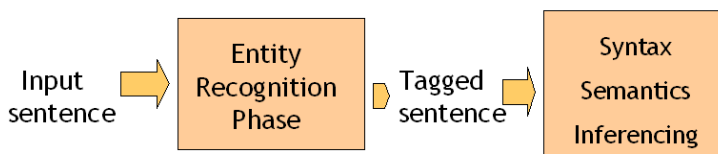


Grammatica: The PathwayArchitect NLP Engine

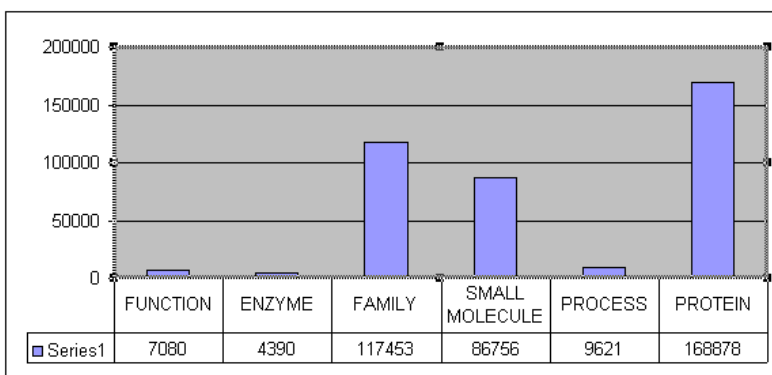
Interactions between biological entities like proteins and small molecules are central to the understanding of biological systems. However, compiling information on interactions is challenging because these are scattered across the breadth of scientific literature. Manual curation of this literature is expensive due to its huge size and massive on-going growth. A practical alternative is to design fully automated natural language processing systems that can extract interactions from literature. *Grammatica* is indeed such a system.

A number of systems have been reported in literature to process text available through sites like MEDLINE. The approaches used by these can be classified into two broad categories, shallow and deep. Shallow approaches use general techniques like pattern matching, machine learning, and statistical methods but limited grammatical structure, examples include Ono et al [1] and Blashke et al [2] (pattern matching) and Park et al [9], Thomas et al [10], Applet et al [2] (shallow parsing possibly clubbed with statistical methods). Deep approaches use detailed grammars to comprehend sentence structure and meaning; these are more laborious to create but provide better accuracy; examples are Yakushiji et al [11], Friedman et al [5], and Novichkova et al [7].

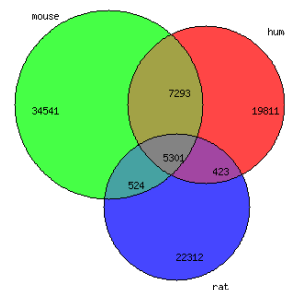
Grammatica is based on deep parsing and driven by an elaborate sentence grammar. The deep approach was chosen to have maximum control over different aspects of the sentence so as to maximize accuracy without compromising recall. Extraction of interactions is done in four main phases namely, named *entity recognition*, *syntax analysis*, *semantic analysis* and *semantic inferencing*. Note that *Grammatica* works sentence by sentence and extracts only interactions which occur completely within a single sentence. Each sentence is put through the four phase pipeline illustrated below.



Entity Recognition: The first step in processing literature is to identify terms in the sentence corresponding to entities of interest. *Grammatica* currently considers the following terms: *proteins*, *protein families*, *enzymes*, *small molecules*, *biological processes*, and *molecular functions*. Other categories of interest, e.g., disease terms, will be added over time. Statistics on the numbers of each term type in the mammalian database are shown in the figure below.



Protein entities were obtained by parsing Entrez Gene data available in ASN format. Human, mouse, rat entries are combined into a single mammal database. In this process, entries that share a common official symbol are merged. Results of the merger are illustrated in the venn diagram here. Protein dictionaries for other organisms are in the process of creation currently. The small molecule dictionary was created using four sources: MeSH 2005 descriptor records, MeSH supplementary concept records, ChEBI and MedLine Name of Substance. Since there is redundancy amongst these sources, records from these sources were merged based on the presence of common aliases and registry numbers. The Enzyme dictionary was obtained from the Expaty site. The Protein Family dictionary was extracted from Medline abstracts using a rule-based algorithm.



For each term category, a detailed list of entity names is collected. Often entities are referred to by multiple names (either alternative names or aliases), so synonyms are identified and clubbed together. Finally, an algorithm for finding these entities names and synonyms in text is needed. A straightforward dictionary lookup is insufficient as there are a number of variations in the way these terms appear in the text, e.g., the presence/absence of hyphens, brackets, and the use of Greek letters like α instead of alpha etc. *Grammatica* uses a tokenization procedure which breaks each term into tokens, identifies matching tokens, and then strings back collections of matching tokens into terms, in the process allowing for some mismatches and variations.

