# NLP and Text Mining in Bioinformatics

Manu Madhavan

# Bioinformatics Literature

- Much literature generated quickly.
  - 11 million citations in MEDLINE.
  - 400,000 added yearly.
- Need methods to deal with data.
  - Query
  - Summarize
  - Organize
  - Understand
- For mankind to benefit from bioinformatics research, the sequence and structure of proteins and other molecules **must be linked to functional genomics and proteomics.**

# PubMed

- The primary store of functional data that links clinical medicine, pharmacology, sequence data, and structure data is in the form of biomedicine documents in online bibliographic databases such as PubMed
- Mining these databases is expected to reveal the relationships between structure and function at the molecular level and their relationship to pharmacology and clinical medicine.
- **Text mining**—automatically extracting this data from documents, which is published in the form of unstructured free text, often in several languages

# The process



IF &lt;protein name&gt;
AND &lt;experimental method name&gt; are in the same sentence
THEN the &lt;experimental method name&gt; refers to the &lt;protein name&gt;

# NLP Research areas

- Information Retrieval
- Information Extraction
- Q&A
- Named Entity Recognition
- Entity Relation Extraction
- Curated Databases
- Text summarization
- Ontology/Knowledge graph

# Information Extraction

# Gene Ontology and Mesh Terms

- Gene Ontology:The Gene Ontology (GO) describes our knowledge of the biological domain with respect to three aspects: Molecular Function. Molecular-level activities performed by gene products.

- Mesh Terms: The Medical Subject Headings (MeSH) thesaurus is a controlled and hierarchically-organized vocabulary produced by the National Library of Medicine.

