

Jul 6th, 12:00 AM

## **Innovative usage of unstructured information sources: From text- and data-mining to model-driven decision-support**

Martin Hofmann-Apitius

*Fraunhofer Institute for Algorithms and Scientific Computing SCAI, Sankt Augustin, Germany*

---

Martin Hofmann-Apitius, "Innovative usage of unstructured information sources: From text- and data-mining to model-driven decision-support." *Proceedings of the IATUL Conferences*. Paper 3.  
<https://docs.lib.purdue.edu/iatul/2015/keynotes/3>

This document has been made available through Purdue e-Pubs, a service of the Purdue University Libraries.  
Please contact [epubs@purdue.edu](mailto:epubs@purdue.edu) for additional information.

---

# “Innovative usage of unstructured information sources: From text- and data-mining to model-driven decision-support”

---

Prof. Dr. Martin Hofmann-Apitius

Head of the Department of Bioinformatics

Fraunhofer Institute for Algorithms and Scientific Computing

**36th Annual IATUL Conference 2015**

July 6 - 9, Hannover, Germany

---

# Where I come from: Fraunhofer Society



- Founded 1949
- Europe´s largest applied research organisation
- 60 Research Institutes (7 Institutes in the US)
- > 23.000 Employees
- Annual Budget > 2 Billion Euro
- Financial model: 2/4 industry collaborations  
1/4 public funding  
1/4 institutional funding

\*Joseph von Fraunhofer (1787 – 1826 )  
Scientist, Inventor and Entrepreneur

# The Fraunhofer Institute Center Schloss Birlinghoven

- Largest research centre for informatics and applied mathematics in Germany
- Around 700 employees, thereof 500 scientists, approx. 200 students and trainees
- University links:
  - Bonn
  - Aachen



# Expertise at the Department of Bioinformatics at SCAI

Fraunhofer SCAI Department of Bioinformatics currently comprises:

- 10 scientists
- 3 scientific software developers
- 7 PhD students
- ~ 5 Master students
- ~ 10 student workers
  
- predominantly computer scientists & biologists
  
- some PhD students via University of Bonn (Bonn-Aachen International Center for Information Technology)



# SCAI Department of Bioinformatics: R&D in a nutshell

## Fraunhofer SCAI Department of Bioinformatics R&D activities:

### 1. Information extraction in the **life sciences**:

- Recognition of named entities and relationships in text
- Large-scale, automated Information Extraction

### 2. Integrative biology; disease modelling

- Focus on neurodegenerative diseases

### 3. eScience, Grid- & Cloud- Computing / HPC (Cluster)

- Focus on scaling of information extraction workflows

*Making Scientific Content  
available for Computing*

---

# Imagine ...

---

- **Cancer Patient, final stage, metastatic pancreas carcinoma**
  - Surgery, chemotherapy without success
  - Distant metastasis in bone marrow, lung and liver
  - Remaining life span: 2 – 4 weeks
- **Last chance: sequencing the cancer genome (< 10k€)**
  - Getting insight into mutations underlying cancer dysregulation
  - Understanding of mechanisms triggering uncontrolled growth
  - Identification of (experimental) compounds that inhibit tumour growth
- **This is not fiction – this is reality**

---

# The challenge ...

---

- **Cancer genome sequencing delivers vast amount of information**

- Tens to hundreds of mutations
- Functional relevance of a significant number of mutations unclear
- Contribution of mutation to tumour growth and metastasis ?

- **How do we assess the biological impact of genetic variation information?**

- Putting genetic variation information into a functional context
- Reasoning over genetic variation information and inference of consequences
- From inference to personalised recommendation .... within 2 weeks of time

- **Let us see where we stand ...**



---

# Semantic Search and Knowledge Discovery in Scientific Literature

---



**SCAIVIEW**



A thick red underline with a curved bottom edge is positioned below the word "SCAIVIEW".

Identification and normalisation of the relevant Life Science terminology is key for information retrieval, information extraction and inferring of knowledge

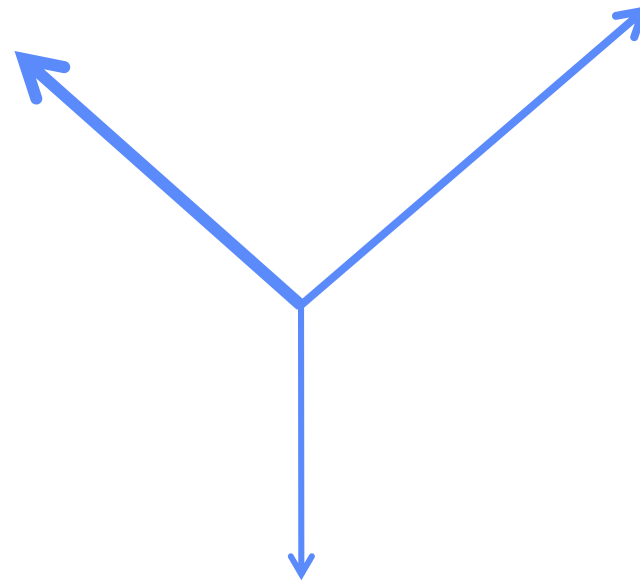
# Named Entity Recognition Technologies @SCAI



Machine learning based approaches  
(CRF based)  
SNP, IUPAC, epigenetic modifications

Dictionary based approaches  
with normalisation and  
embedding in hierarchies  
  
Genes/Proteins, Disease, Drugs,  
Cells...

Regular expression  
(partly to normalise)  
Chrom. Locations, rs numbers...



Combined entities  
Combined normalisation  
SNP, Histonmodification

# Different input output formats



MEDLINE abstracts



ASCII text

XML text,  
e.g. patents

Defined  
format

HTML text e.g.  
Journal articles,  
Web pages



PDF text e.g. Journal articles

## 1. Association of breast cancer resistance protein/ABCG2 phenotypes and novel promoter and intron 1 single nucleotide polymorphisms.

PubMed 18180275 **Authors:** Balasubramanian Poonkuzhali, Jatinder Lamba, Stephen Strom, Alex Sparreboom, Kenneth Thummel, Paul Watkins, Erin Schuetz, **Date:** 2008-04- **Journal:** Drug metabolism and disposition: the biological fate of chemicals **SciMag:** 0.43

Statistics

The hypothesis was tested that sequence diversity in breast cancer resistance protein (BCRP)'s cis-regulatory region is a significant determinant of BCRP expression. The BCRP promoter and intron 1 were resequenced in lymphoblast DNA from the polymorphism discovery resource (PDR) 44 subset. BCRP single nucleotide polymorphisms (SNPs) were genotyped in donor human livers.

The screenshot shows a journal article page from Matrix Biology. At the top, there are logos for Elsevier, ScienceDirect, and Matrix Biology. The article title is "Ucma — A novel secreted factor represents a highly specific marker for distal chondrocytes". The authors listed are Andreas Tagariello, Julia Luther, Melanie Streiter, Lydia Didt-Kozziel, Manuela Wuelling, Cordula Surmann-Schmitt, Michael Stock, Nadia Adam, Andrea Vorkamp, and Andreas Winterpacht. The article includes affiliations for the Institute of Human Genetics, Center for Medical Biotechnology, and Department of Experimental Medicine I. It also shows the abstract, which discusses the growth and development of the vertebrate skeleton and the isolation of a novel transcript from a human fetal growth plate cartilage cDNA library.

# Named Entity Recognition and Normalisation

Chromosomal Locations  Drug Names  Protein/Gene  iupac  OMIM Reference  
 @neurIST  non Normalized SNP  Normalized SNP  MeSH Disease  Relations  
 Cell Lines  Genetic Association  CRF-laura  Brenda

## 1. Association of breast cancer resistance protein/ABCG2 phenotypes and novel promoter and intron 1 single nucleotide polymorphisms.

PubMed 18180275 Authors: Balasubramanian Poonkuzhali, Jatinder Lamba, Stephen Strom, Alex Sparreboom, Kenneth Thummel, Paul Watkins, Erin Schuetz, Date: 2008-04- Journal: Drug metabolism and disposition: the biological fate of chemicals SciMago: 0.43

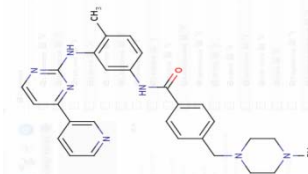
Statistics

The hypothesis was tested that sequence diversity in breast cancer resistance protein (BCRP)'s cis-regulatory region is a significant determinant of BCRP expression. The BCRP promoter and intron 1 were resequenced in lymphoblast DNA from the polymorphism discovery resource (PDR) 44 subset. BCRP single nucleotide polymorphisms (SNPs) were genotyped in donor human livers, intestines, and lymphoblasts quantitatively phenotyped for BCRP mRNA expression. Carriers of the -15622C>T SNP had lower BCRP expression in multiple tissues. The intron 1 SNP 16702C>T was associated with high expression in livers; 1143G>A was associated with low expression in intestine; 12283T>C was associated with higher expression in the PDR44 and White livers. The -15994C>T promoter SNP was significantly associated with higher BCRP expression in multiple tissues. Patients with the -15994C>T genotype had substantially higher clearance of p.o. imatinib. We next determined whether BCRP expression was related to polymorphic alternative splicing or alternative promoter use. Liver polymorphically expressed an alternatively spliced mRNA [splice variant (SV) 1] skipping exon 2. Although SV1+ livers did not uniformly carry the exon 2 G34A allele, 90% of G34A livers expressed SV1 (versus 4% of 34GG livers). BCRP mRNA was significantly lower among Hispanic livers with the G34A variant genotype and may be due, in part, to polymorphic exon 2 splicing. Analysis of allele expression imbalance (AEI) showed that PDR44 samples with AEI had lower BCRP mRNA expression; however, no linked cis-polymorphisms were identified. BCRP used multiple promoters, and livers differentially using alternative exon 1b had lower BCRP. In conclusion, BCRP expression in lymphoblasts, liver, and intestine is associated with novel promoter and intron 1 SNPs.

GeneID: 9429: ABCG2  
GO:0016887: ATP hydrolase  
KEGG:02010: ATP transporter

Neoplasms by Site  
Breast Neoplasms (D001943)  
Breast Neoplasms, Male  
Carcinoma, Ductal, Breast  
Phyllodes Tumor

## DB00619: imatinib



refSNP ID: rs2231137  
ABCG2-Position:18897  
Avg Het: 0.203+/-0.246

multiple sclerosis

Subcorpus Statistics Server Statistics

The following entities relating to 'multiple sclerosis' were found in 1244 documents.

- Human Genes / Proteins
- Chromosomal Location
- STS Marker
- non Normalized SNP
- Normalized SNP
- Normalized CRF SNP
- Drug Names
- ALIMENTARY TRACT AND METABOLISM
- ANTIINFECTIVES FOR SYSTEMIC INFECTIONS
- ANTINEOPLASTIC AND IMMUNOSUPPRESSANT DRUGS
- ANTIPARASITIC PRODUCTS, INSECTICIDES AND VERMIFUGES
- BLOOD AND BLOOD FORMING COMPOUNDS
- CARDIOVASCULAR SYSTEM
- DERMATOLOGICALS
- GENITO URINARY SYSTEM AND MUSCULO-SKELETAL SYSTEM
- + NERVOUS SYSTEM
- RESPIRATORY SYSTEM
- SENSORY ORGANS
- SYSTEMIC HORMONAL PREPARATIONS
- VARIOUS
- IUPAC-like
- Corpora
- Human miRNA
- Mouse Genes
- MeSH Disease
- @neurIST Ontology
- GO Component
- GO Function
- GO Process

SCAIVIEW

Filter next search in the current entity class exclusively to this entity

multiple sclerosis

Result for 'multiple sclerosis', NER run 'DrugBank' for entity **Carbamazepine**, Page 0 with 10 documents per page, totals to 81 and took 281 ms.

- Human Genes / Proteins
- Chromosomal Location
- STS Marker
- non Normalized SNP
- Normalized SNP
- Normalized CRF SNP
- Drug Names
- ALIMENTARY TRACT AND METABOLISM
- ANTIINFECTIVES FOR SYSTEMIC INFECTIONS
- ANTINEOPLASTIC AND IMMUNOSUPPRESSANT DRUGS
- ANTIPARASITIC PRODUCTS, INSECTICIDES AND VERMIFUGES
- BLOOD AND BLOOD FORMING COMPOUNDS
- CARDIOVASCULAR SYSTEM
- DERMATOLOGICALS
- GENITO URINARY SYSTEM AND MUSCULO-SKELETAL SYSTEM
- + NERVOUS SYSTEM
- RESPIRATORY SYSTEM
- SENSORY ORGANS
- SYSTEMIC HORMONAL PREPARATIONS
- VARIOUS
- IUPAC-like
- Corpora
- Human miRNA
- Mouse Genes
- MeSH Disease
- @neurIST Ontology
- GO Component
- GO Function
- GO Process

Toggle Abstracts

SCAIVIEW

Fulltext Query for search, or select saved Searches under blue down arrow

multiple sclerosis

Subcorpus Statistics Server Statistics

The following entities relating to 'multiple sclerosis' were found in 1244 documents.