Syntax Analysis: The next phase is the syntax analysis phase in which the syntactic structure of the sentence is derived using a *context free grammar* for English, i.e., a set of transformation rules as illustrated in the picture on the right. The syntactic structure, also called the *syntax tree*, is a hierarchical breakup of the sentence into its underlying

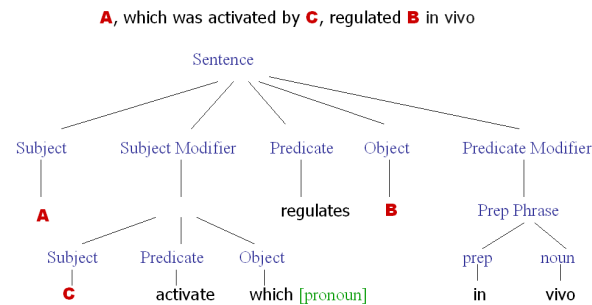
```

Sentence → Subject Predicate Object
Sentence → Subject SubjectMod Predicate Object
Subject → noun
Subject → pronoun
Object → noun | 'that' Sentence | ...
SubjectMod → NounPhrase | Sentence | ...

```

•~200 such rules.

linguistic constituents like nouns, verbs, and adverbial/prepositional clauses and phrases. The syntax tree also captures the functional roles of different parts of the sentence like the subject, the object, the subject modifier, the object modifier and the predicate modifier. The grammar rules are specified in the form of *Augmented Transition Networks (ATN)* (Allen [1], Jurafsky et al [6]). The syntax analyzer or the context free parser is based on the *chart parsing* algorithm (Allen [1]). An example syntax tree is shown below.



Before entity recognized sentences (the stretch containing an entity term is marked up) can be parsed for syntax tree construction, a lexical analysis step is needed. In this step, each word in the sentence is subjected to morphological analysis using a dictionary of English words derived from UMLS Lexicon. The purpose of the morphological analyzer is to identify a word match with the dictionary handling variations caused by change of tense, number and person.

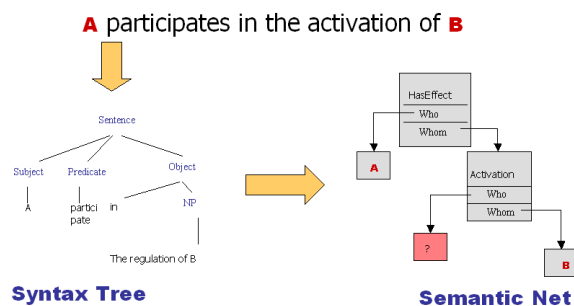
Note that syntax analysis could lead to multiple parse trees for each sentence due to inherent ambiguity in language as well as ambiguity in the grammar used to parse the sentence. All the parse trees resulting from a sentence share portions and are packed into a space efficient structure for further processing.

Also note that several different syntactical structures and sentence formations could carry the same semantic content, as illustrated in the sentences in the figure on the right. The Semantic Analysis and Semantic Inferencing phases aim to round off syntactical structures so different syntactic structures carrying the same meaning are identified as one.

- A** regulates **B**.
 - B** is regulated by **A**.
 - A** plays a role in the regulation of **B**.
 - A** belongs to the family of regulators of **B**.
 - B** activity was found to be modulated by the addition of **A**
- All mean **A** $\xrightarrow{\text{Reg}}$ **B**

Semantic Analysis: The conversion from syntax to semantics is driven by a manually created *semantic dictionary* which maps all words of interest to corresponding semantic concepts. For instance, the word *modulate* would map to the concept *regulate*. The semantic tree would then need to identify what regulates what. It uses the sentence structure imposed by the syntax tree to identify these agents. This process involves elaborate rules for each type of syntactical structure found in a sentence, i.e., prepositional phrases, subordinate clauses etc. Each syntax tree of the sentence could possibly define a different semantic tree; as in the case of syntax trees, these multiple semantic trees are

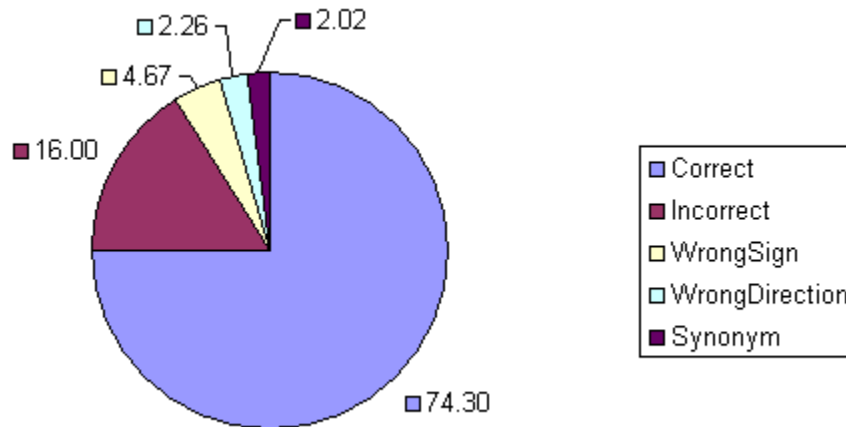
packed into a space efficient structure for further processing. An example semantic tree is shown below.