# Gene Ontology

http://geneontology.org/docs/ontology-documentation/


Gene Enrichment Analysis: http://geneontology.org/

# Mesh

# Resources

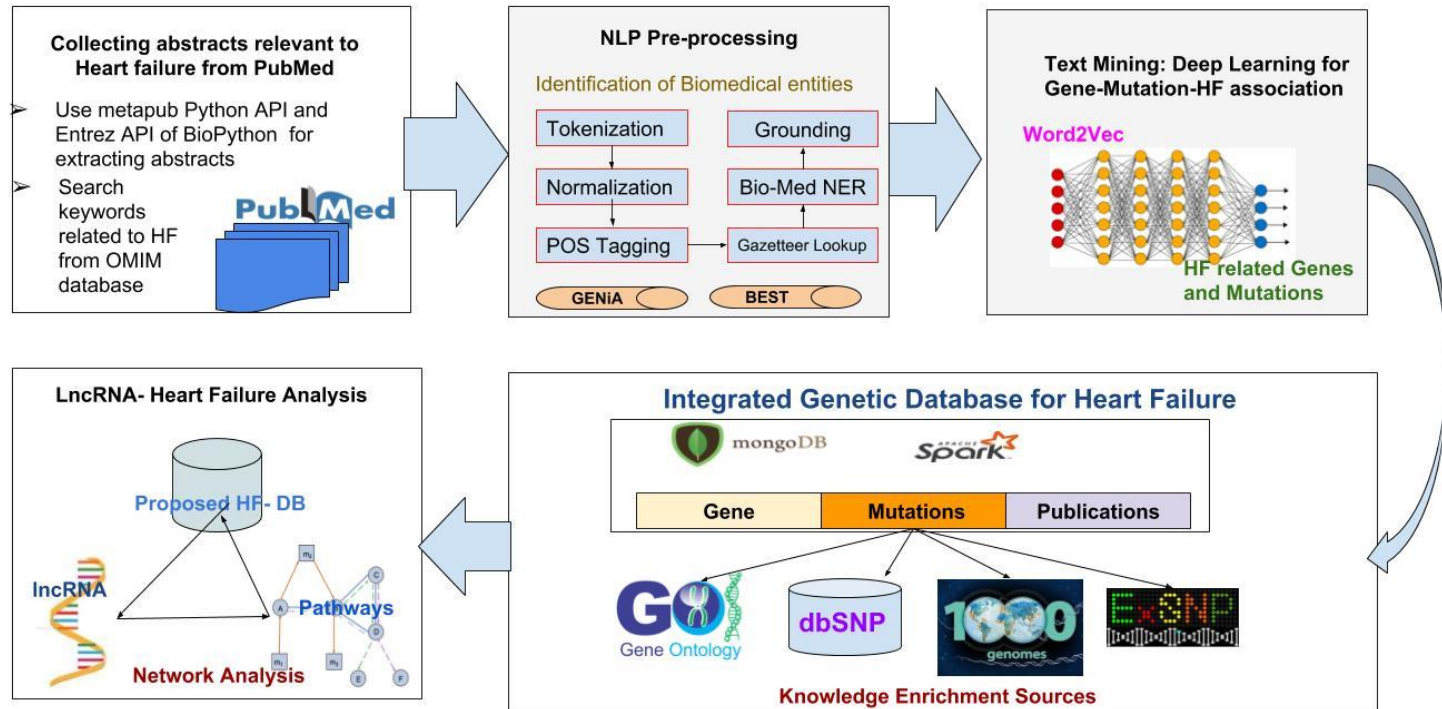| | |
|---|---|
| Informatics for Integrating Biology and the Bedside (i2b2 - https://www.i2b2.org/) | National Center for Biomedical Computing with focus on translational research t facilitates and proves data sets for clinical natural language processing research |
| Gene Ontology (https://www.geneontology.org) | Controlled vocabulary with relationships including partonymy and inheritance, designed for describing gene functions, broadly construed |
| Entrez Gene (https://www.ncbi.nlm.nih.gov/gene) | Source for gene names, symbols, and synonyms; also the source for GeneRIFs ar SUMMARY fields |
| PubMed/MEDLINE (https://www.ncbi.nlm.nih.gov/pubmed) | The National Library of Medicine's database of abstracts of biomedical publicatic (MEDLINE) and search interface for accessing them (PubMed) |
| Unified Medical Language System (https://www.nlm.nih.gov/research/umls/) | Large lexical and conceptual resource, including the UMLS Metathesaurus, which aggregates a large number of biomedical and some genomic vocabularies |
| SWISSPROT (https://www.uniprot.org/) | Database of information about proteins with literature references, useful as a go standard |
| PharmGKB (https://www.pharmgkb.org/) | Database of relationships between a number of clinical, genomic, and other enti with literature references, useful as a gold standard |
| Comparative Toxicogenomics Database (https://ctdbase.org/) | Database of relationships between genes, diseases, and chemicals, with literatur references, useful as a gold standard |

Various terminological resources, data sources, and gold-standard databases for biomedical natural language processing.
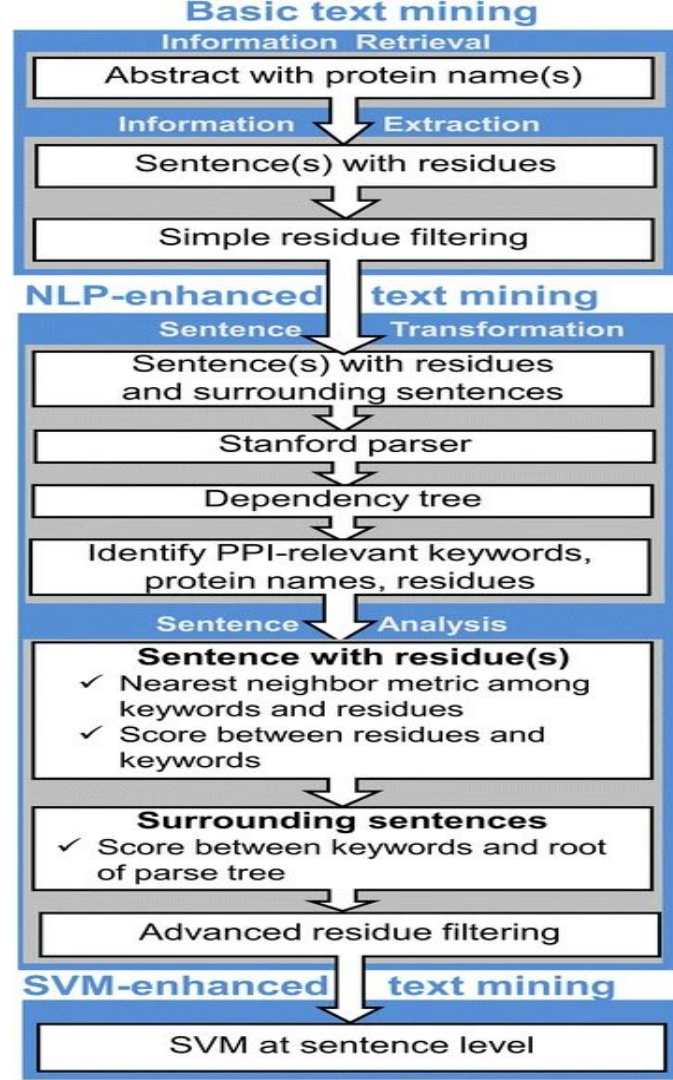doi:10.1371/journal.pcbi.1003044.t001

