# Text Mining Special Interest Group

[Therese Vachon, NIBR]

Informatics and Knowledge Management/Information and Knowledge Engineering

Novartis Institute for Biomedical Research, Cambridge, MA 4-6th October 2004

# At NIBR, Text Mining is considered a solution for

- Analyzing, tagging, annotating, exploring, structuring and classifying textual data and large document sets (internal and external)
- Identifying meaningful concepts from text and relationships between concepts
- Improving the quality of text retrieval methods
- Detecting novel patterns
- Detecting similarities across textual data
- Improving navigation across data sources and document sets

# What are the Critical Business Needs that Text Mining Could Address?

- Discovery of unexpected relationship and relevant information
- Generation of new hypotheses
- Assisted annotation and curation
- Unified view of heterogeneous sources
- Analysis of trends and patterns
- Analysis of complex relationships between data elements
- Detection of unexpected or emerging information
- Knowledge inference
- Contextual and semantic navigation

# What Questions do Users Want to Ask ?

- Find information about specific products, targets, genes, proteins, companies, etc.
- Link literature and experimental data
- Identify biological interactions
- Semantic text analysis of full text papers
- Dynamic data flow analysis and categorization
- Monitoring competitors' capabilities and activities
- Discovering networks of associations

The **PRISM** Forum

# Text Mining Applications/Components

| Function | Current | Future |
|---|---|---|
| Text retrieval methods | Keyword searching, fuzzy search, semantic search | Case-based reasoning<br>Query inference |
| Knowledge representation | Thesauri, taxonomies, ontologies | UIMA (unstructured information management architecture), … |
| Annotation | Manual annotation | Assisted annotation |
| Categorization & Clustering | K-means partitioning, hierarchical clustering, SOMs, rule-based categorization | SVM, Latent semantic analysis |
| NLP or OBIIE | Concept extraction<br>IE 1st generation | IE 2nd generation<br>Analysis of positive and negative associations among data objects |
| Full text access | Indexing , access and mining of full text internal documents | Indexing, access and mining of full text literature articles |
| Visualisation | Interactive network of co-occurences (XML SVG)<br>Heat maps, factorial maps, etc. and graphlets technology | Network of typed relations |

# Current Text Mining Applications in Novartis/NIBR?

- Knowledge Space Portal
- Novartis Knowledge Miner (Ulix)
- Competitive Intelligence Analysis Platform
- News Analysis Platform
- Research EIS
- Text mining components

# What Future Applications are Planned?

- Text mining in genomics
- Text mining in chemoinformatics
- Text mining of full text publications
- Text mining for marketing and sales
- Mining of medical enquiries
- Expert location systems
- Generic text analysis platform
- Generic categorical variable analysis platform

The **PRISM** Forum

# Knowledge Space Portal

- Provide key elements for efficiently accessing Novartis-internal and external information relevant to daily decision in the drug discovery and development process:
    - Data integration across heterogeneous data sources and applications (internal and external)
    - Consistent user interface for data retrieval, exploration and analysis across all data types
    - Contextual (ultralink), tree-based (static or dynamic taxonomies) and semantic (knowledge map) navigation
    - Data exploration and analysis methods
    - Personalized views
    - Collaborative, annotation and information sharing tools
    - Alerting

The **PRISM** Forum

# Display-Navigation-Ultralink
Protease modulator in Literature DB (Medline-Embase)

Easy navigation in record titles

Sort capabilities

Ranking value and access to document

Analysis tools

Search report: Number of Docs, Key-words extracted

**Search** Small screen

Text | Chem. Str. | Concept

- RDS Minutes
- RDS Registration
- RDS Scientific Rationale
- RDS Animal Models
- RDS Competitive Assessment
- Porto
- PREDI
- COMP. INTELLIGENCE
- LITERATURE
  - Medline-Embase
  - eNova
- NEWS
- PATENTS
- CHEMISTRY & BIOLOGY
- BIOINFORMATICS
- INTERNET

