# Text Mining in der biomedizinischen Forschung

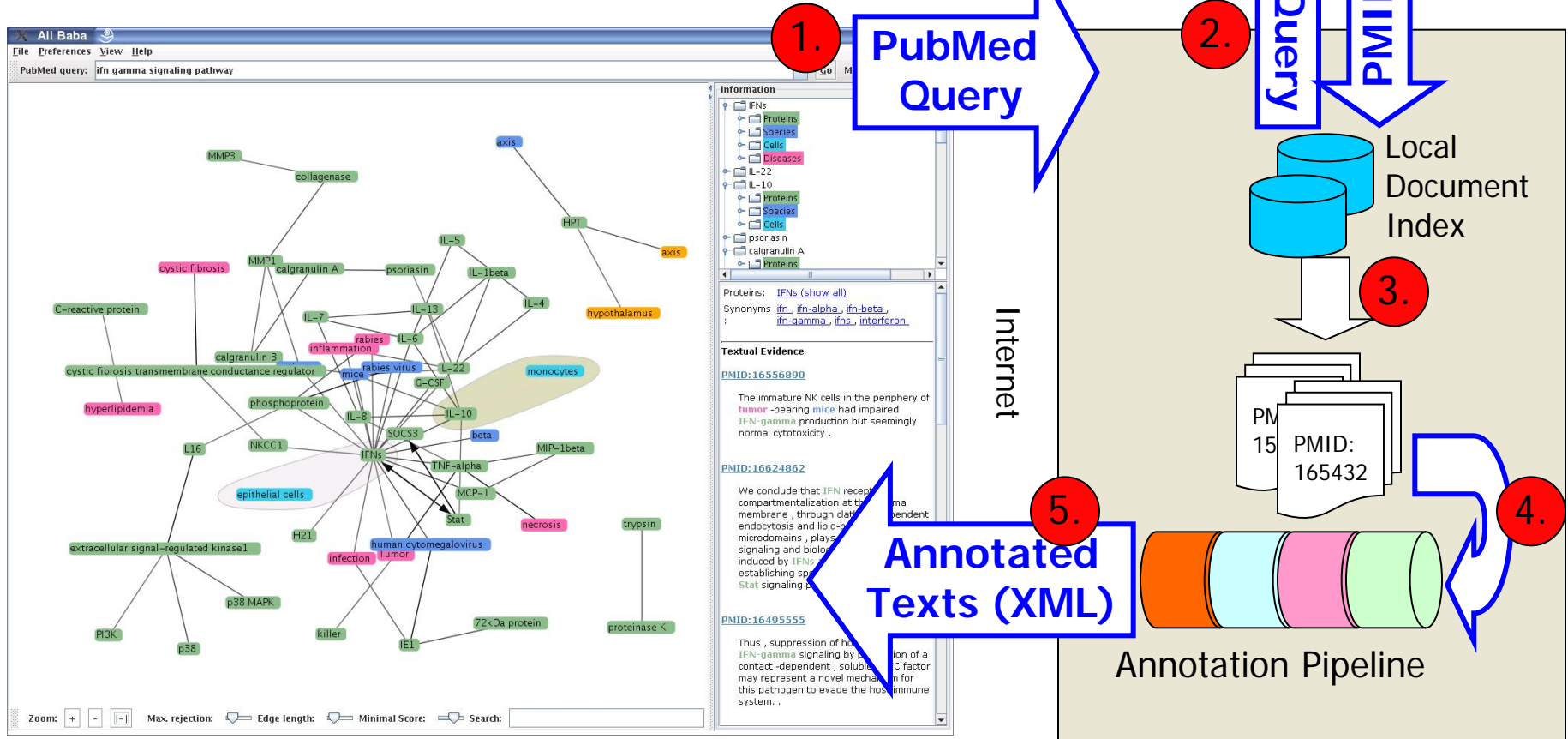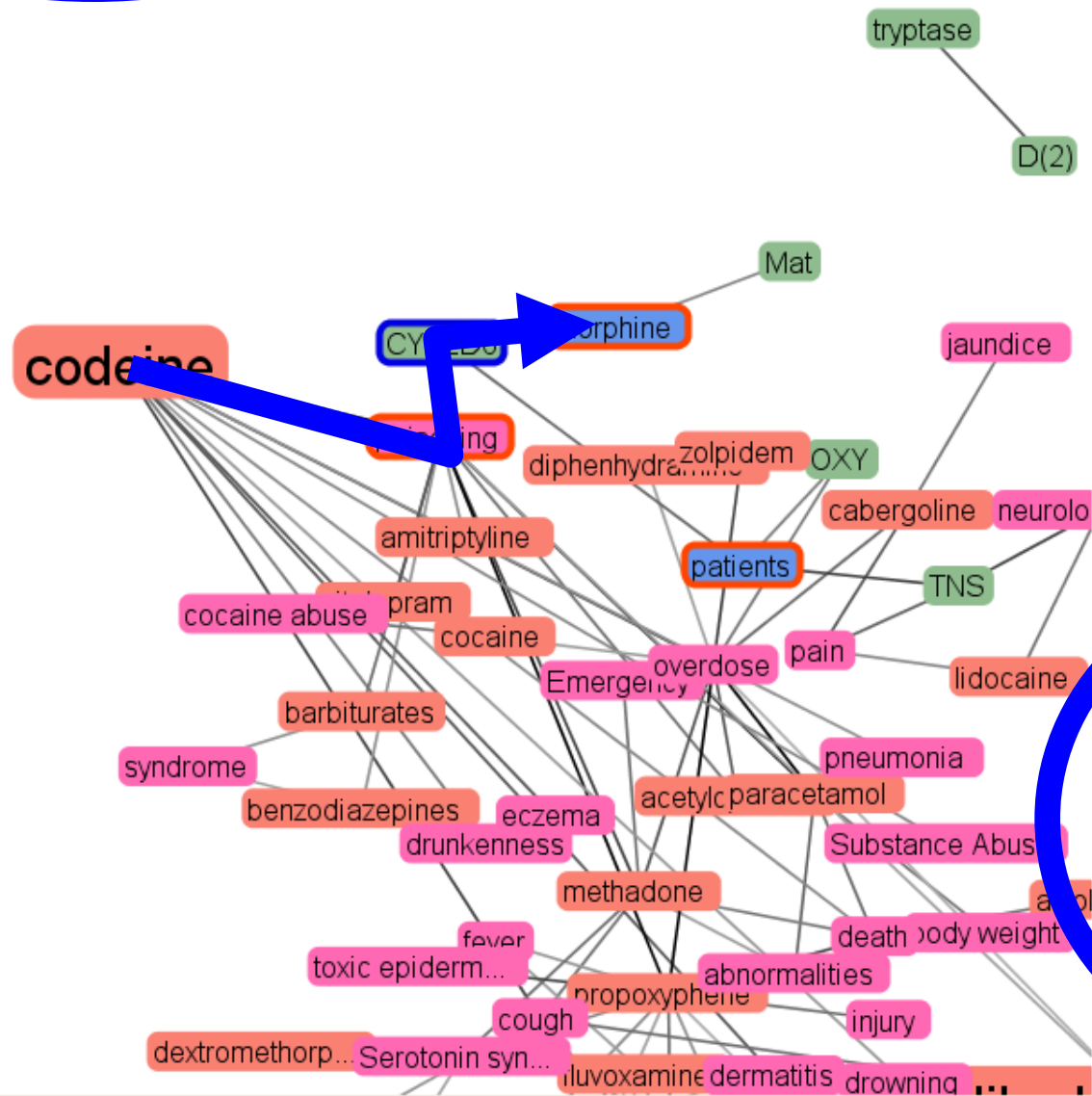Ulf Leser, Humboldt-Universität zu Berlin

# Case Report

- Patient with pneumonia and cough
- Normal dosage of codeine
- Patient not responding any more at day 4
- What's going on?
  - PubMed „Codeine intoxication" -> 70 abstracts
  - Aren't there better ways?