- Human Genes / Proteins
- Chromosomal Location
- STS Marker
- non Normalized SNP
- Normalized SNP
- Normalized CRF SNP
- Drug Names
- ALIMENTARY TRACT AND METABOLISM
- ANTIINFECTIVES FOR SYSTEMIC INFECTIONS
- ANTINEOPLASTIC AND IMMUNOSUPPRESSANT DRUGS
- ANTIPARASITIC PRODUCTS, INSECTICIDES AND VERMIFUGES
- BLOOD AND BLOOD FORMING COMPOUNDS
- CARDIOVASCULAR SYSTEM
- DERMATOLOGICALS
- GENITO URINARY SYSTEM AND MUSCULO-SKELETAL SYSTEM
- + NERVOUS SYSTEM
- RESPIRATORY SYSTEM
- SENSORY ORGANS
- SYSTEMIC HORMONAL PREPARATIONS
- VARIOUS
- IUPAC-like
- Corpora
- Human miRNA
- Mouse Genes
- MeSH Disease
- @neurIST Ontology
- GO Component
- GO Function
- GO Process

- Chromosomal Locations
- IUPAC
- Epigenetics

1. Pharmacologic

PubMed 20568832

Multiple sclerosis (MS) is an autoimmune disease that directly relieves MS symptoms of patients. Nocebo-pair direct consequence of a less other symptoms. Central ne triptyclic antidepressants are randomized, placebo-control treatment perceived an improved pain syndrome. The antispasmodic medications are often unrecognized by clinic highlight shortcomings in clinical practice.

MeSH: Analgesics therapie

2. Tuberos sclerosis

PubMed 20845641 and Paediatrics, University of

BACKGROUND: Tuberos mental retardation. The aim primary two pupil of jaw tibia for six years. He also had a has been on Carbamazepine on the malar area of the face tomography showed multiple in the cortical white matter c clinic. CONCLUSION: Tube

MeSH: Analgesics; Non-Na Tuberos Sclerosis complic

3. Voltage-gated

PubMed 20298965

Besta, Milano, Italy.

Select Confidence: 1 2 3 4

From MEDLINE®, Pub  
PubMed Central®, datat  
U.S. National Library o  
PubMed PubMed

Select Confidence: 1 2 3 4 5

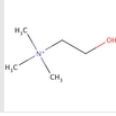

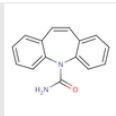
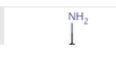
From MEDLINE®, PubMed®, and  
PubMed Central®, databases of the  
U.S. National Library of Medicine.  
PubMed PubMedCentral

multiple sclerosis

Subcorpus Statistics Server Statistics

The following entities relating to 'multiple sclerosis' were found in 1244 documents.

- Human Genes / Proteins
- Chromosomal Location
- STS Marker
- non Normalized SNP
- Normalized SNP
- Normalized CRF SNP
- Drug Names
- ALIMENTARY TRACT AND METABOLISM
- ANTIINFECTIVES FOR SYSTEMIC INFECTIONS
- ANTINEOPLASTIC AND IMMUNOSUPPRESSANT DRUGS
- ANTIPARASITIC PRODUCTS, INSECTICIDES AND VERMIFUGES
- BLOOD AND BLOOD FORMING COMPOUNDS
- CARDIOVASCULAR SYSTEM
- DERMATOLOGICALS
- GENITO URINARY SYSTEM AND MUSCULO-SKELETAL SYSTEM
- + NERVOUS SYSTEM
- RESPIRATORY SYSTEM
- SENSORY ORGANS
- SYSTEMIC HORMONAL PREPARATIONS
- VARIOUS
- IUPAC-like
- Corpora
- Human miRNA
- Mouse Genes
- MeSH Disease
- @neurIST Ontology
- GO Component
- GO Function
- GO Process

Select	Structure Image	Entity	Relative Entropy	Ref. Doc Count	Doc Count	Date Reported	Links
<input type="checkbox"/>		Choline	0.2579	17936	108	2010-11	SNP D SP
<input type="checkbox"/>		Iron	0.2267	90528	161	2010-11	D
<input type="checkbox"/>		Carbamazepine	0.2144	9757	81	2010-07	SNP D SP
<input type="checkbox"/>	No Structure Available	Calcium	0.2141	206066	217	2010-10	SNP D SP
<input type="checkbox"/>							

Select Confidence: 1 2 3 4 5

From MEDLINE®, PubMed®, and  
PubMed Central®, databases of the  
U.S. National Library of Medicine.  
PubMed PubMedCentral

# SCAIVIEW functionalities

- **Document View** - Displays all the documents retrieved based on the search query. Entity classes can be selected that you want to highlight. By default, Documents are displayed according the date (newest on the top).
- **Entity View** - Displays named entities under the column entities and are linked to corresponding abstracts
- **Export** - PMID and Entity tables can be exported to text files or excel sheets

The screenshot shows the SCAIVIEW interface with the search query "Human Genes / Proteins 'EGFR'". The results are displayed in a list view, sorted by date. The top result is a document titled "Short-term change in kidney function and risk of end-stage renal disease" by Wabnitz et al., published in 2012.

**BACKGROUND:** It is unclear what degree of change in the **eGFR** over a 1-year period indicates clinically significant progression, and whether the change adds additional information beyond that obtained by a single **eGFR** measurement.

**METHODS:** We included 568 397 adults who had at least two independent **eGFR** measurements (at least 8 months apart) during 1-year accrual period in Alberta, Canada. Change in kidney function (using the first and last **eGFR**) was defined by change in kidney function category with confirmation based on percent (%) change in **eGFR** (less  $\geq 25\%$ , the **eGFR**  $\geq 100$ ). The groups for change in kidney function were thus defined as: certain drop (drop in CKD category with  $\geq 25\%$  increase in the **eGFR**), uncertain drop (drop in CKD category with  $\geq 25\%$  decrease in the **eGFR**), stable (no change in CKD category), uncertain rise (rise in CKD category with  $\geq 25\%$  rise in the **eGFR**) and certain rise (rise in CKD category with  $\geq 25\%$  increase in the **eGFR**). Adjusted end-stage renal disease (ESRD) rates (per 1000 person-years) for each group of change in kidney function were calculated using Poisson regression. Adjusted rate of ESRD associated with change in kidney function, in reference to stable kidney function, were estimated.

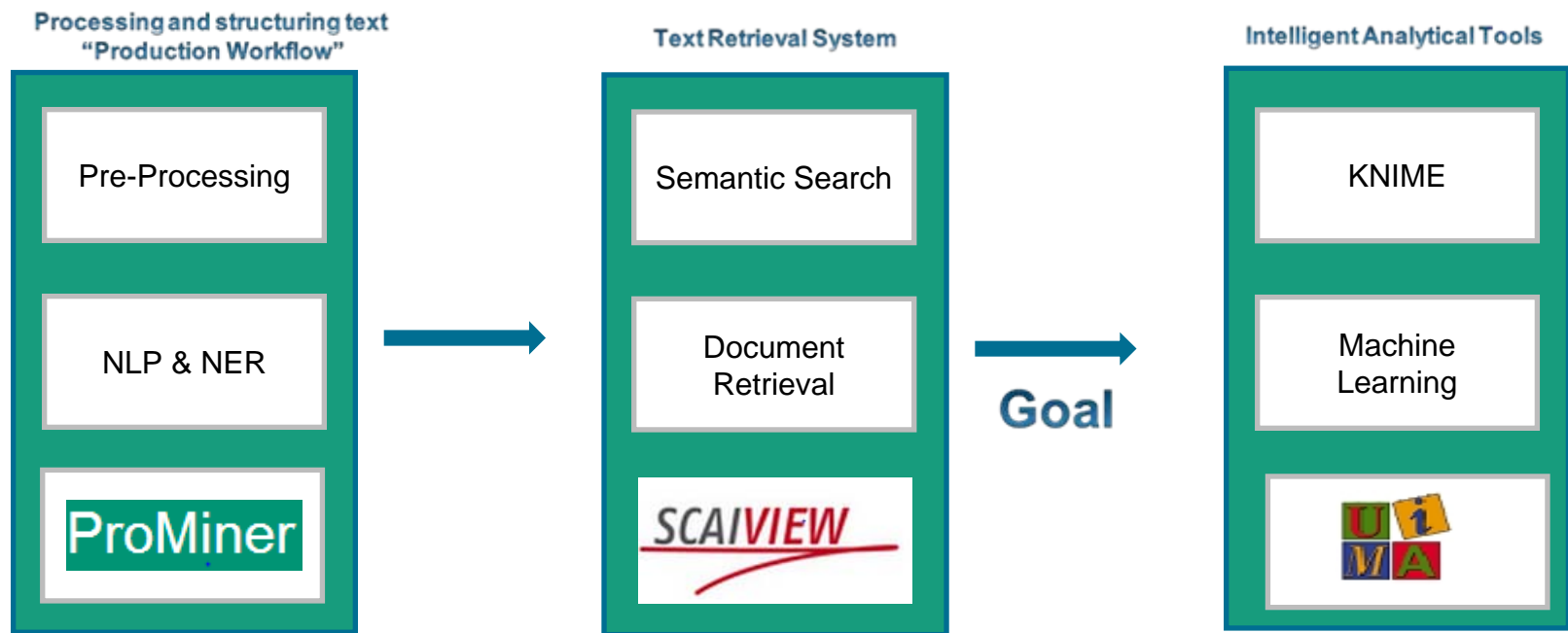
**RESULTS:** Among the 568 397 participants, 74.6% ( $n = 447 570$ ) had stable (no change in CKD category), 3.2% ( $n = 19 581$ ) had a certain drop and 3.7% ( $n = 22 171$ ) had a certain rise in kidney function. Participants who experienced a certain change in kidney function (both drop and rise) were older, more likely to be female, and had a higher prevalence of comorbidities, in comparison with those with stable kidney function. There were 1960 (2.7%) ESRD events over a median follow-up of 3.5 years. Compared with participants with stable kidney function, after adjustment for covariables, and the **eGFR** measurement, those with certain drop had 5-fold increased rate of ESRD (HR: 5.11; 95% CI: 4.58-5.71), whereas those with an uncertain drop had 2-fold increased rate (HR: 2.13; 95% CI: 1.84-2.47). After adjustment for the **eGFR** and covariables at the last visit, neither a certain nor uncertain drop in the **eGFR** was associated with an increased ESRD rate. The ESRD risk associated with the last **eGFR**  $\geq 100$ , adjusted for the slope over time, were 2.88 (95% CI: 2.35-3.51), 10.38 (95% CI: 8.69-13.87), 25.20 (95% CI: 21.95-44.32) and 147.96 (116.92-187.23) for categories 2, 3a, 3b and 4, respectively, in reference to category 1.

The screenshot shows the SCAIVIEW interface with the search query "Human Genes / Proteins 'EGFR'". The results are displayed in a table view, sorted by date. The table shows the following entities and their associated document counts and dates:

Entity	Doc Count	Date Reported
Neoplasms	17082	2014-10
Carcinoma Non-Small-Cell Lung	3651	2014-08
Breast Neoplasms	5599	2014-07
Lung Neoplasms	4938	2014-06
Adenocarcinoma	3287	2014-08
Carcinoma Squamous Cell	2135	2014-08
Carcinoma	3548	2014-08
Neoplasm Metastasis	3213	2014-10
Glioblastoma	1343	2014-05
Colorectal Neoplasms	1569	2014-06



# Motivation



# UIMA workflow and UI applications

- UIMA (Unstructured Information Management Architecture) is a software architecture for deploying and developing unstructured information management application.
- Originally developed by IBM, now open source.
- Unstructured information application may be defined as a software system designed to analyse large volume of unstructured information in order to discover, organise, and deliver knowledge to the end user
- Thus this architecture provides analytical platform by converting unstructured text to structured information .



