web interface to Lexical FreeNet, and the results of a shortest path query.

window of three to ten seconds. The results screen summarizes the query and allows the user to re-submit it with modifications, improving the ease of database “navigation” over having to return to the title screen.

Feedback from the Web site indicates that the system has been used as an aid in writing poetry and lyrics; devising product names; generating puzzles for elementary school language arts classes; writing greeting cards; devising insults and compliments; and, above all, just exploring. Following are selected examples of the system’s output in various configurations.

Shortest path queries

The *shortest path* query is the primary vehicle for establishing connections between words and concepts:

- Shortest path queries that allow all lexical relations can be used to aid in generating puns and quips involving the two endpoint concepts. For example, below is the shortest path between **Clinton** and **Lewinsky** using all relations:

$$\begin{array}{c} \text{CLINTON} \xrightarrow{\text{TRG}} \text{HOUSE} \xrightarrow{\text{GEN}} \text{CABIN} \xrightarrow{\text{TRG}} \\ \text{KACZYNSKI} \xrightarrow{\text{RHY}} \text{LEWINSKY} \end{array}$$

- Shortest path queries allowing only the hyponymy relations can connect any two nouns in the WordNet hyponymy tree through their least common ancestor. For example, animals can be connected taxonomically, as in the shortest path between **potto** and **langur** using only the specialization ($\overset{\text{SPC}}{\implies}$) and generalization ($\overset{\text{GEN}}{\implies}$) and relations:

$$\begin{array}{c} \text{POTTO} \xrightarrow{\text{SPC}} \text{LEMUR} \xrightarrow{\text{SPC}} \text{PRIMATE} \xrightarrow{\text{GEN}} \\ \text{MONKEY} \xrightarrow{\text{GEN}} \text{OLD_WORLD_MONKEY} \xrightarrow{\text{GEN}} \text{LANGUR} \end{array}$$

- Shortest path queries allowing only the meronymy relations can connect many noun pairs. For example, geographical connections can be made between place names to find the largest enclosing region, as in the shortest path between **Saskatoon** and **Winnipeg** using only the comprise ($\overset{\text{COM}}{\implies}$) and part-of ($\overset{\text{PAR}}{\implies}$) relations:

$$\begin{array}{c} \text{SASKATOON} \xrightarrow{\text{PAR}} \text{SASKATCHEWAN} \xrightarrow{\text{PAR}} \\ \text{CANADA} \xrightarrow{\text{COM}} \text{MANITOBA} \xrightarrow{\text{COM}} \text{WINNIPEG} \end{array}$$

- It is counter-intuitive but true that most common words can be connected using only the synonym relation ($\overset{\text{SYN}}{\implies}$). This demonstrates the high degree of polysemy exhibited by familiar words. Consider the shortest synonym path between **one** and **zero**, a computer scientist’s favorite antonym pair. Every successive word pair exhibits a different sense:

$$\begin{array}{c} \text{ZERO} \xrightarrow{\text{SYN}} \text{CIPHER} \xrightarrow{\text{SYN}} \text{CALCULATE} \xrightarrow{\text{SYN}} \\ \text{DIRECT} \xrightarrow{\text{SYN}} \text{LEAD} \xrightarrow{\text{SYN}} \text{STAR} \xrightarrow{\text{SYN}} \text{ACE} \xrightarrow{\text{SYN}} \text{ONE} \end{array}$$

- Using only the trigger ($\overset{\text{TRG}}{\implies}$) relation, one can connect concepts that occur in the domain of the data used to train the trigger pairs, in this case broadcast news:

$$\begin{array}{c} \text{SMOKING} \xrightarrow{\text{TRG}} \text{CIGARETTES} \xrightarrow{\text{TRG}} \text{MACHINES} \xrightarrow{\text{TRG}} \\ \text{COMPUTERS} \end{array}$$

- The trigger relation enriches the WordNet-derived vocabulary of common nouns with topical proper names, as in the shortest paths shown below. Trigger pairs are often expressible in terms of a sequence of one or more WordNet-derived relations. In many cases, however, news-based triggers defy any fixed set of hand-coded lexical relations.