Note that the relationships captured in a semantic tree are only those which are immediately and directly indicated by the sentence structure. Only a fraction of the interactions are captured directly by such relationships. A good fraction of the interactions are captured indirectly and extracting these interactions requires making inferences across multiple relationships present in the semantic tree; for instance, in the semantic tree in the figure above, who activates B is not directly available unless one runs an inference across the Activate and HasEffect concepts. The next phase of information inferencing does exactly this.

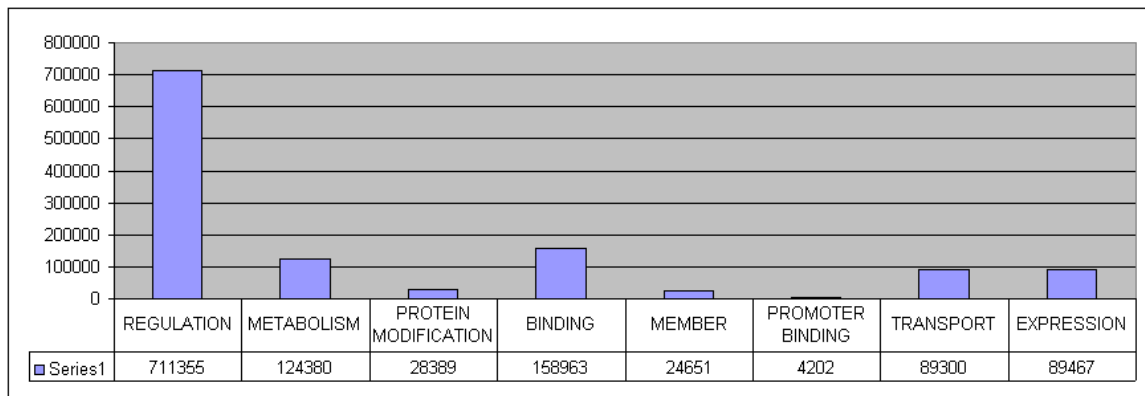
Inferencing: A good example of a sentence requiring inferencing is "Treatment by A causes an increase in the amount of B". The semantic tree would yield the following concepts: *cause, treatment, increase, amount*. To conclude that A regulates B positively, we would need to make inferences across these four semantic concepts, often using agents in one relationship to fill in missing holes in other relationships. This process is again driven by a set of domain-specific inference rules which make several passes through the semantic tree trying to unify the various concepts present in the semantic tree and inferring new concepts. In this process, it is possible that we augment the semantic network with new semantic nodes to represent the inferred concept. The inference rules are encoded as *tree transformation* rules, which specify the mapping from a source semantic subtree to a destination semantic subtree. Finally we extract interactions by searching the resulting semantic network for interaction nodes with all its arguments filled.

Performance and Statistics: The *Grammatica* NLP pipeline described above was run on all of the Medline abstracts available at the beginning on December 2005. The run was performed on multinode Linux clusters. The resulting interactions were sampled at random and examined for correctness. Initial examination showed the need for a *post-processing step* to eliminate interactions obtained from false tagging (e.g., words like ml (milliliter) coming up as proteins, Ca being confused as both protein and small molecule, etc). A post-processing module was then implemented and rules were added to it repeatedly to remove gross errors. Finally, about 500 interactions were drawn at random from the results and manually evaluated for accuracy. The pie chart below shows the accuracy results.