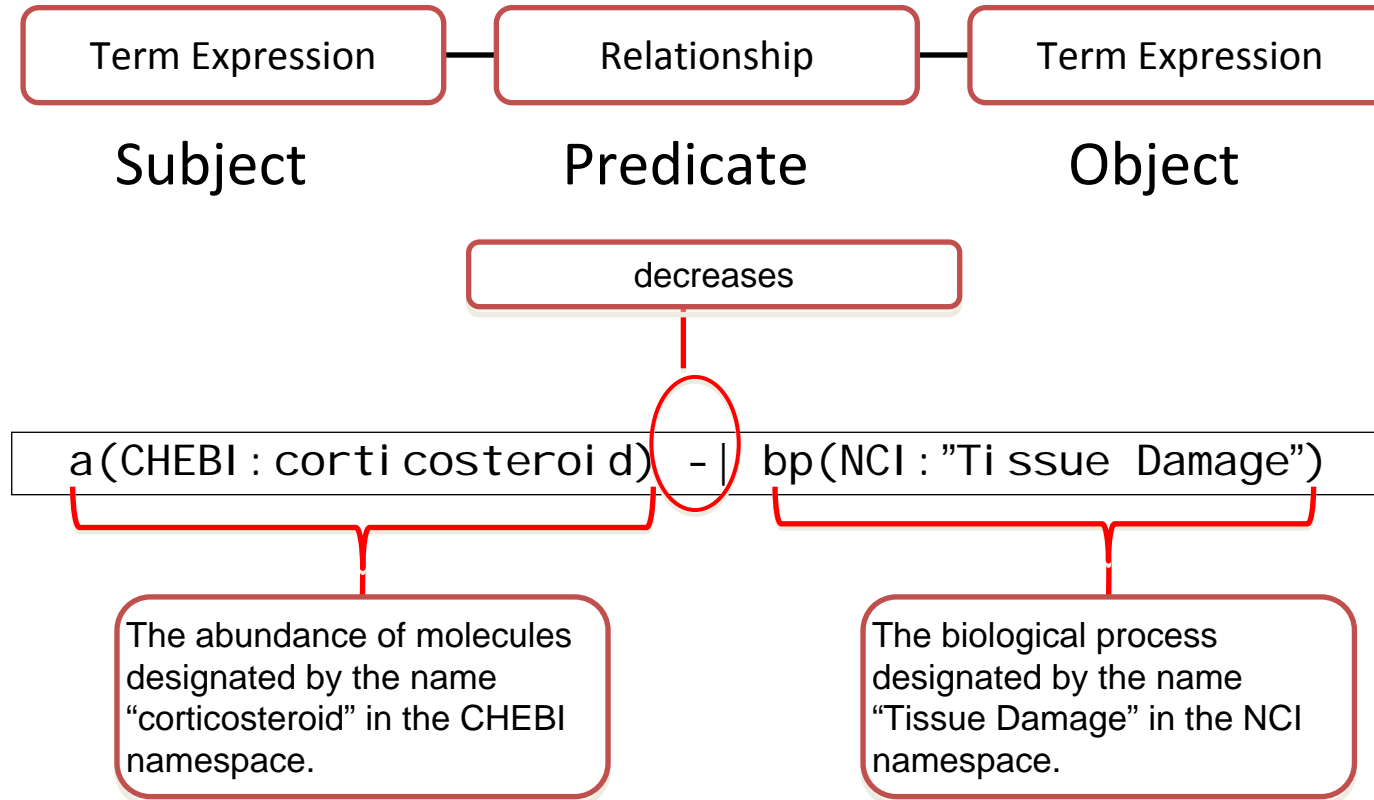


# UIMA based analysis at SCAI

- BEL like Statement Extraction
- Co-occurrence and Tri-occurrence based relationship extraction
- Machine Learning based relationship extraction
- Topic Modelling
- Term Frequency based Analysis
- ...



# Capturing Knowledge on Causes and Effects: OpenBEL





## OpenBEL: Capturing of Knowledge and “encoding” of data

Phosphorylation of **glycogen synthase kinase 3beta** at **Threonine, 668** **increases** the **degradation** of **Amyloid precursor protein**.

**p** (HGNC:**GSK3B**, **pmod** (P,T,668)) -> **deg** (**p** (HGNC:**APP**))

**BEL Functions**

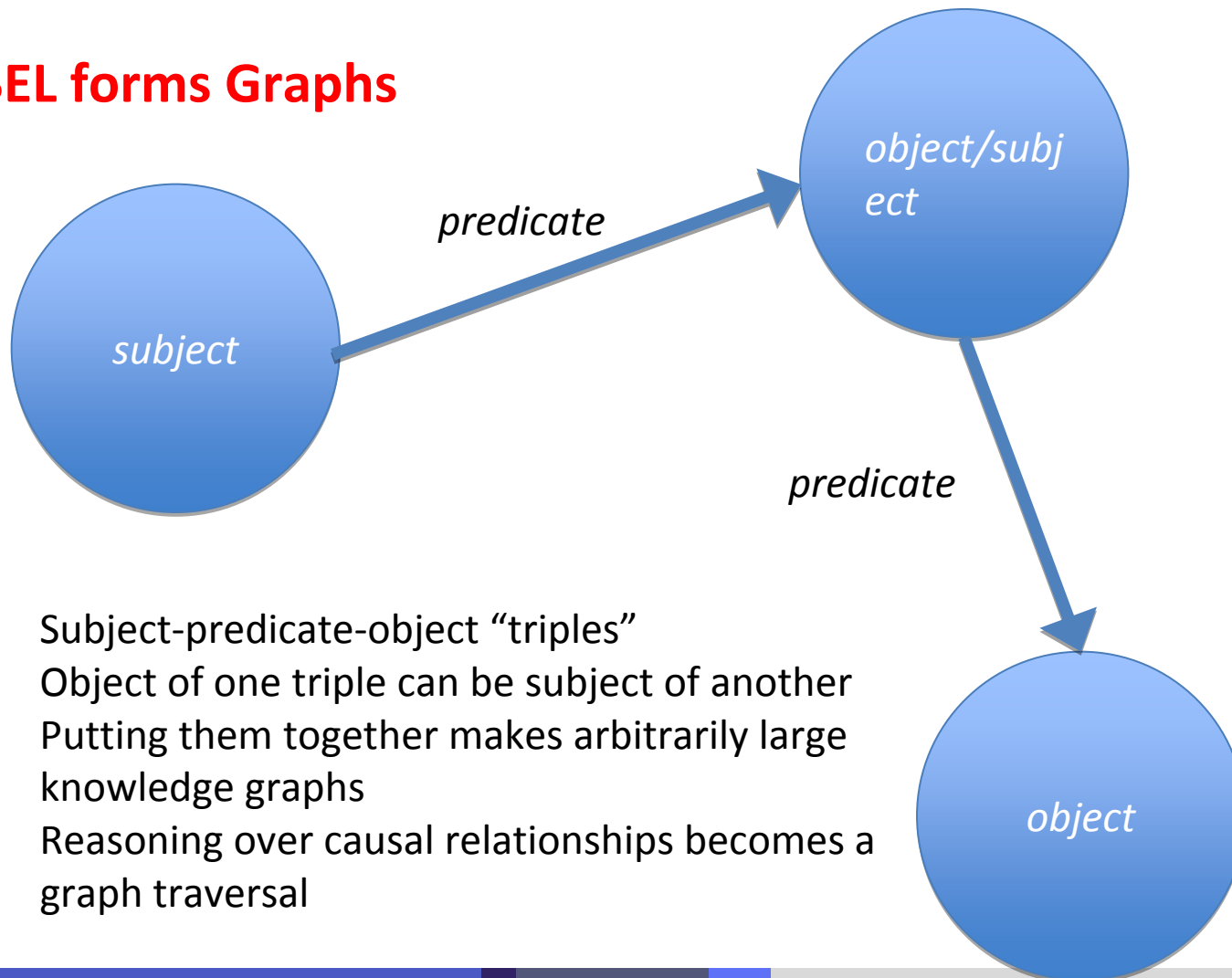
**Namespace Identifiers**

**Entity Definitions**





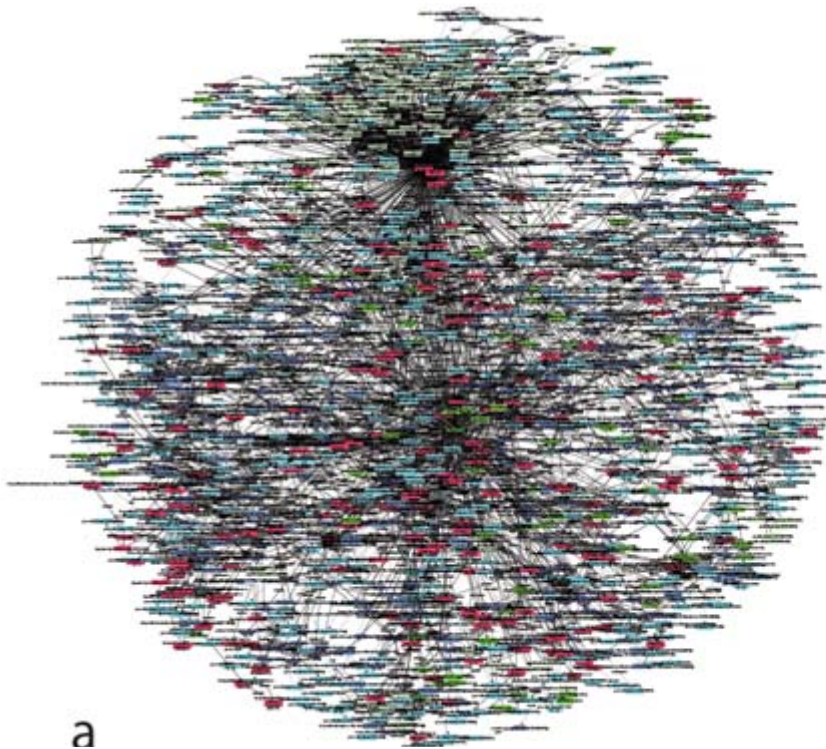
## BEL forms Graphs



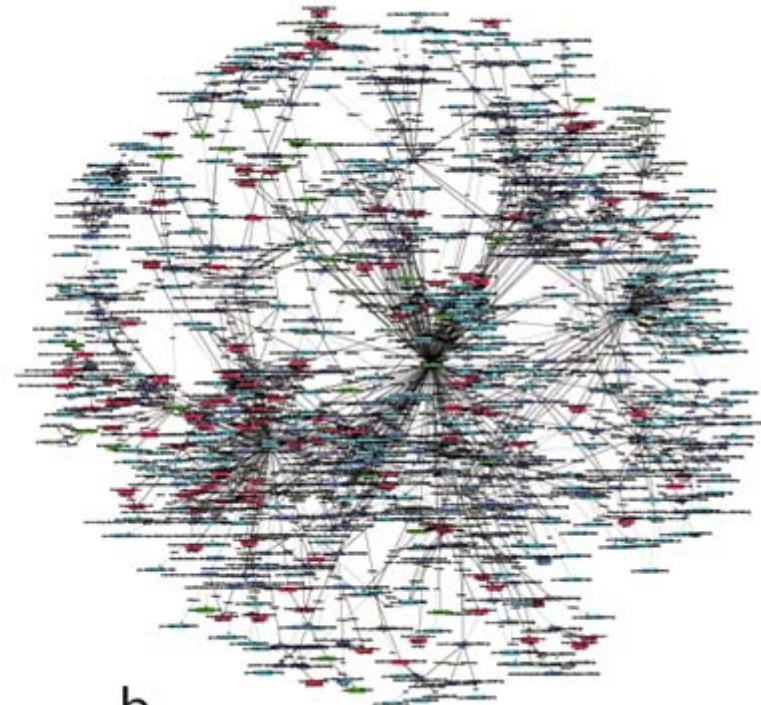
- Subject-predicate-object “triples”
- Object of one triple can be subject of another
- Putting them together makes arbitrarily large knowledge graphs
- Reasoning over causal relationships becomes a graph traversal



# The World's largest Computable Model for Alzheimer's Disease



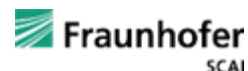
a



b



Kodamullil, Alpha Tom, et al. "Computable cause-and-effect models of healthy and Alzheimer's disease states and their mechanistic differential analysis." *Alzheimer's & Dementia* (2015).