# Case Studies

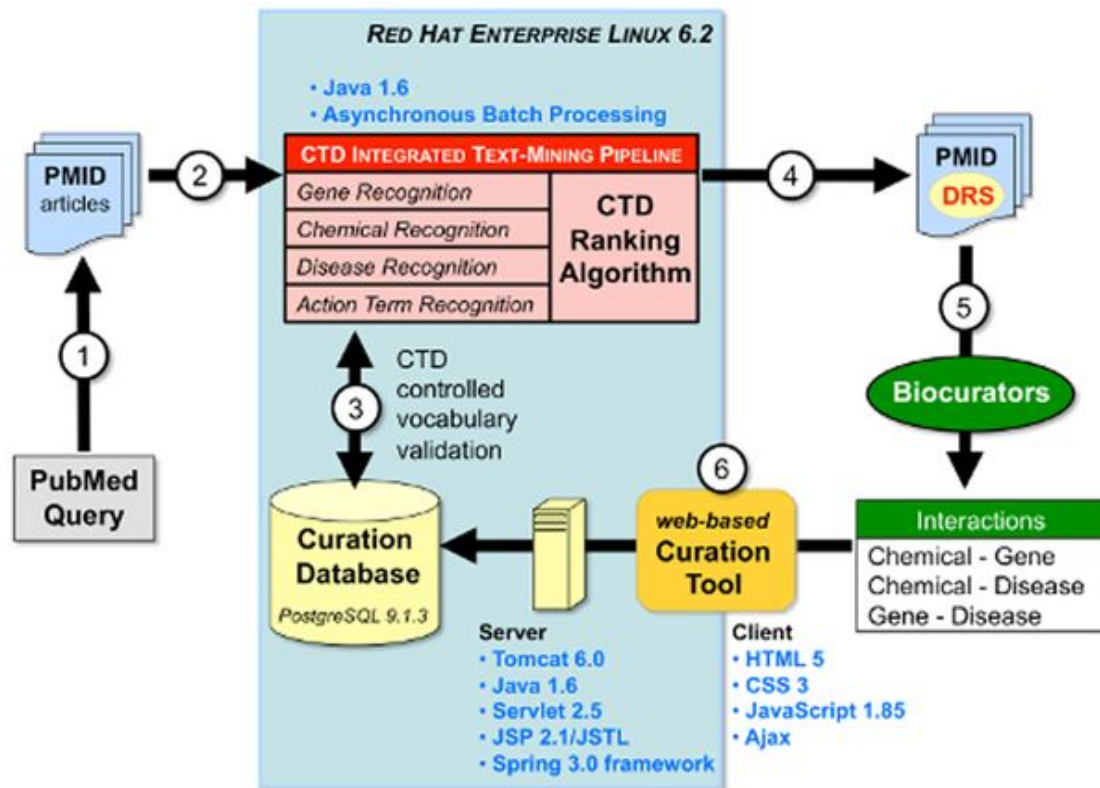# Text Minded database of Mutations in HF