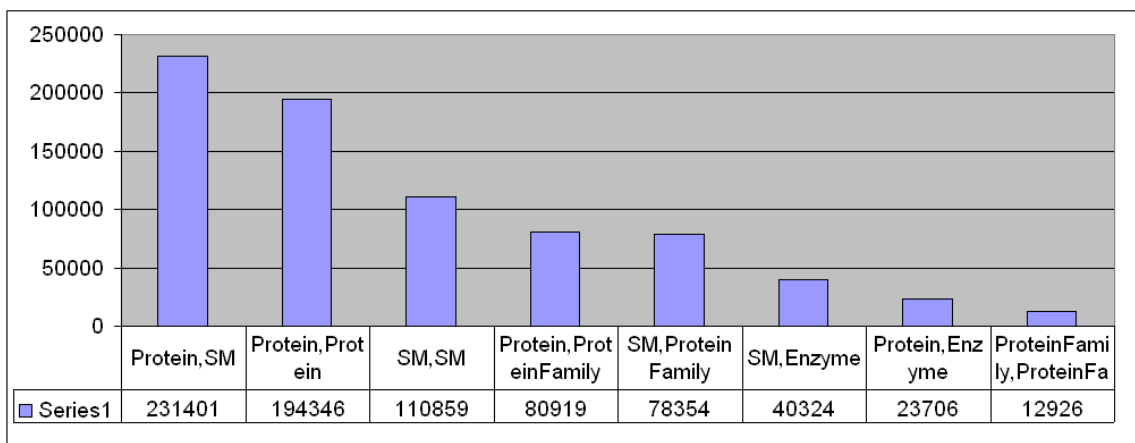
Accuracy: 74.3% of the roughly 500 interactions examined were found correct. An additional 4.67% had only a wrong sign (i.e., positive regulation became negative regulation or vice versa) but were otherwise correct. 2.26% had only a wrong direction (i.e., A regulates B became B regulates A). 16% were clearly incorrect and came about due to a variety of reasons.

Interaction Statistics: The following chart shows the number of interactions of each type that were extracted. The total number of interactions were about 1.2million of which more that half were of the generic regulation variety. Metabolism and Binding were well above 100,000 in number while Transport and Expression were just below 100,000.



Entity Pair Statistics: Note that the same pair of entities (proteins or small molecules etc) could appear multiple times in the above counts as they could participate in different types of interactions or have multiple interactions with the same interaction type but with different effects (positive/negative/unknown regulation). Statistics on just *pairs of interacting entities (irrespective of interaction type or*

effect) is presented in the figure below. In all there are about 850,000 distinct pairs of entities detected as interacting. The break up across entity types is given on the chart below.



Future Plans: In the immediate future, we plan to generalize *Grammatica* to create interaction databases for E.Coli, Arabidopsis, C. Elegans, Yeast, and Drosophila. The goal for the near future is to make the engine customizable to answer a wider range of queries, e.g., what cleavage sites are known about a protein, what genes have been patented for what function etc. The goal in the near future is also to expand the pool of full text articles on which *Grammatica* runs for greater coverage. Finally, the goal in the near future is to improve the accuracy by tightening mis-parsing and incorrect entity tagging.

References:

1. James Allen. Natural Language Understanding. Benjamin/Cummings, 2nd edition, 1995.
2. D. Applet, J. Hobbs, J. Bear, D. Israel, M. Kamayaqma, A. Kehler, D. Martin, K. Meyers, and M. Tyson. SRI International {FASTUS} system: Muc-6 test results and analysis. In {\em Proceedings of the 6th Message Understanding Conference}, pages 237--248, 1995.
3. C. Blaschke, M. A. Andrade, C. Ouzounis, and A. Valencia. Automated extraction of biological information from scientific text: protein-protein interactions. ISMB, pages 60--67, 1999.
4. I. Donaldson, J. Martin, B. deBruijn, C. Wolting, V. Lay, B. Tuekam, S. Zhang, B. Baskin, G. D. Bader, K. Michalickova, T. Pawson, and C. W. Hogue. Pre-BIND and textomy-mining the biomedical literature for protein-protein interactions using a support vector machine. BMC Bioinformatics, 4(1), 2003.
5. C.~Friedman, P. Kra, H. Yu, M. Krauthammer, and A.~Rzhetsky. GENESIS: A natural language processing system for the extraction of

molecular pathways from journal articles. *Bioinformatics*, 17 (Suppl.1):S74--S82, 2001.

6. Daniel Jurafsky and James H. Martin. *Speech and language processing*. Prentice-Hall, 2000.

7. S. Novichkova, S. Egorov, and N. Daraselia. Medscan, a natural language engine processing engine for {MEDLINE} abstracts. *Bioinformatics*, 19:1699--1706, 2003.

8. T. Ono, H. Hishikagi, A. Tanigami, and T. Takagi. Automated extraction of information on protein-protein interactions from the biological literature. *Bioinformatics*, 17:155--161, 2001.

9. J. C. Park, H. S. Kim, and J. J. Kim. Bidirectional incremental parsing for automatic pathway identification with combinatorial categorical grammar. *Pac. Symp. Biocomput.*, 6:396--407, 2001.

10. J. Thomas, D. Milward, C. A. Ouzounis, S. Pulman, and M. Carroll. Automated extraction of protein interactions from scientific abstracts. *Pac. Symp. Biocomput.*, pages 541--552, 2000.

11. A. Yakushiji, Y. Tateisi, Y. Miyao, and J. Tsujii. Event extraction from biomedical papers using a full parser. *Pac. Symp. Biocomput.*, 6:408--419, 2001.