# Project Goals of the Work of Gurnoor Singh

- Implement a solution integrating **document retrieval via SCAIView and analytical tools** which extracts information/knowledge based on a UIMA workflow
- A **generic** solution which works well on any analysis workflow(‘wrapper’)
- A **well distributed, flexible, and efficient solution** for multitasking
- Show application by, performing an exemplary analysis which measures the difference of **information gain** between abstract representation and full text representation of Biomedical journals

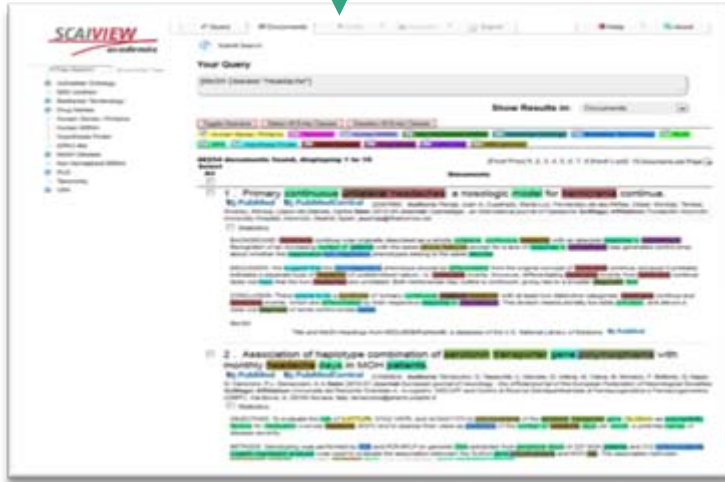


# Generic Workflow

Text Processing Daemon(s)



## Query Submission



## Document Selection

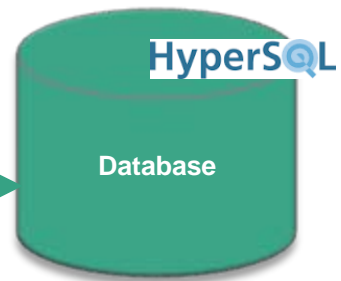


First Name	Gurnoor
Last Name	Singh
Email	gurnoor1990@gmail.com
Output format	archive
UMA workflow	BEL
<input type="button" value="Submit Workflow"/>	

Content	Number
Total Documents selected	5
Total Documents allowed to process per hour	1000
Number of empty documents	0
Number of documents exceeding in length	0
Number of extra documents removed	0
Number of documents processing in UMA	5

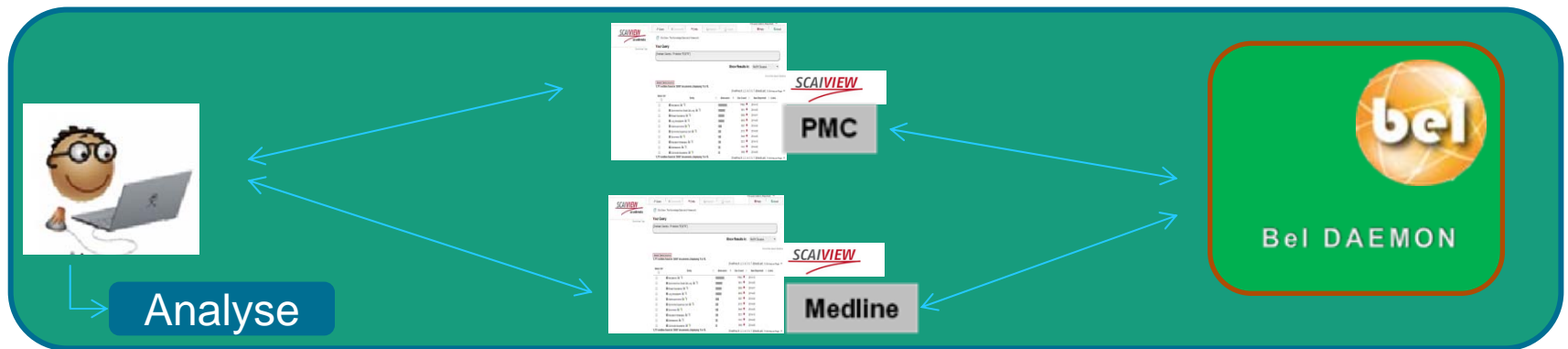
## Credentials & Job Details



## Email Notification & Download

# Exemplary Application

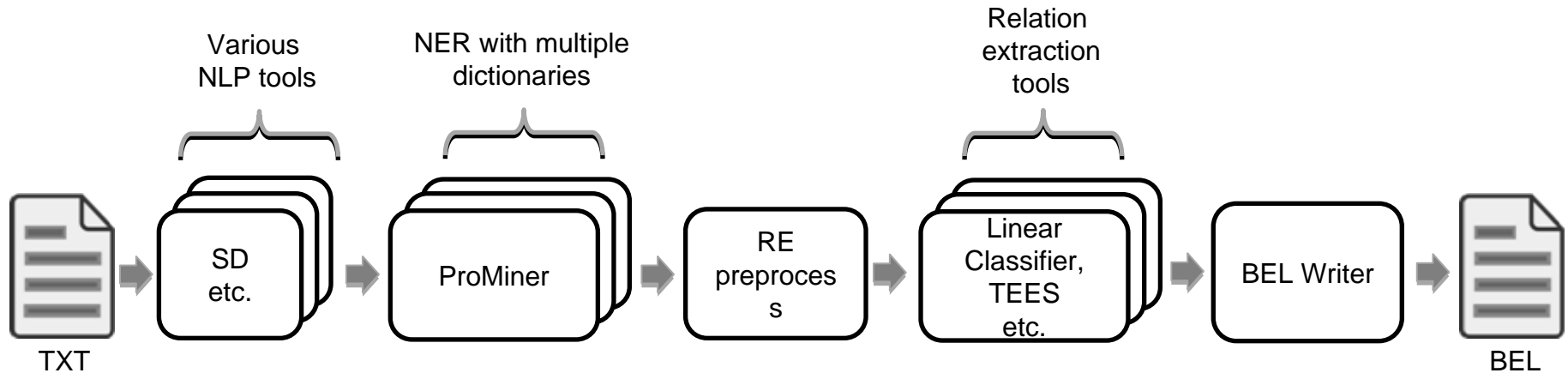
- **Biological Research Question** : “Is there a difference in **Information Gain** between abstract and full text of a document?”
- Information Gain can be measured as number of unique BEL like statements.
- Install daemon for analysing BEL Like Statement.



- **Corpus** : Collection of PMID as defined by user query in SCAIView



# BELIEF Workflow



**Unstructured  
Information Management  
Architecture**

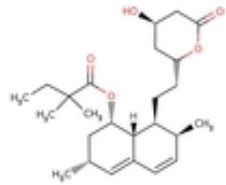
*An Apache Project*

# What text mining can deliver – NER & Normalisation

## 3. Simvastatin inhibits induction of matrix metalloproteinase-9 in rat alveolar macrophages

### Simvastatin:

*Simvastatin - A derivative of lovastatin and potent competitive inhibitor of 3-hydroxy-3-methylglutaryl coenzyme A reductase (hydroxymethylglutaryl COA reductases), which is the rate-limiting enzyme in cholesterol biosynthesis*



D

smoke extract.

**Central** 19299917 **Authors:** Kim, Sang Eun; Thanh Thuy, Tran Thi; Lee, Ji Hyun; Ro, Jai Youl; Bae, Young An; Kong, ; Oh, Yeon Mock; Lee, Sang Do; Lee, Yun Song **Date:** 2009-0 **Journal:** Experimental & molecular medicine **Affiliation:** y, Sungkyunkwan University School of Medicine, Samsung Biomedical Research Institute, Suwon 440-746, Korea.

MMP-9) may play an important role in emphysematous change in chronic obstructive pulmonary disease (COPD), one of the morbidity worldwide. We previously reported that simvastatin, an inhibitor of HMG-CoA reductase, attenuates emphysematous lungs of rats exposed to cigarette smoke. However, it remained uncertain how cigarette smoke induced MMP-9 and how smoke-induced MMP-9 expression in alveolar macrophages (AMs), a major source of MMP-9 in the lungs of COPD patients. signaling for MMP-9 induction and the inhibitory mechanism of simvastatin on MMP-9 induction in AMs exposed to isolated rat AMs, CSE induced MMP-9 expression and phosphorylation of ERK and Akt. A chemical inhibitor of MEK1/2 or ERK or Akt, respectively, and also inhibited CSE-mediated MMP-9 induction. Simvastatin reduced CSE-mediated MMP-9 induction and inhibition was reversed by farnesyl pyrophosphate (FPP) or geranylgeranyl pyrophosphate (GGPP). Similar to transferase or GGPP transferase suppressed CSE-mediated MMP-9 induction. Simvastatin attenuated CSE-mediated activation

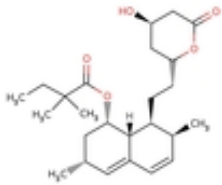
of RAS and phosphorylation of ERK, Akt, p65, I kappa B, and nuclear AP-1 or NF-kappa B activity. Taken together, these results suggest that simvastatin may inhibit CSE-mediated MMP-9 induction, primarily by blocking prenylation of RAS in the signaling pathways, in which Raf-MEK-ERK, PI3K/Akt, AP-1, and I kappa B-NF-kappa B are involved.

# What text mining can deliver – NER & Normalisation

## 3. Simvastatin inhibits induction of matrix metalloproteinase-9 in rat alveolar macrophages

### Simvastatin:

*Simvastatin - A derivative of lovastatin and potent competitive inhibitor of 3-hydroxy-3-methylglutaryl coenzyme A reductase (hydroxymethylglutaryl COA reductases), which is the rate-limiting enzyme in cholesterol biosynthesis*



D

smoke extract.

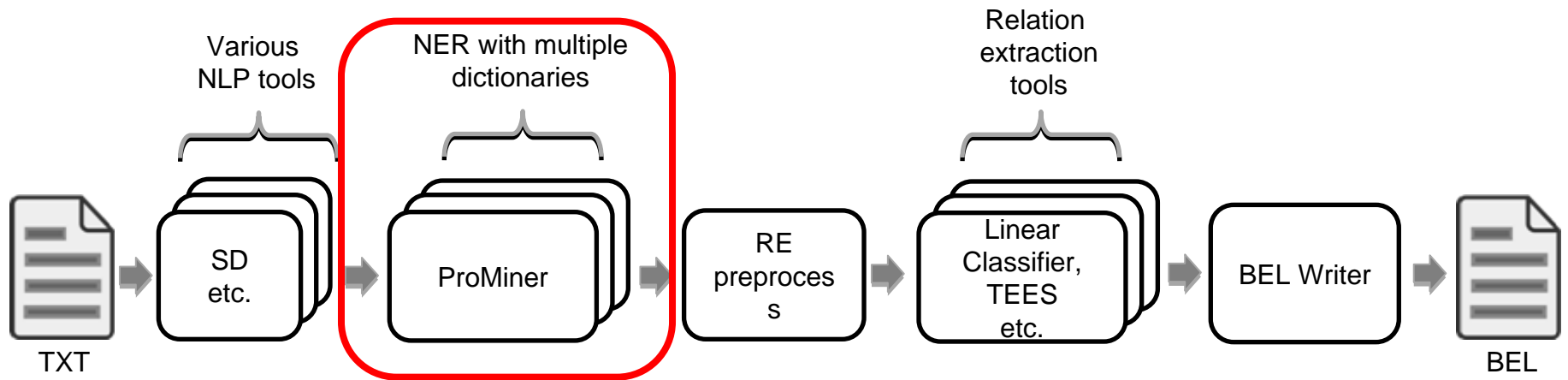
**Central** 19299917 **Authors:** Kim, Sang Eun; Thanh Thuy, Tran Thi; Lee, Ji Hyun; Ro, Jai Youl; Bae, Young An; Kong, Y.; Oh, Yeon Mock; Lee, Sang Do; Lee, Yun Song **Date:** 2009-0 **Journal:** Experimental & molecular medicine **Affiliation:** Y, Sungkyunkwan University School of Medicine, Samsung Biomedical Research Institute, Suwon 440-746, Korea.

