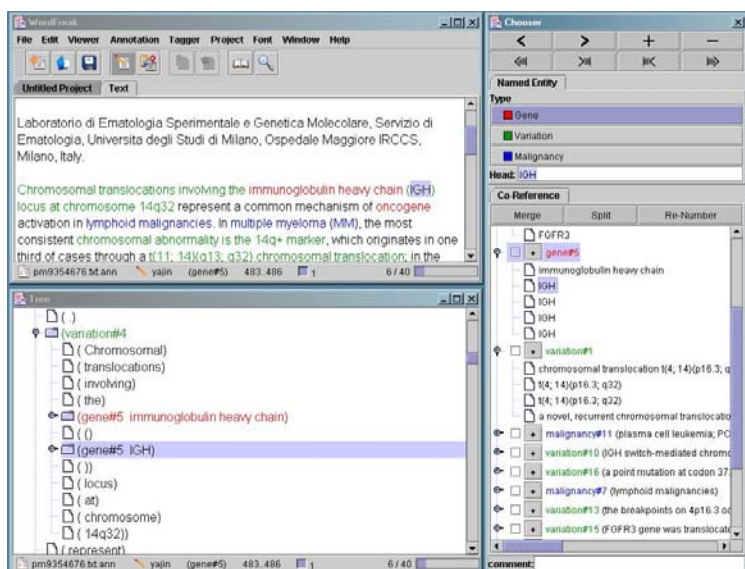EIA-0205448: "Mining the Bibliome: Information Extraction from the Biomedical Literature"
University of Pennsylvania
Investigators: Aravind Joshi, Susan Davidson, Mark Liberman, Mitch Marcus, Martha Palmer
Biomedical collaborators: Peter White (CHOP, eGenome), Paula Matuszek (GSK)

**Websites:** http://www.ldc.upenn.edu/myl/ITR
http://www.cis.upenn.edu/~mamandel/term.htm

Our goal is qualitatively better methods for automatically extracting information from the biomedical literature, relying on three techniques: high-accuracy parsing, shallow semantic analysis, and integration of existing databases. An initial step is to create annotated corpora in collaboration with biomedical researchers: two test cases are gene variations in pediatric oncology, and inhibition of CYP450 enzymes.

The screen shot below shows part of a scientific abstract that we have annotated for entities and relations relevant to a project (based at Children's Hospital of Pennsylvania) to construct a database of cancer-associated mutational events.

The many millions of biomedical publications available in electronic form contain a vast quantity of scientific information. Researchers would like access to this information structured in terms of well-defined relations (like "inhibition" or "mutation") among entities of interest (like "gene", "compound" or "cell line"). Recent techniques from computational linguistics can make more of this information accessible to biomedical researchers; at the same time, large existing biomedical databases promise an unparalleled depth of semantic support to make linguistic analysis more accurate.

To experiment with new textual information-extraction techniques, training and testing corpora are needed. To this end, we've developed or adapted software tools that allow human experts to annotate biomedical texts for relevant entities and relations ("entity tagging"), to mark syntactic structure (producing a "tree bank"), and to indicate shallow semantic structure, such as co-reference relations and predicate-argument relations (producing a "proposition bank"). We've defined compatible forms of "stand-off markup" (that is, annotations kept separate from the texts annotated), in order to allow multiple independent forms of text annotation to be created and used in an integrated way. We're producing annotation specifications and training materials for syntactic and semantic annotation of biomedical text, and an annotation crew is testing these tools and specifications by multiple annotation of a set of initial test sets, in order to ensure adequate inter-annotator agreement. Once enough annotated material has been created, automatic taggers and parsers are trained for each task, and the annotators' job becomes correcting the programs' output, or creating initial material for a new application or a new dimension of analysis.

All of the software developed in this project is available under open-source license, and all annotated texts and resulting databases will be published for the use of other researchers.