TITANIC $\xrightarrow{\text{TRG}}$ SANK $\xrightarrow{\text{TRG}}$ SHIP $\xrightarrow{\text{TRG}}$ VALDEZ $\xrightarrow{\text{TRG}}$
COFFEE

NADER $\xrightarrow{\text{TRG}}$ REGULATIONS $\xrightarrow{\text{TRG}}$
ENVIRONMENTAL $\xrightarrow{\text{TRG}}$ GORE

FALWELL $\xrightarrow{\text{TRG}}$ CHRISTIAN $\xrightarrow{\text{TRG}}$
CONSERVATIVE $\xrightarrow{\text{TRG}}$ GINGRICH

- But when the WordNet-derived semantic relations are permitted in addition to the trigger relation, shortest paths become shorter, overcoming the inherent limitations of the data-derived triggers. In the case below, the pair (**relativity**, **physics**) did not occur sufficiently often in training data for the pair to make the grade as a trigger.

EINSTEIN $\xrightarrow{\text{TRG}}$ RELATIVITY $\xrightarrow{\text{PAR}}$ PHYSICS $\xrightarrow{\text{TRG}}$
VELOCITY $\xrightarrow{\text{GEN}}$ SPEED_OF_LIGHT

- For amusement, the phonetic relations, rhymes-with ($\xrightarrow{\text{RHY}}$) and sounds-like ($\xrightarrow{\text{SIM}}$), can be used alone to produce “word ladders” of sequentially similar words, as in the example below. In combination with the semantic relations, the phonetic relations can aid in creating rhymed poetry and puns.

KNIFE $\xleftrightarrow{\text{SIM}}$ NINE $\xleftrightarrow{\text{SIM}}$ SPINE $\xleftrightarrow{\text{SIM}}$ SPOON

Intersection queries

Intersection queries can be used in Lexical FreeNet to find the set of concepts and words that two inputs both directly relate to in some way. We use the notation $(w_1 \xrightarrow{r_1}, w_2 \xrightarrow{r_2})w_3$ to mean that “ w_1 is related to w_3 by relation r_1 , and w_2 is related to w_3 by relation r_2 .”

- For concrete nouns, the results are often expected but sometimes subtle:

(FROG $\xrightarrow{\text{TRG}}$, TURTLE $\xrightarrow{\text{TRG}}$) POND

(ORANGE $\xrightarrow{\text{TRG}}$, APPLE $\xrightarrow{\text{TRG}}$) JUICE

(BANANA $\xrightarrow{\text{TRG}}$, ONION $\xrightarrow{\text{TRG}}$) PEEL

(BOOK $\xrightarrow{\text{TRG}}$, TELEVISION $\xrightarrow{\text{TRG}}$) STORY

(TREE $\xrightarrow{\text{COM}}$, TOOTH $\xrightarrow{\text{COM}}$) CROWN

- Triggers can be a useful tool for discovering what two names in the news have in common, or two names in history:

(STARR $\xrightarrow{\text{TRG}}$, MCDUGAL $\xrightarrow{\text{TRG}}$) WHITEWATER

(CHURCHILL $\xrightarrow{\text{TRG}}$, STALIN $\xrightarrow{\text{TRG}}$)
HITLER, ROOSEVELT, TRUMAN, POTSDAM

- In some cases, identification questions can be formulated as intersection queries. For example, “What’s the name of that congresswoman from Colorado I’m always hearing about?” can be asked as an intersection query with arguments (**congresswoman**, **Colorado**). “What’s the capital of the state of Nebraska?” can be asked as an intersection query with arguments (**Nebraska**, **state_capital**):

(COLORADO $\xrightarrow{\text{TRG}}$, CONGRESSWOMAN $\xrightarrow{\text{TRG}}$)
SCHROEDER

(NEBRASKA $\xrightarrow{\text{PAR}}$, STATE_CAPITAL $\xrightarrow{\text{GEN}}$) LINCOLN

Rhyme coercion queries

The phonetic relations in Lexical FreeNet are particularly useful for finding rhyming words with certain target meanings. The coercion function on the Web interface is hardcoded such that the relation **re1** (see Section 2.3) is simply the union of all semantic relations, and **re2** is the union of all phonetic relations. Thus, given two endpoint words (w_1, w_2) , the system tries to find words (w'_1, w'_2) , with respectively related meanings, that rhyme or sound alike. For example, if you wanted to write a poem about petting a lion, you might do a coercion query with the words **touch** and **lion**. Amongst a few others, you’ll get back the suggestions (**RUB**, **CUB**), since **TOUCH** $\xrightarrow{\text{GEN}}$ **RUB** and **LION** $\xrightarrow{\text{TRG}}$ **CUB**; and (**PAT**, **CAT**), since **TOUCH** $\xrightarrow{\text{GEN}}$ **PAT** and **LION** $\xrightarrow{\text{TRG}}$ **CAT**. Most rhyme coercion queries to the online system have produced at least one result in this manner.