MMP-9) may play an important role in emphysematous change in chronic obstructive pulmonary disease (COPD), one of the most morbidly worldwide. We previously reported that simvastatin, an inhibitor of HMG-CoA reductase, attenuates emphysematous change in the lungs of rats exposed to cigarette smoke. However, it remained uncertain how cigarette smoke induced MMP-9 and how cigarette-induced MMP-9 expression in alveolar macrophages (AMs), a major source of MMP-9 in the lungs of COPD patients. We investigated the signaling for MMP-9 induction and the inhibitory mechanism of simvastatin on MMP-9 induction in AMs exposed to cigarette smoke. In isolated rat AMs, CSE induced MMP-9 expression and phosphorylation of ERK and Akt. A chemical inhibitor of MEK1/2 or ERK or Akt, respectively, and also inhibited CSE-mediated MMP-9 induction. Simvastatin reduced CSE-mediated MMP-9 induction and inhibition was reversed by farnesyl pyrophosphate (FPP) or geranylgeranyl pyrophosphate (GGPP). Similar to FPP or GGPP transferase suppressed CSE-mediated MMP-9 induction. Simvastatin attenuated CSE-mediated activation

of RAS and phosphorylation of ERK, Akt, p65, IκappaB, and nuclear AP-1 or NF-κappaB activity. Taken together, these results suggest that simvastatin may inhibit CSE-mediated MMP-9 induction, primarily by blocking prenylation of RAS in the signaling pathways, in which Raf-MEK-ERK, PI3K/Akt, AP-1, and IκappaB-NF-κappaB are involved.

- Recall and Precision rates are between 70% and 90% for biomedical NER
- **Recall:** How many of the existing names does the system detect
- **Precision:** How many of the detected names are correct

# BELIEF Workflow



# Current dictionaries included

Entity class	Resources	BEL namespace
Human Genes/Proteins	EntrezGene/ Uniprot	HGNC
Mouse Genes/Proteins	EntrezGene/ Uniprot	MGI
Rat Genes/Proteins	EntrezGene/ Uniprot	RGD
Protein family names	OpenBEL	PFH
Protein complex names	OpenBEL	NCH
Protein complex names	Gene Ontology	GOCCTER M
Chemical names	OpenBEL	SCHEM
Chemical names	ChEBI	CHEBI
Chemical names	ChEMBL	SCHEM
Disease names	MeSH	MESHD
Anatomy names	MeSH	MESHA

# Use case: relation between small molecules (mainly protein inhibitors) and their targets

Dictionary	Recall rate initial version	Recall rate application adapted
Genes/Protein: (HGNC)	80 %	93 %
Chemical compounds: ChEBI	15 %	66 %
Chemical compounds: SCHEM	30 %	75 %
Chemical compounds: ChEBI + SCHEM+ ChEMBL	not determined	91 %
Selventa-human-complex	40 %	46 %
GO-Complex	not determined	64 %
Selventa-human-complex + Complex	not determined	82 %
GO-Function	22 %	not determined
Selventa-human-families	8 %	77 %

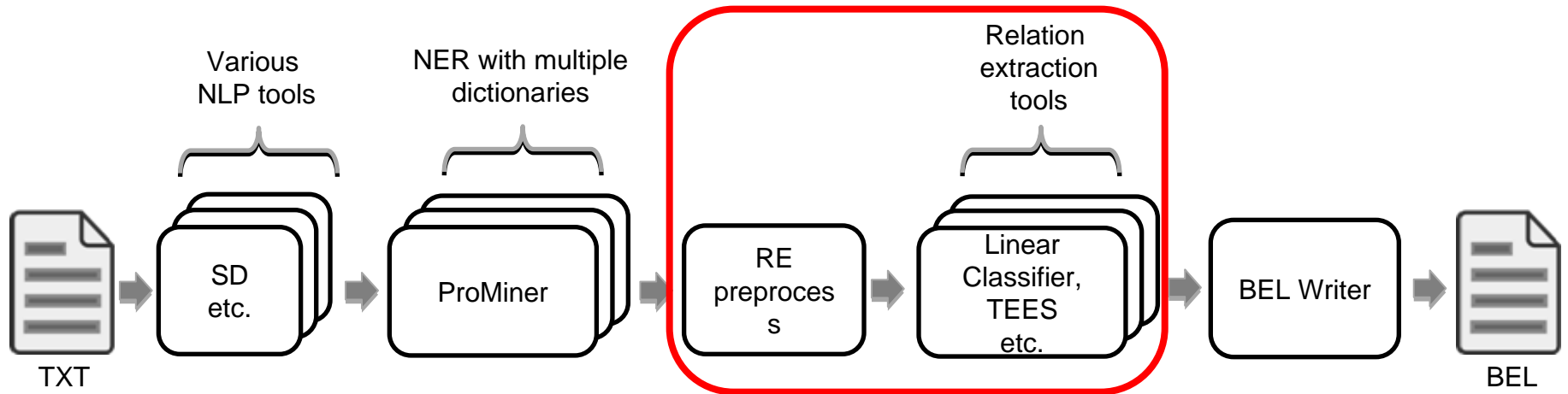
# Recognition and normalization of terminology

- Normalisation is needed!
- Use external and internal (OpenBEL) resources for named entity recognition (Mapping!)
- Combine various resources
- Adapt terminology to use cases (OpenBel namespaces provide no synonyms)
- Offer curators the annotation of different concepts

**For relation extraction high recall is a precondition!!!**

---

# BELIEF Workflow





# Relation Extraction

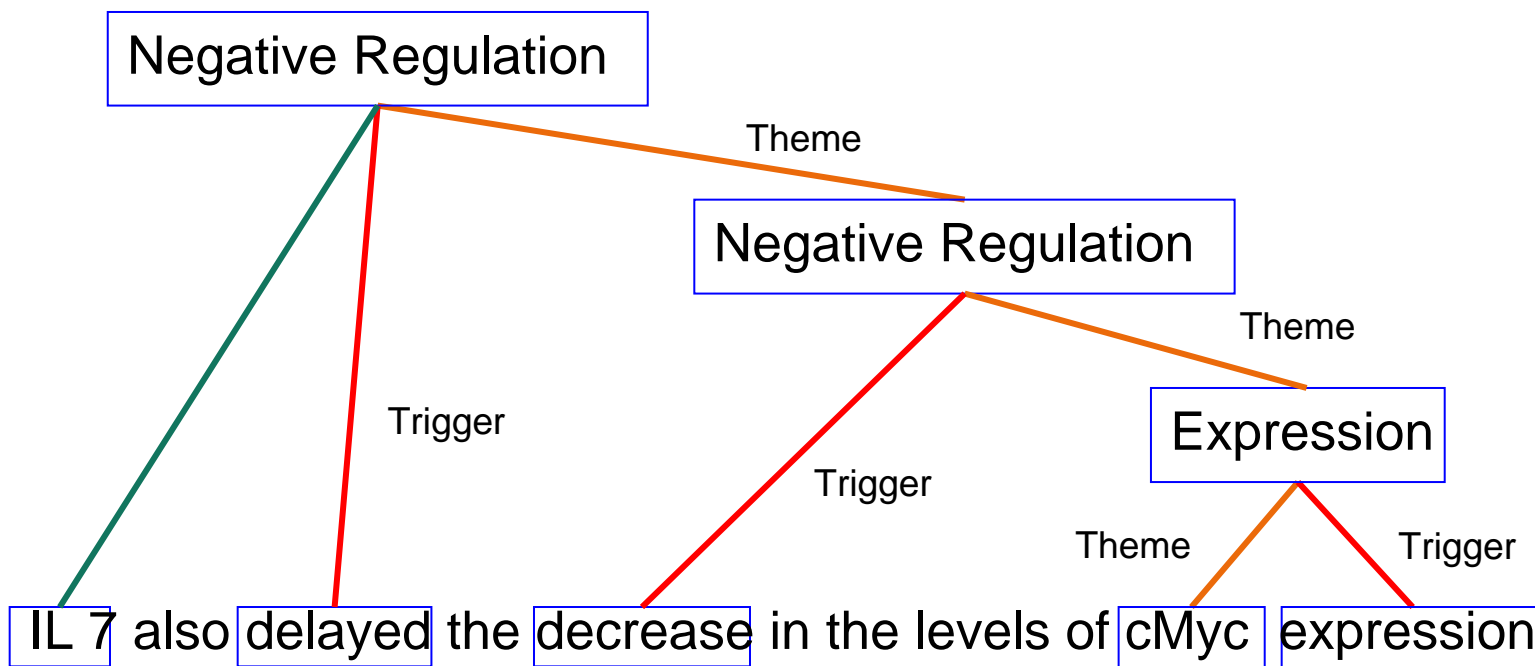
Two kinds of relationship extraction tools are available which are tested and compared on common benchmark sets:

- The BioNLP shared tasks deliver a very detailed annotation for relationship extraction similar to the information needed for BEL

# Relation Extraction

Two kinds of relationship extraction tools are available which are tested and compared on common benchmark sets:

- The BioNLP shared tasks deliver a very detailed annotation for relationship extraction similar to the information needed for BEL



# Relation Extraction

Two kinds of relationship extraction tools are available which are tested and compared on common benchmark sets:

- The BioNLP shared tasks deliver a very detailed annotation for relationship extraction similar to the information needed for BEL

- Simpler binary classification:

IL 7 also delayed the decrease in the levels of cMyc ....

IL7 – cMyc Relation: Yes

Classifies if a relation between 2 entities is existing but gives no information about the direction or type

# Technology – Performance

- NLP (Sentence Detection ~6% error) 94
- NLP (Tokenization ~8% error) 86
- NER (Different Classes ~15% error ) 73
- RelationExtraction (Multi-step ~25% error) 54

94

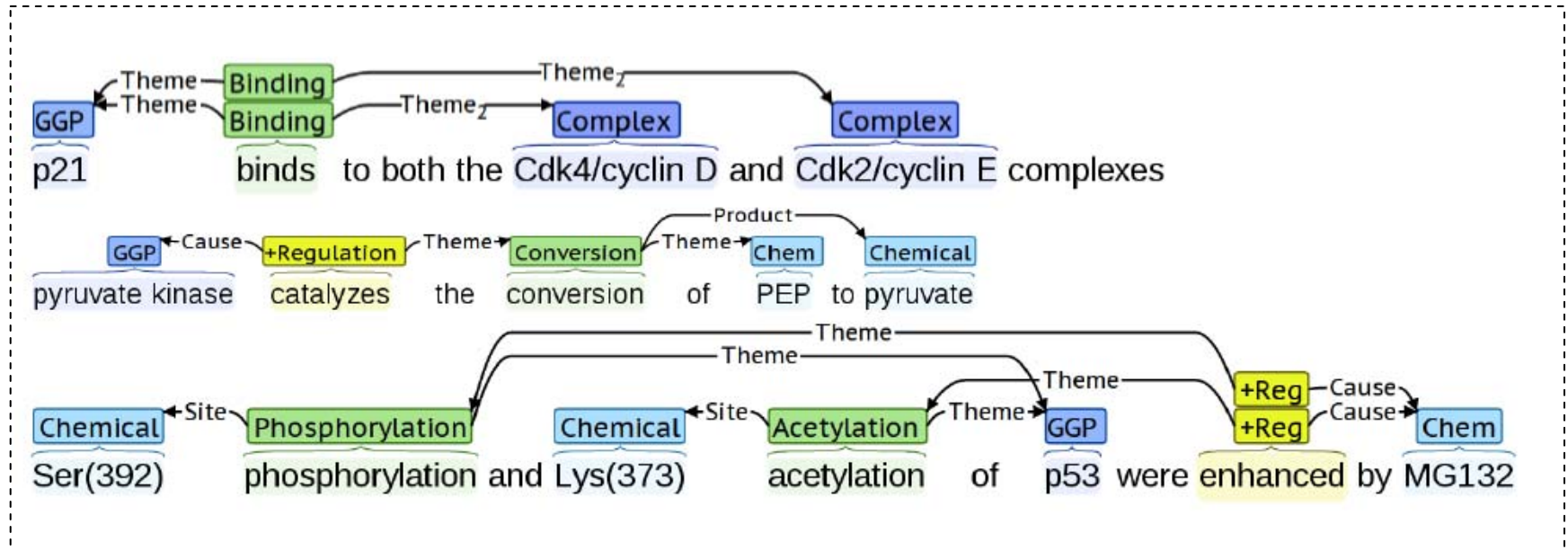
86

73

complexity

Propagated error!