**For :**

protease modulator

GO | Clear

| Rank | Title | Source | Publication year |
|---|---|---|---|
|  | ...eir function in neuronal | Medline-Embase | 2004 |
| 69 | A solid-phase assay for identification of modulators of prion protein interactions | Medline-Embase | 2003 |
| 68 | ClpS, a substrate modulator of the ClpAP machine | Medline-Embase | 2002 |
| 68 | Protease-Activated Receptor-1 (Thrombin Receptor) is Expressed in Mesenchymal Portions of... Follicle | Medline-Embase | 2003 |
| 68 | Portrait of a killer: the mitochondrial ap... emerges from the shadows. | Medline- | |
| 68 | ...family | | |
| 68 | ...family | Embase | |
| 68 | ...an intra-golgi ...the human... protease HsApg4A | | |
| 68 | Modulation of protein-protein interactions with small ...ules | Medline-Embase | 2003 |
| | ...etween antiretroviral drugs and drugs ...herapy of the metabolic complications ...durin HIV infection. | Medline-Embase | 2002 |
| | ...rminus of GATE-16, an ...ulator is cleaved by the ...g4A. | | |
| | ...l medical countermeasur... | | |
| | ...ators of cell-cell and cell... | | |

**Proteinase activated receptor 1**
Search in OMIM
Search in GENBANK
Search in RefSeq
Search in SwissProt/TrEMBL
Requery this term inside KSP

**HIV infection**
HIV infection [diseases] Portfolio Analysis
Products in development for HIV infection
Search in Ensembl
Search in Harrison
Search in HON
Knowledge Map
Launched Products
Search in OMIM
Requery this term inside KSP

Products (30)
Diseases (38)
- Alzheimer disease
- angioneurotic edema
- atherosclerosis
- breast cancer
- cancer
- chemoprophylaxis
- dermatitis
- diabetic foot

**Resultset Tools**
- Filtering
- Clustering
- Data Analysis
- Graph navigator

**Search Statistics**

**Search Returned**
168 documents for "protease modulator"

**Documents by source:**
Medline-Embase      168

**Attributes:**
EMTREE Index term      1863
Author      716
Targets      277
Diseases      136

# Graph Navigator –

Protease modulators in CI DBs July 2004 - ADIS & Pharmaprojects



Graph navigator

# Clustering

# Data Analysis –

## Protease modulators in CI DBs July 2004 - ADIS & Pharmaprojects
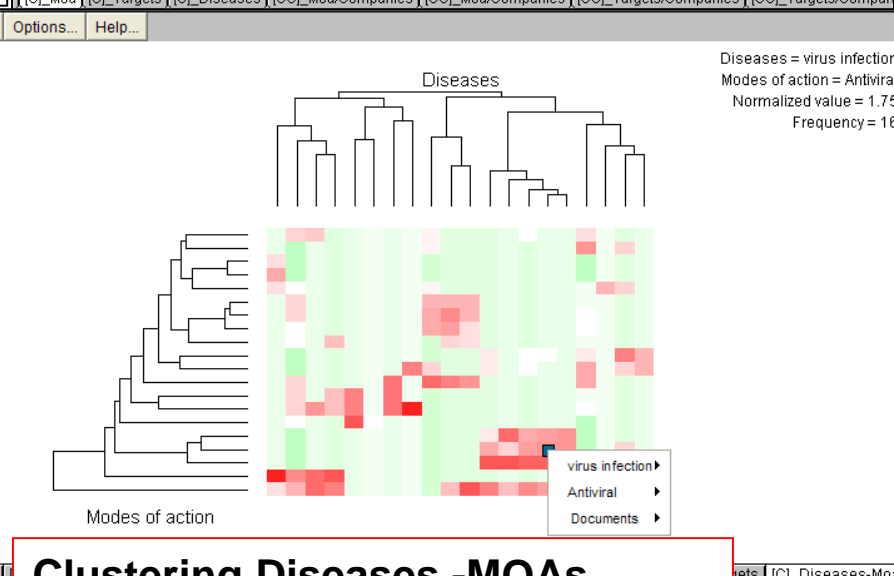
Data Analysis



**Univariate - Companies**

**Univariate - MOA**

**Univariate - Diseases conditionned by Companies**

**Clustering Diseases -MOAs**

# Conclusions

- Text mining techniques have already been implemented in relevant areas
- Additional techniques need to be developed/tested/implemented especially in the field of information extraction/NLP and in the field of categorization
- Collaboration with external partners/research institutes needed
- Text mining components to be applied to all applications dealing with text

The **PRISM** Forum