4.2 Text processing

Here we shall somewhat speculatively discuss applications of Lexical FreeNet to applications in which the client is an algorithm rather than a human being.

The behavior of shortest path queries on randomly generated word pairs (Figure 4) suggests that most pairs occurring in text can be joined using only semantic relations, and that both the trigger relation and the WordNet-derived relations play a significant role in shortening connections.

Path length	Triggers	Random pairs		
	Query A	Query A	Query B	Query C
1	2.13%	0.02%	0.05%	0.10%
2	5.20%	0.14%	0.44%	2.27%
3	12.12%	1.55%	7.13%	24.97%
4	11.60%	5.00%	24.99%	48.90%
5	8.17%	5.19%	25.82%	17.61%
6	3.27%	2.94%	10.20%	2.52%
7	0.83%	1.16%	2.23%	0.41%
8	0.22%	0.29%	0.26%	0.05%
9	0.00%	0.14%	0.05%	0.03%
≥ 10	0.02%	0.05%	0.05%	0.05%
No path	56.44%	83.52%	28.78%	3.08%

Figure 4: Summary of shortest path query experiments on selected word pairs. In the second column we report the distribution of lengths (in links) of the shortest paths found for the top 10,000 most mutually-informative trigger pairs. In the other columns we report the lengths of the shortest paths found for a fixed set of 10,000 word pairs generated randomly from the vocabulary. Query “A” permits all semantic relations in the path except the trigger relation. Query “B” permits all semantic relations in the path, including the trigger relation. Query “C” permits all ten semantic, phonetic, and orthographic relations. The “No path” line counts the number of pairs between which there exists no connecting path in the multigraph with the chosen relations. The large number of unconnected pairs in the experiments with Query “A” can be attributed to inflected forms, which are not modeled in the WordNet-derived relations. Triggers capture inflected forms, and so does the “sounds like” phonetic relation, used in the fourth experiment.

Segmentation

Topic segmentation is the task of automatically partitioning a stream of natural language text or multimedia into a sequence of topically cohesive segments.

Machine learning algorithms that learn to identify topic changes in the lexical stream are a fairly recent innovation. An algorithm that employs maximum entropy modeling (Beeferman et al., 1997b) to place topic breaks based on binary *features* of the context of a boundary candidate offers an ideal test bed for this work. Candidate features already in the system include lexical questions of the form “Does the word X appear to the left of the boundary candidate” and “Does an unrecognized word appear on both sides of the boundary candidate.” It is possible to extend this framework to include questions parameterized by *two* words X and Y , and ask questions of the form “Does X appear to the left of the boundary candidate *and* Y appear to the right?” An interesting choice of (X, Y) pairs to try in this experiment is the set of Lexical FreeNet relations, or the set of word pairs that can be joined by a path of at most k semantic relations, for some k . The feature induction framework described in (Beeferman et al., 1997b) offers a mechanism by which features can be selected in the order of their predictive value to the system on training data, and so the value of the various FreeNet relations to this task can be made transparent: the ranking produced by the feature induction algorithm can be analyzed to see how often, and how early, pairs derived from each relation are used in the system. With this experimental setup

we can learn to what extent knowledge-derived links supplement or reinforce data-derived links.

Summarization

Summarization is the task of automatically mapping a topically cohesive segment of natural language to a more compact representation that preserves the meaning of the discourse as much as possible. Today’s summarizers are mostly sentence- or phrase- *selective* in nature; that is, they do not generate novel text, but instead identify the most pertinent regions of the input and output or highlight them. Others extract keywords automatically.

The topic graph (Niwa et al., 1997) is an intriguing graph-theoretic tool that models a discourse as a graph of content word nodes and lexical-relational edges. One could imagine viewing a discourse as a *subgraph* of a semantic net, in which all the content words appearing in the document are present as nodes; and in which all the word pairs appearing in the document within some temporal window that are also related in the semantic net are present as edges. We can then inspect the graph for high-contention nodes, *i.e.* concepts in the document that appear to be connected to many other concepts in the document. The locations of these high-activity areas are likely to be good candidates for selection by a summarizer. The role of FreeNet in this algorithm would simply be to provide answers on which pairs of concepts are related, and how strongly.