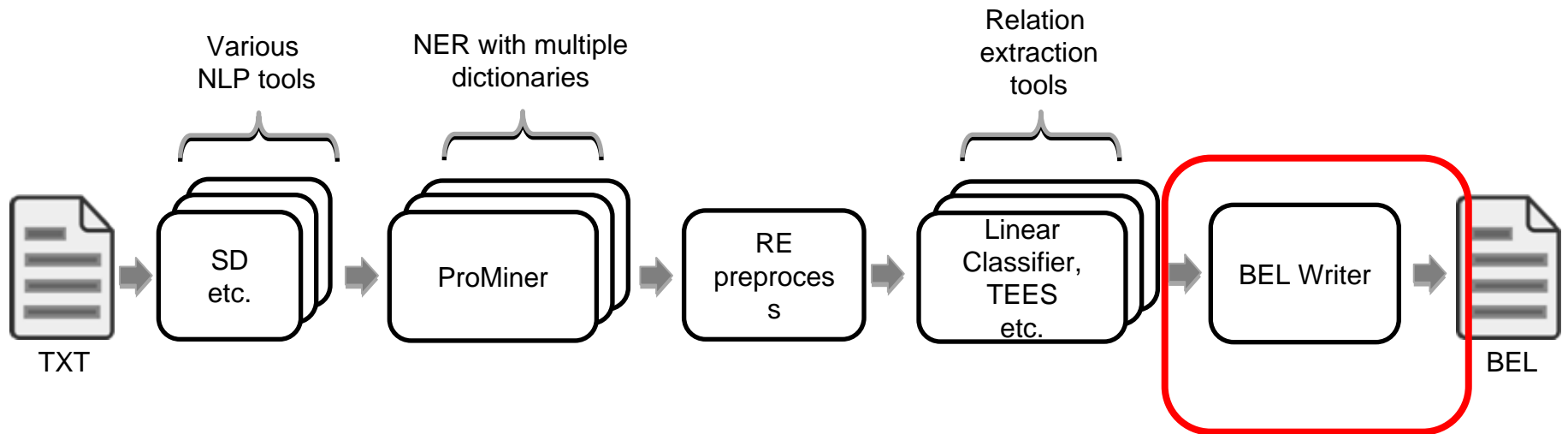
# What text mining can deliver: Relation Extraction – Example from BioNLP shared Task



- Recall ~30% and Precision ~50% for regulation events – binary classification has higher recall and precision rates

That seems not very promising but many relations might be redundant!

# BELIEF Workflow



# BioNLP Shared Task to openBEL conversion<sup>1</sup>

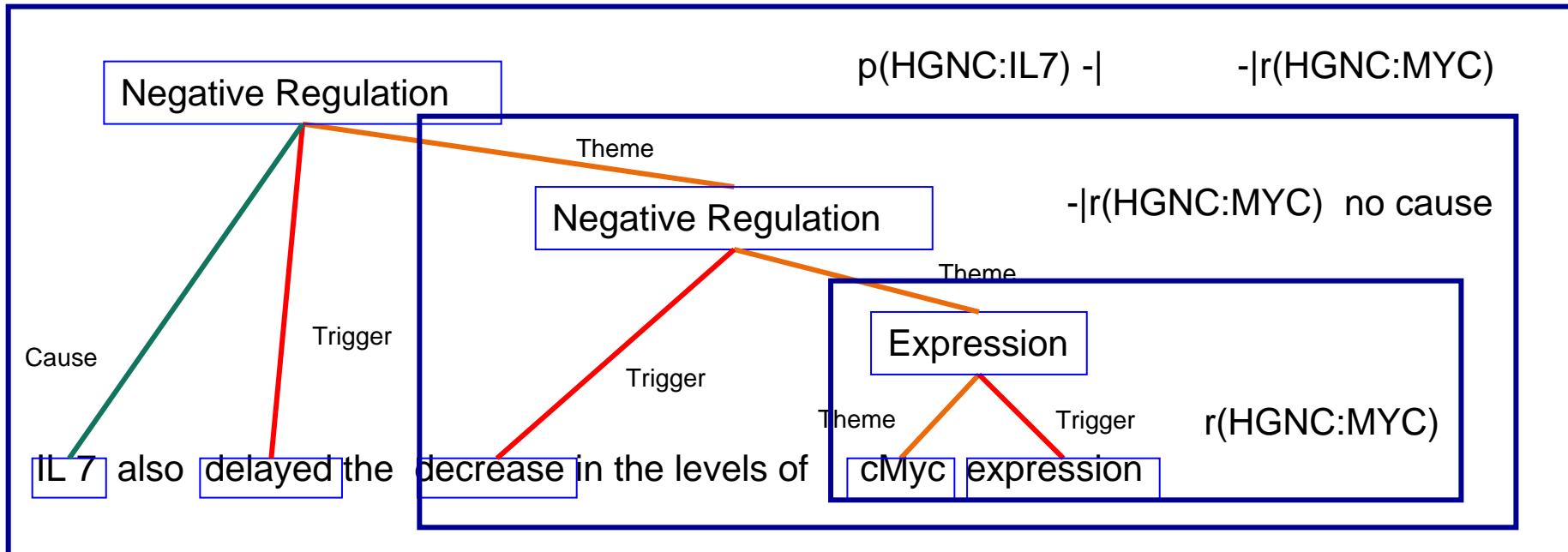
We implemented a rule set translating BioNLP SharedTask to BEL

Abundances	p(),...	Positive Regulation	->
Modifications	pmod()	Negative Regulation	-
Catabolism	deg()	Regulation/Association	--
Location	tloc()		
Binding/Complex	complex()		

<sup>1</sup>BEL networks derived from qualitative translations of BioNLP Shared Task annotations.

The Association for Computational Linguistics (ACL) Sofia 2013  
Fluck, J.; Klenner, A.; Madan, S.; Ansari, S.; Bobic, T.; Hoeng, J.; Hofmann-Apitius, M.; Peitsch, M.;

# Relation extraction example result



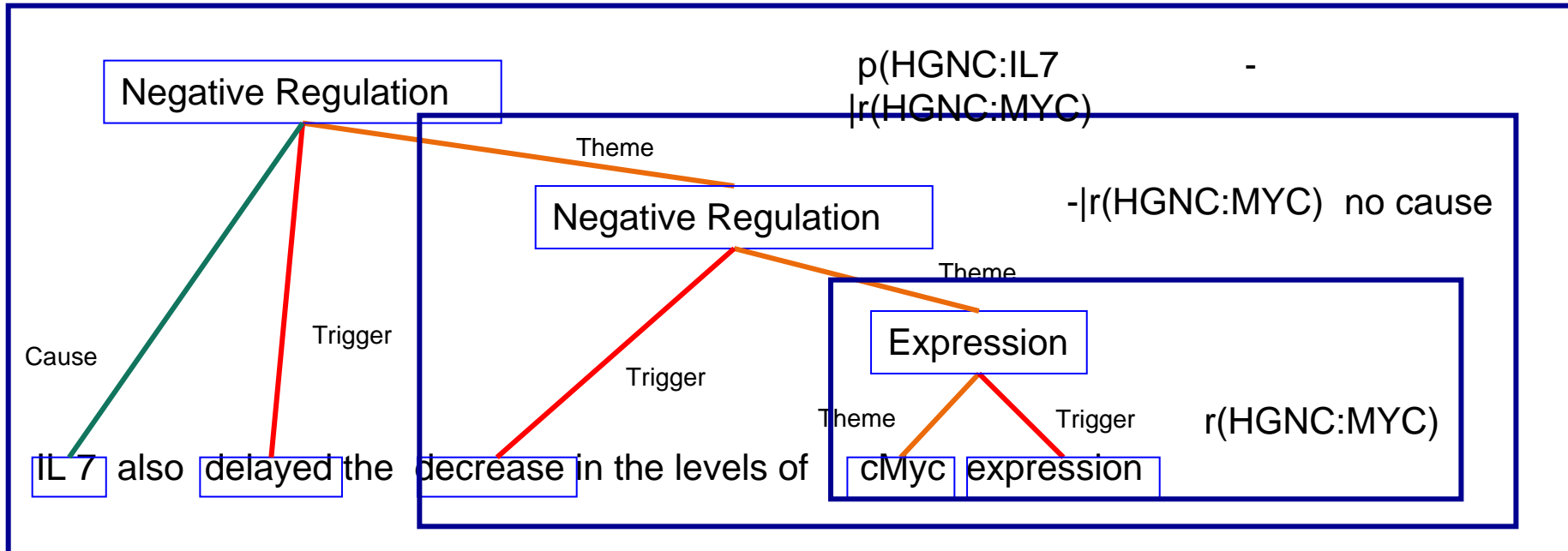


# Relation extraction example result

Automatic extension to full statements in workflow:

Fixme -|r(HGNC:MYC)

p(HGNC:IL7) -| Fixme -|r(HGNC:MYC)



# Relation extraction example result

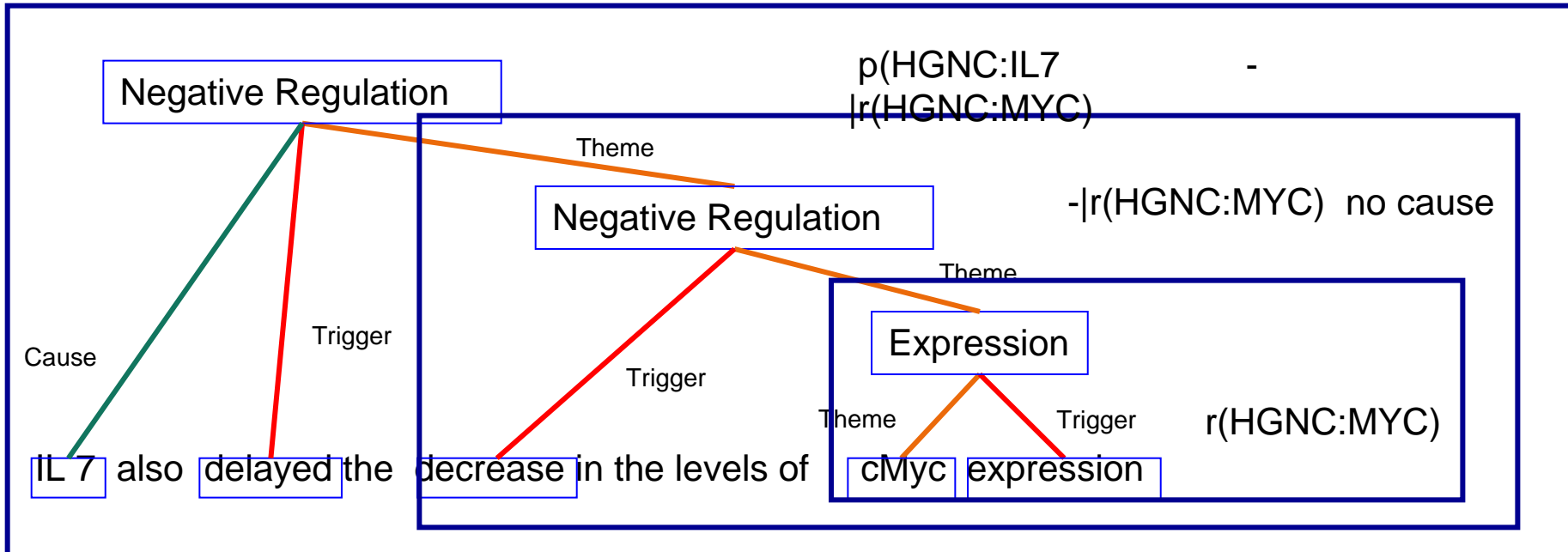
Automatic extension to full statements in workflow:

Fixme -|r(HGNC:MYC)

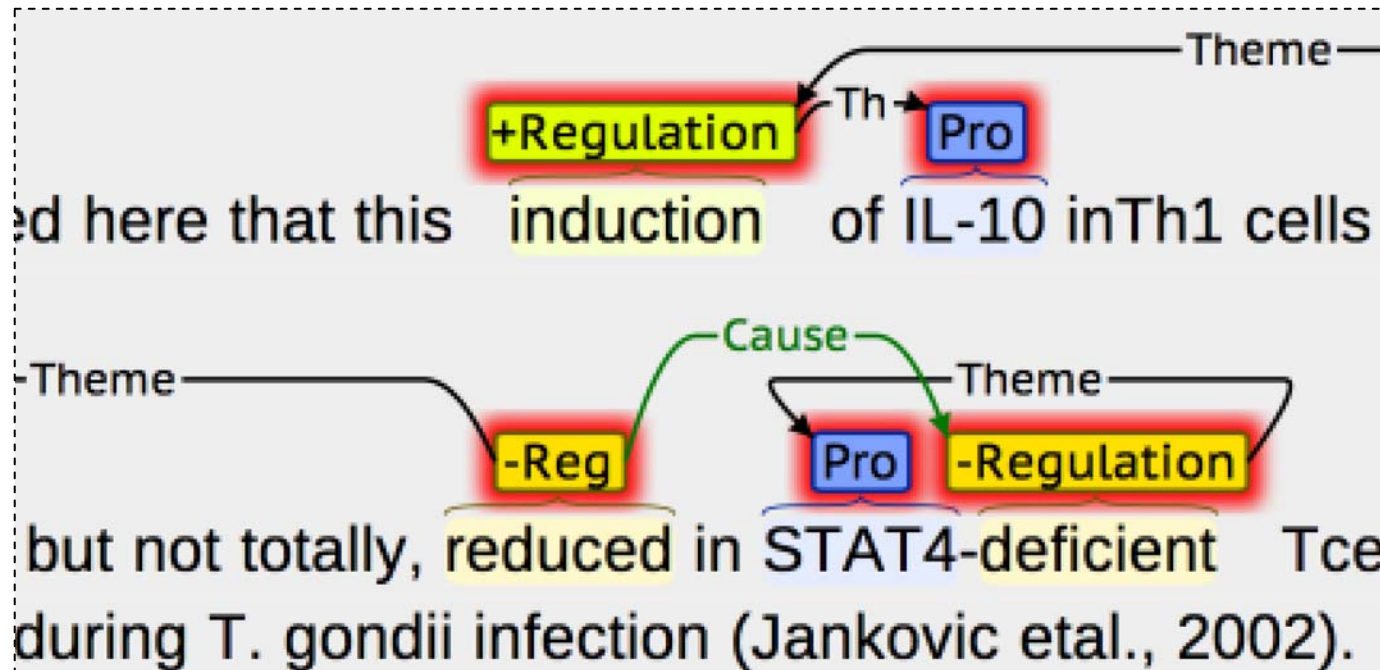
Binary classification:

p(HGNC:IL7) -| Fixme -|r(HGNC:MYC)

p(HGNC:IL7) -- r(HGNC:MYC)