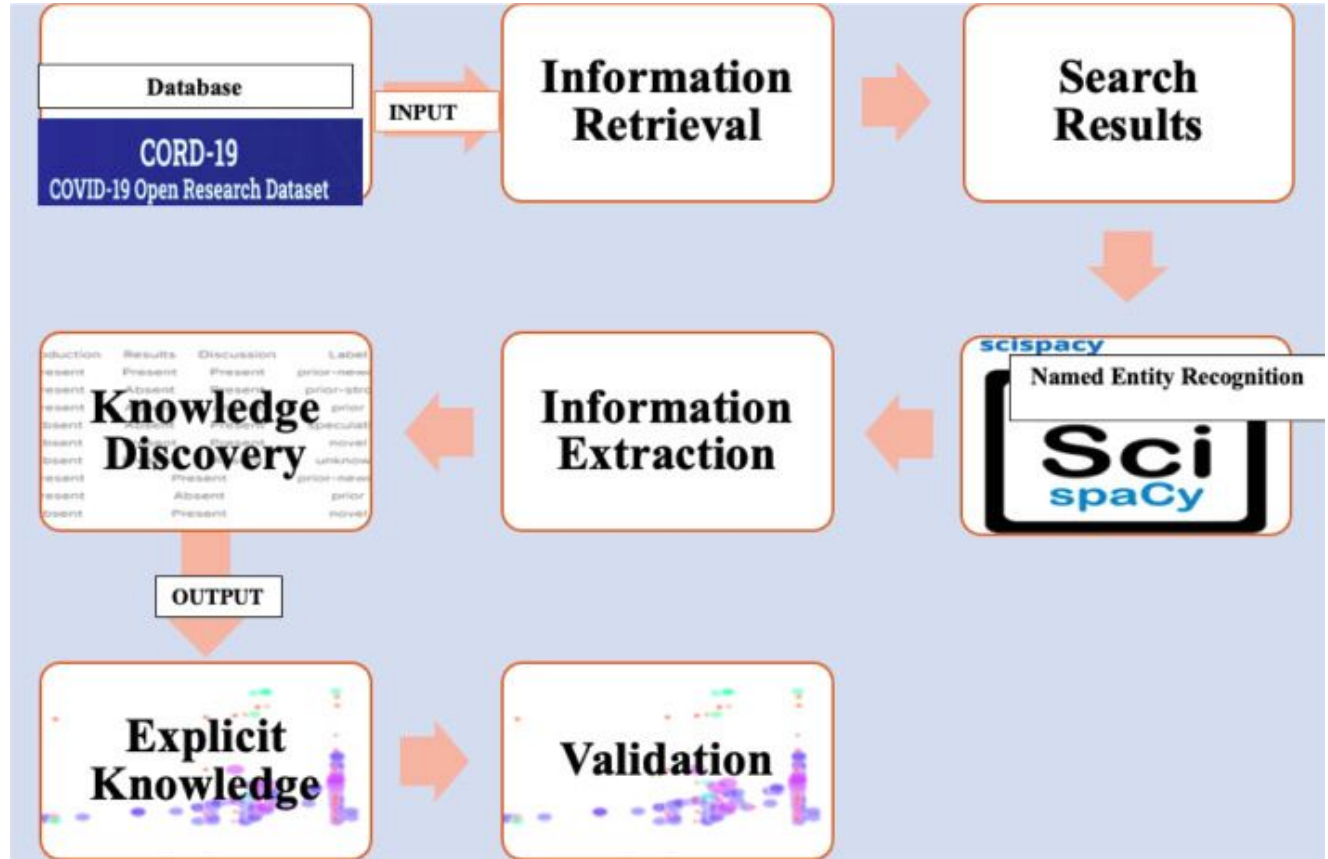# Text mining for structural modeling of protein complexes

- TM tool, which utilizes natural language processing (NLP) for analyzing the context of the residue occurrence
- TM procedure for extracting protein-protein binding site residues from the PubMed abstracts was significantly advanced by the deep parsing (NLP techniques for contextual analysis) in purging of the initial pool of the extracted residues.

# Biocuration

# Covid Dataset

# Pubtator Demo

https://www.ncbi.nlm.nih.gov/research/pubtator/

# NCBI text mining tools

https://www.ncbi.nlm.nih.gov/research/bionlp/Tools/

# More Examples

1. Pathway extraction and reasoning
2. Gene prioritization and gene function prediction
3. Pharmacology
4. Drug repurposing