Other approaches involving *lexical chains* (Barzilay and Elhadad, 1997) could benefit from the FreeNet engine. An intriguing question to ask about

a discourse unit is “What is the maximum-length chain (or subsequence, but not necessarily contiguous) of lexically-related words appearing in this region?” An approximate answer to such a “longest path” question, known to be NP-complete, could nonetheless fuel a novel sort of keyword extraction.

5 Conclusion

We have introduced a database system called FreeNet that facilitates the description and exploration finite binary relations, and also an instance of the system called Lexical FreeNet that supports a range of lexical semantic applications. The program has proven itself to be a useful and entertaining resource for lexical discovery by end users. We hope to employ the system as a common algorithmic core for three text processing applications as well—segmentation, summarization, and information extraction.

Acknowledgments

The author thanks Michael Turniansky for early feedback on this work; Adam Berger for developing the Trigger Toolkit; Carl Burch for help with the phonetic and orthographic edit distance functions; Bob Harper and John Lafferty for useful discussions; and the many users of the World Wide Web interface who have provided entertaining feedback on the system.

References

- R. Barzilay and M. Elhadad. 1997. Using lexical chains for text summarizations. In *Proc. ACL/EACL-97 Workshop of Summarization*.
- D. Beeferman, A. Berger, and J. Lafferty. 1997a. A model of lexical attraction and repulsion. In *Proceedings of the ACL*, Madrid, Spain.
- D. Beeferman, A. Berger, and J. Lafferty. 1997b. Text segmentation using exponential models. In *Proc. Second Conference on Empirical Methods in NLP*.
- A. Blum, G. Konjevod, R. Ravi, and S. Vempala. to appear. Semi-definite relaxations for minimum bandwidth and other vertex-ordering problems. In *Proc. of the 30th ACM Symposium on the Theory of Computing*, pages 95–100.
- CMU. 1997. Carnegie Mellon University Pronouncing Dictionary v0.6d. <http://www.speech.cs.cmu.edu/cgi-bin/cmudict>.
- C. Fellbaum. 1990. English verbs as a semantic net. *International Journal of Lexicography*, 3,4:278–301.
- D. Gross and K. Miller. 1990. Adjectives in WordNet. *International Journal of Lexicography*, 3,4:265–277.
- H. Kozima and T. Furugori. 1993. Similarity between words computed by spreading activation on an english dictionary. In *Proc. EACL*, pages 232–239.
- LDC. 1997. DARPA Continuous Speech Recognition Corpus-IV: Radio Broadcast News (CSRIV Hub-4). <http://morph ldc.upenn.edu/>.
- I. Mani and E. Bloedorn. 1997. Multi-document summarization by graph search and matching. In *Proc. AAAI-97*, pages 69–79.
- G. Miller. 1985. WordNet: a dictionary browser. In *Proc. First International Conference on Information in Data*.
- G. Miller. 1990a. Introduction to WordNet: An on-line lexical database. *International Journal of Lexicography*, 3,4:235–244.
- G. Miller. 1990b. Nouns in WordNet: a lexical inheritance system. *International Journal of Lexicography*, 3,4:245–264.
- Y. Niwa, S. Nishioka, M. Iwayama, and A. Takano. 1997. Topic graph generation for query navigation: Use of frequency classes for topic extraction. In *Proc. of the Natural Language Processing Pacific Rim Symposium*, pages 95–100.
- C. Papadimitriou. 1976. The NP-completeness of the bandwidth minimization problem. *Computing*, 16:263–270.
- S. Richardson, L. Vanderwende, and W. Dolan. 1993. Combining dictionary-based and example-based methods for natural language analysis. In *Proc. Fifth International Conference on Theoretical and Methodological Issues in Machine Translation*, pages 69–79.
- S. Richardson. 1997. *Determining Similarity and Inferring Relations in a Lexical Knowledge Base*. Ph.D. thesis, The City University of New York.
- R. Rosenfeld. 1994. *Adaptive Statistical Language Modeling: a Maximum Entropy Approach*. Ph.D. thesis, Carnegie Mellon University, April.
- J. Rosenzweig. 1998. WordNet bibliography. See <http://www.cis.upenn.edu/josephr/wn-biblio.html>. Includes 146 citations to publications on applications of WordNet.
- D. Sleator and D. Temperley. 1991. Parsing English with a link grammar. Technical Report CMU-CS-91-196, School of Computer Science, Carnegie Mellon University.