# What text mining currently does not deliver: Interpretations



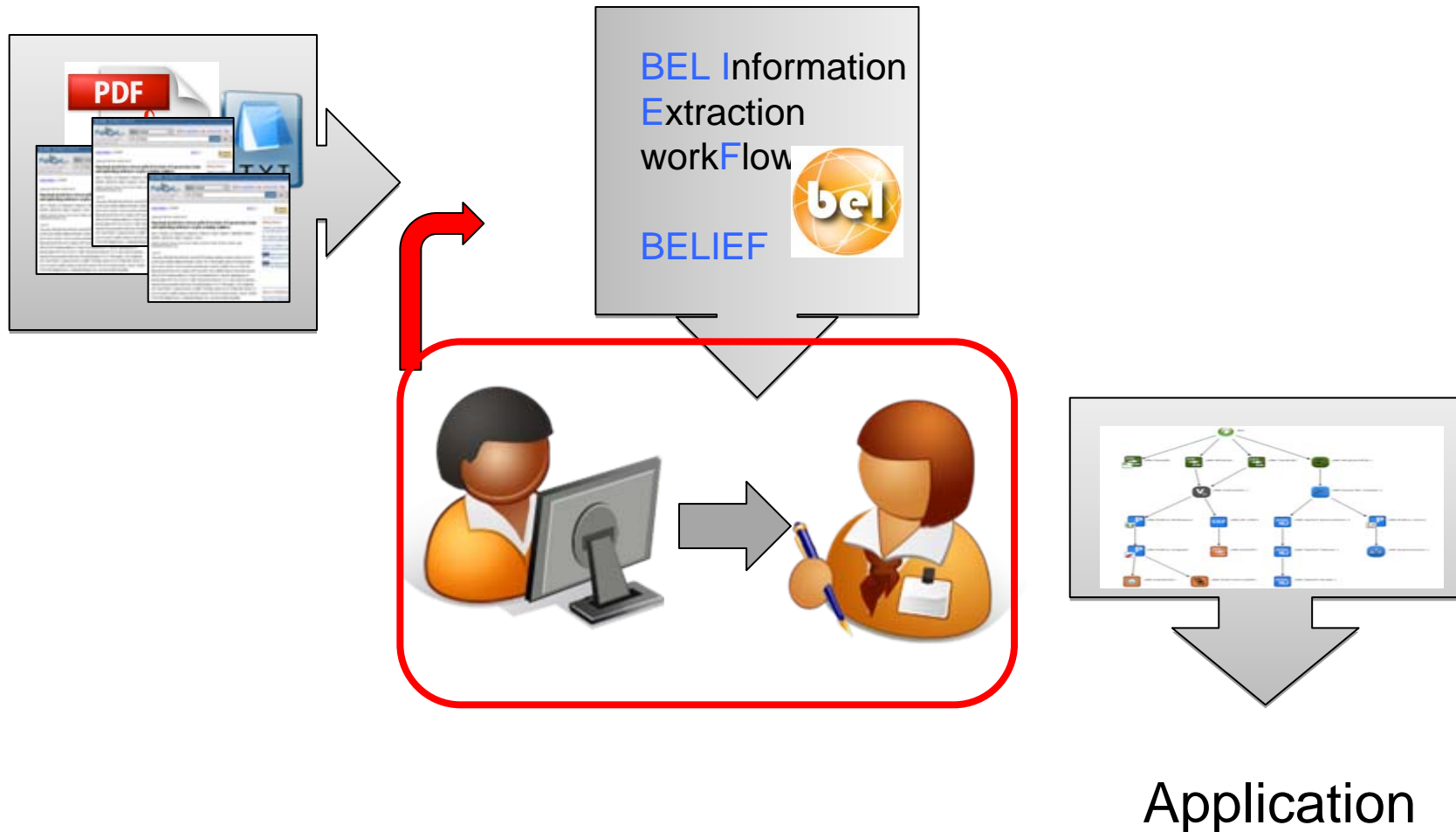
- Generated Statement:

- $p(\text{FIXME}) \neg p(\text{HGNC:STAT4}) \neg p(\text{FIXME}) \rightarrow p(\text{HGNC:IL10})$

- Manual Statement:

- $p(\text{HGNC:STAT4}) \rightarrow p(\text{HGNC:IL10})$

# Semi automatic BEL Knowledge Extraction Pipeline



# BeliefDashboard

Projects management and multiple document upload

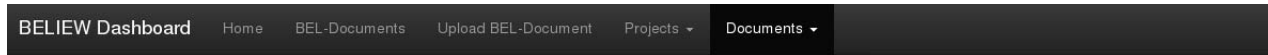
- Create, delete and list functionality



## Create Project

Name

Description



## Document Lists

- List
- Create

Id	Pubmed Id	Title	Created at	Processed at	Curate	Show details	Delete
1	2234234	Myc-Max-Mad network	2014-07-16T11:41:19Z	Queued	---	<a href="#">show details</a>	
2	12165281	Smoking causes a dose-dependent increase in granulocyte-bound L-selectin.	2014-07-16T11:41:19Z	Queued	---	<a href="#">show details</a>	

Showing 1 to 2 of 2 entries

– Previous 1 Next –

# BELIEF Dashboard Curation Interface

Detected concepts:  
-Highlighted text entity in mouse over  
-Provide a fast overview of all entities in the evidence

The screenshot displays the BELIEF Dashboard Curation Interface. At the top, a text snippet reads: "SIRT1 knockdown caused a slight yet highly significant decrease in LXRalpha expression in both fed and fasted livers." Below this, a smaller text block states: "However, SIRT1 knockdown decreased expression of PGC-1beta only the fasted state." Evidence is noted as "Evidence: 1/2" with a "Delete evidence" button. A sidebar on the right, titled "Detected concepts", lists: SIRT1 (HGNC:SIRT1, MGI:Sirt1, :Sirt1), LXRalpha (MGI:Nr1h3, :Nr1h3), and livers (MeSH\_A:Liver). Below the text, a section labeled "' low' confidence:" contains a table for context annotations. The table has columns for ID, BEL statement, and checkboxes for Edit, Delete, and Export. The first row shows ID 1 (id:976) with the BEL statement "p(FIXME) -| (p(HGNC:SIRT1) -| r(MGI:Nr1h3))". The second row shows ID (id:1895) with "BodyRegion" and "Liver" in the BEL statement column. Below the table, there are input fields for "CellLine" and "Enter annotation value".

ID	BEL statement	Edit	Delete	Export
1 (id:976)	p(FIXME) -  (p(HGNC:SIRT1) -  r(MGI:Nr1h3))	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
(id:1895)	BodyRegion   Liver	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

BEL statement

Context annotation

Edit/Delete/Export

# BEL Editor – Update document information

## Document information

Pubmed Id:

**PMID:** [11350768](#)

**Title:** Nicotine infusion alters leptin and uncoupling protein 1 mRNA expression in adipose tissues of rats.

**Journal:** American journal of physiology. Endocrinology and metabolism; Vol. 280; Iss. 6

**Authors:** K Arai, K Kim, K Kaneko, M Iketani, A Otagiri, N Yamauchi, T Shibasaki

**Published:** Jun/2001

BELIEW Dashboard

Curate document Arai 200

Title: Nicotine infu...  
Abstract:  
**Nicotine infu...**  
[View full text](#)

An J Physiol Endocrinol Metab. 280: E867-E876, 2001. We attempted to clarify whether leptin and uncoupling protein 1 (UCP1) are involved in the action of nicotine on the energy balance. Male Wistar rats were infused subcutaneously with nicotine (12 mg/kg/21 dday21) for 4 or 14 days.

Evidence: 1/11

Statements with "low" confidence:

1 (id:3299): (s[SCHEM:Nicotine])--r[HGNC:UCP1]

(id:6341): Species: Rattus norvegicus

(id:6911): BodyRegion: Adipose Tissue

CellLine:

Statements with "very low" confidence:

2 (id:3298): (s[SCHEM:Nicotine])--r[SCHEM:Leptin]

(id:6338): Species: Rattus norvegicus

(id:6912): BodyRegion: Adipose Tissue

CellLine:

Export

**Detected concepts**

Nicotine  
SCHEM:Nicotine  
Chemspid:Nicotine  
CHEBI:nicotine  
CHEBI:\_S\_nicotine  
leptin  
SCHEM:leptin  
MGI:Leptin  
HGNC:LEP  
uncoupling protein  
HGNC:UCP1  
uncoupling protein 1  
mRNA  
CHEBI:messenger RNA  
tissues  
MeSH:A:Tissues  
rats  
SubTaxonomy:Rattus norvegicus

**Document information**

Pubmed Id:

**PMID:** [11350768](#)

**Title:** Nicotine infusion alters leptin and uncoupling protein 1 mRNA expression in adipose tissues of rats.

**Journal:** American journal of physiology. Endocrinology and metabolism; Vol. 280; Iss. 6

**Authors:** K Arai, K Kim, K Kaneko, M Iketani, A Otagiri, N Yamauchi, T Shibasaki

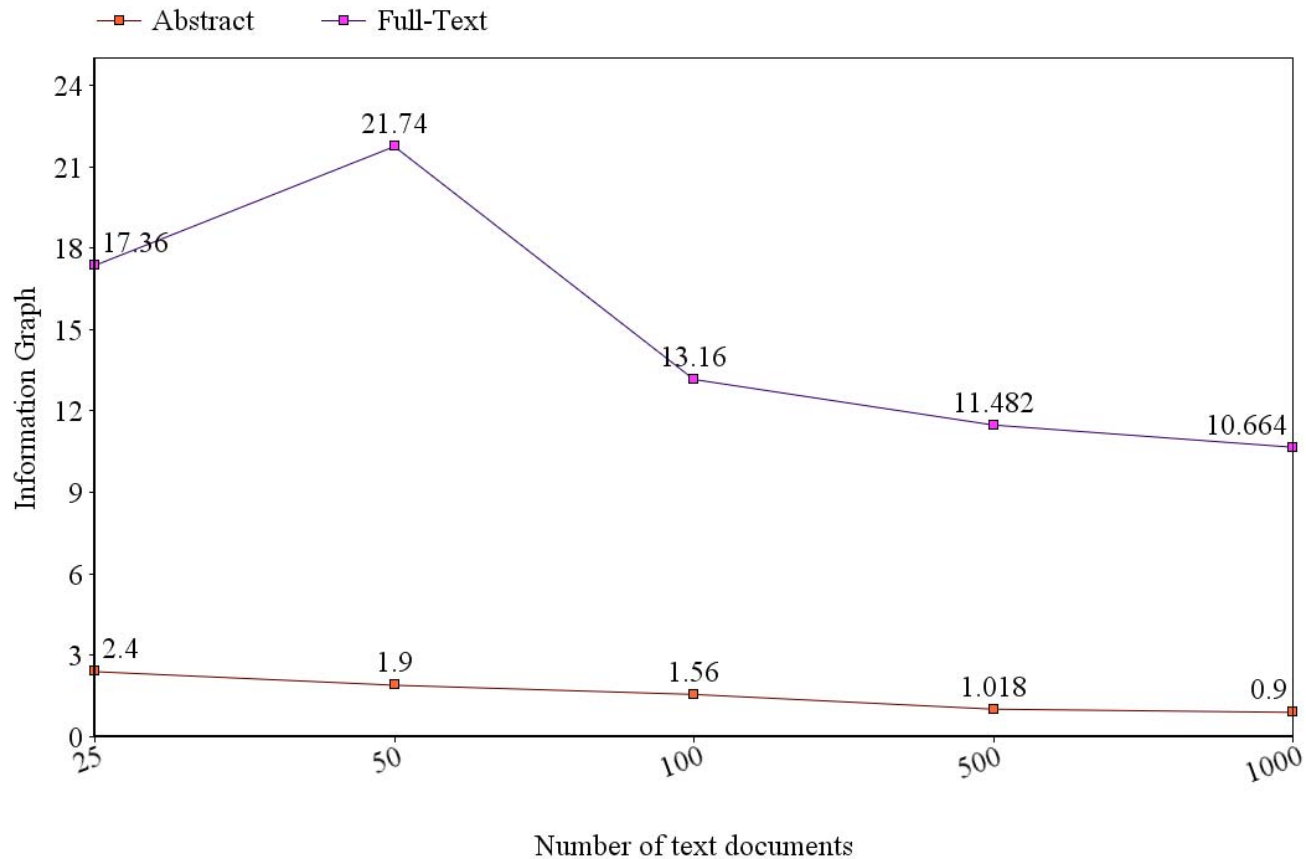
**Published:** Jun/2001

# Back to Gurnoor Singh: Experimental Setup

- Query disease under study: Alzheimer's Disease (AD)
- A total 10 jobs selecting top 25, 50, 100, 500, and 1000 text documents from both the SCAIView systems were exported to BEL processing daemon.
- BEL documents retrieved via SCAIview were analyzed using KAM navigator and Cytoscape
- Biological networks were further narrowed down to Protein-Protein Interactions networks

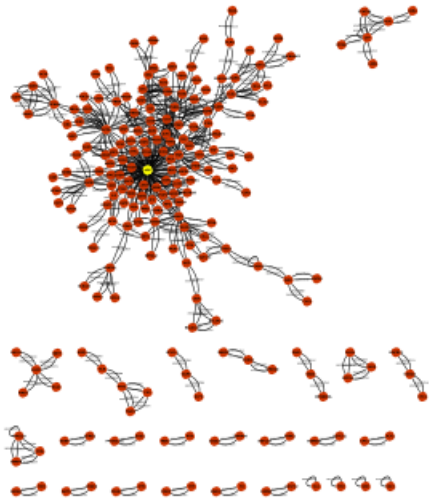


# Information Graphs



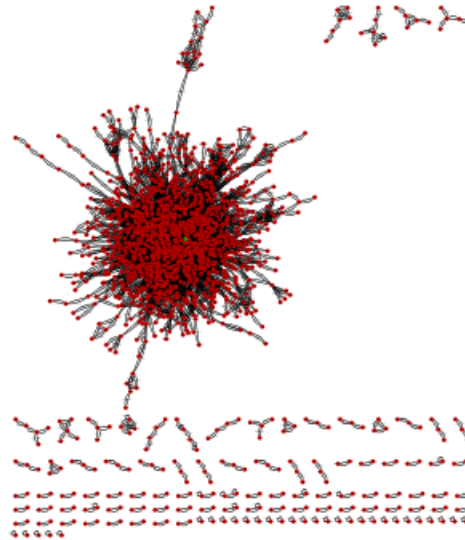
# Graph Topology

## 500- Abstract



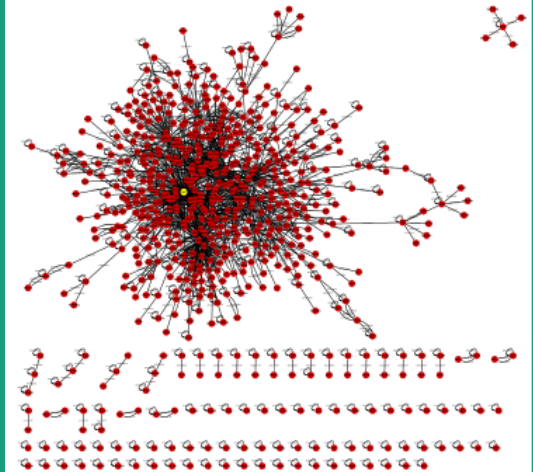
Documents: 500  
Nodes: 215  
Edges: 509

## 500- Full Text



Documents : 500  
Nodes : 1490  
Edges : 5741

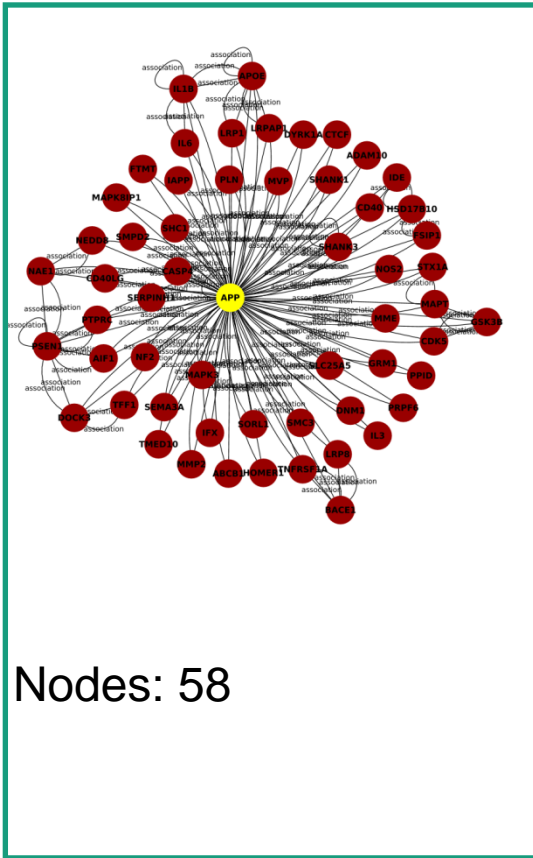
## Manually Curated



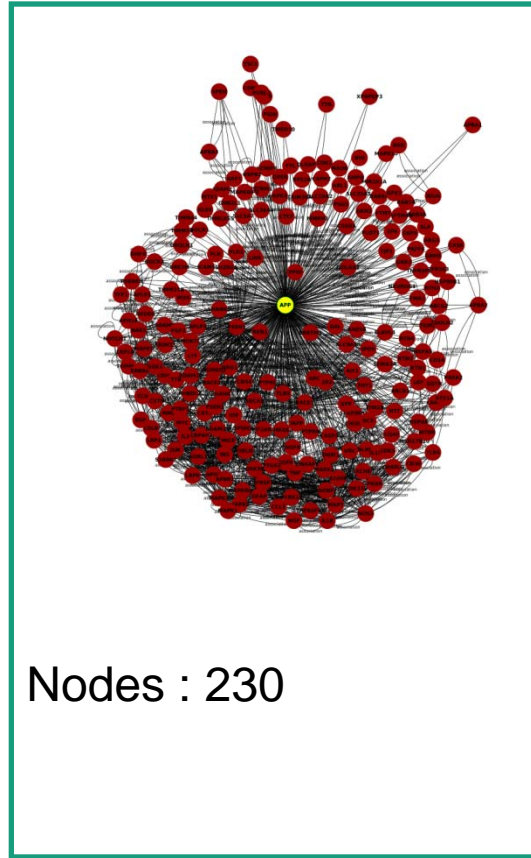
Documents : 500  
Nodes : 553  
Edges : 3525

# Plotting APP (yellow) and its First Neighbors

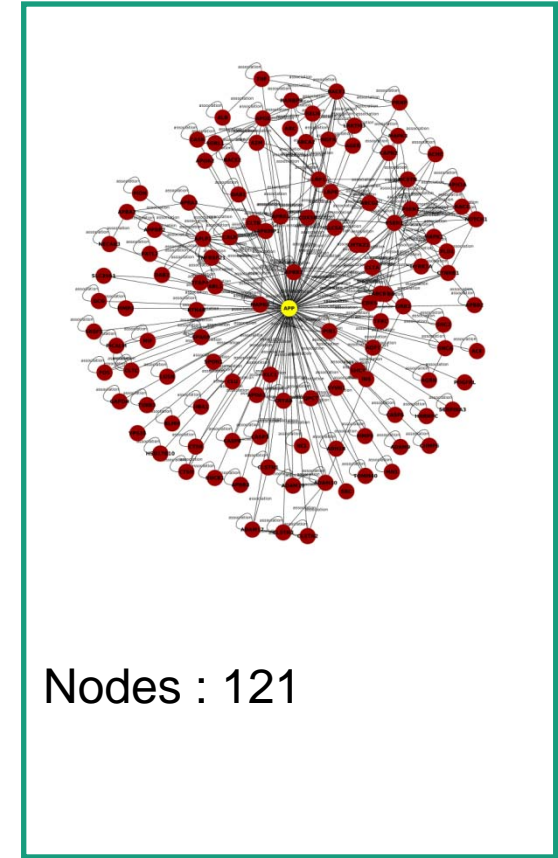
## 500- Abstract



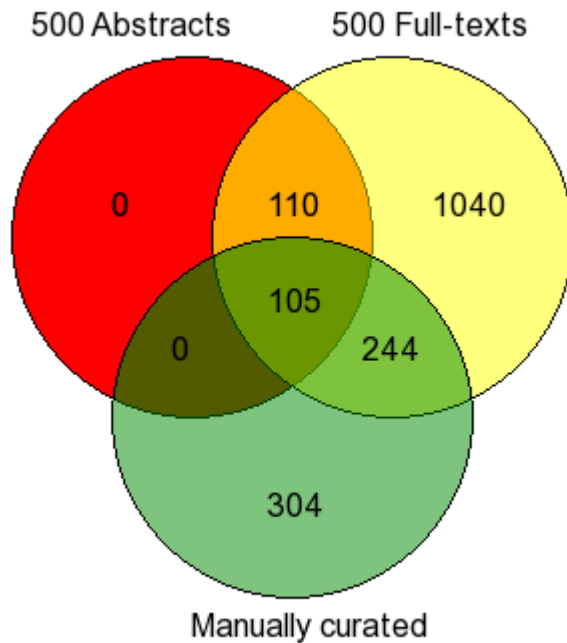
## 500- Full Text



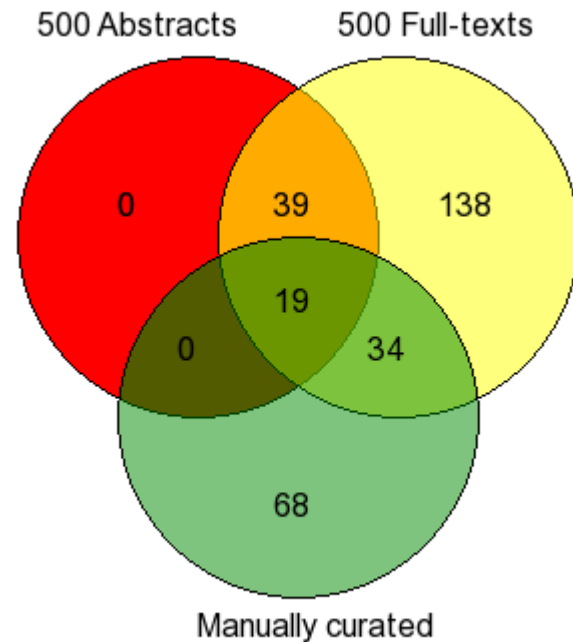
## Manually Curated



# Venn Diagrams



All triples



APP triples

---

# Text Mining & Decision Support

---

- **Our cancer patient may die, because:**

- Nobody can read all the papers that contain relevant information in only two weeks
- The publishing industry does not permit machines to do the job

- **How will we reason over genetic variation information in a functional context?**

- See Naz et al., Briefings in Bioinformatics, in press

- **Time for a “GRAND CHALLENGE”**

- Let us work together to organise a Grand Challenge that demonstrates how our cancer patient could be saved if automated text mining would be supported by the publishing industry

# Take home message .....

---

## Innovative Text Mining and Decision Support

- We can use unstructured text like any database
- We can extract useful and interesting facts, such as triples that represent causal relationships in biomedicine
- We can use these semi-automated information extraction processes to generate a knowledge base in languages such as BEL
- In the future, such knowledge bases will enable decision support in life saving, time-critical scenarios