Case report from Univ. Hospital Geneva, thanks to Christian Meisel, Roche

# AliBaba (Plake et al. 2006, Hakenberg et al. 2010)

# What we Need to do

Z-100 is an arabinomannan extracted from Mycobacterium tuberculosis that has various immunomodulatory activities, such as the induction of interleukin 12, interferon gamma (IFN-gamma) and beta-chemokines. The effects of Z-100 on human immunodeficiency virus type 1 (HIV-1) replication in human monocyte-derived macrophages (MDMs) are investigated in this paper. In MDMs, Z-100 markedly suppressed the replication of not only macrophage-tropic (M-tropic) HIV-1 strain (HIV-1JR-CSF), but also HIV-1 pseudotypes that possessed amphotropic Moloney murine leukemia virus or vesicular stomatitis virus G envelopes. Z-100 was found to inhibit HIV-1 expression, even when added 24 h after infection. In addition, it substantially inhibited the expression of the pNL43lucDeltaenv vector (in which the env gene is defective and the nef gene is replaced with the firefly luciferase gene) when this vector was transfected directly into MDMs. These findings suggest that Z-100 inhibits virus replication, mainly at HIV-1 transcription. However, Z-100 also downregulated expression of the cell surface receptors CD4 and CCR5 in MDMs, suggesting some inhibitory effect on HIV-1 entry. Further experiments revealed that Z-100 induced IFN-beta production in these cells, resulting in induction of the 16-kDa CCAAT/enhancer binding protein (C/EBP) beta transcription factor that represses HIV-1 long terminal repeat transcription. These effects were alleviated by SB 203580, a specific inhibitor of p38 mitogen-activated protein kinases (MAPK), indicating that the p38 MAPK signalling pathway was involved in Z-100-induced repression of HIV-1 replication in MDMs. These findings suggest that Z-100 might be a useful immunomodulator for control of HIV-1 infection.

# Find Entities

*Z-100* is an *arabinomannan* extracted from Mycobacterium tuberculosis that has various immunomodulatory activities, such as the induction of **interleukin 12**, **interferon gamma** (**IFN-gamma**) and beta-chemokines. The effects of *Z-100* on human immunodeficiency virus type 1 (HIV-1) replication in human **monocyte-derived macrophages** (**MDMs**) are investigated in this paper. In **MDMs**, *Z-100* markedly suppressed the replication of not only macrophage-tropic (M-tropic) HIV-1 strain (HIV-1JR-CSF), but also HIV-1 pseudotypes that possessed amphotropic Moloney murine leukemia virus or vesicular stomatitis virus G envelopes. *Z-100* was found to inhibit HIV-1 expression, even when added 24 h after infection. In addition, it substantially inhibited the expression of the pNL43lucDeltaenv vector (in which the *env* gene is defective and the *nef* gene is replaced with the *firefly luciferase* gene) when this vector was transfected directly into **MDMs**. These findings suggest that *Z-100* inhibits virus replication, mainly at HIV-1 transcription. However, *Z-100* also downregulated expression of the cell surface receptors **CD4** and **CCR5** in **MDMs**, suggesting some inhibitory effect on HIV-1 entry. Further experiments revealed that *Z-100* induced **IFN-beta** production in these cells, resulting in induction of the 16-kDa **CCAAT/enhancer binding protein** (**C/EBP**) **beta transcription factor** that represses HIV-1 long terminal repeat transcription. These effects were alleviated by SB 203580, a specific inhibitor of **p38 mitogen-activated protein kinases** (**MAPK**), indicating that the **p38 MAPK** signalling pathway was involved in *Z-100*-induced repression of HIV-1 replication in **MDMs**. These findings suggest that *Z-100* might be a useful immunomodulator for control of HIV-1 infection.
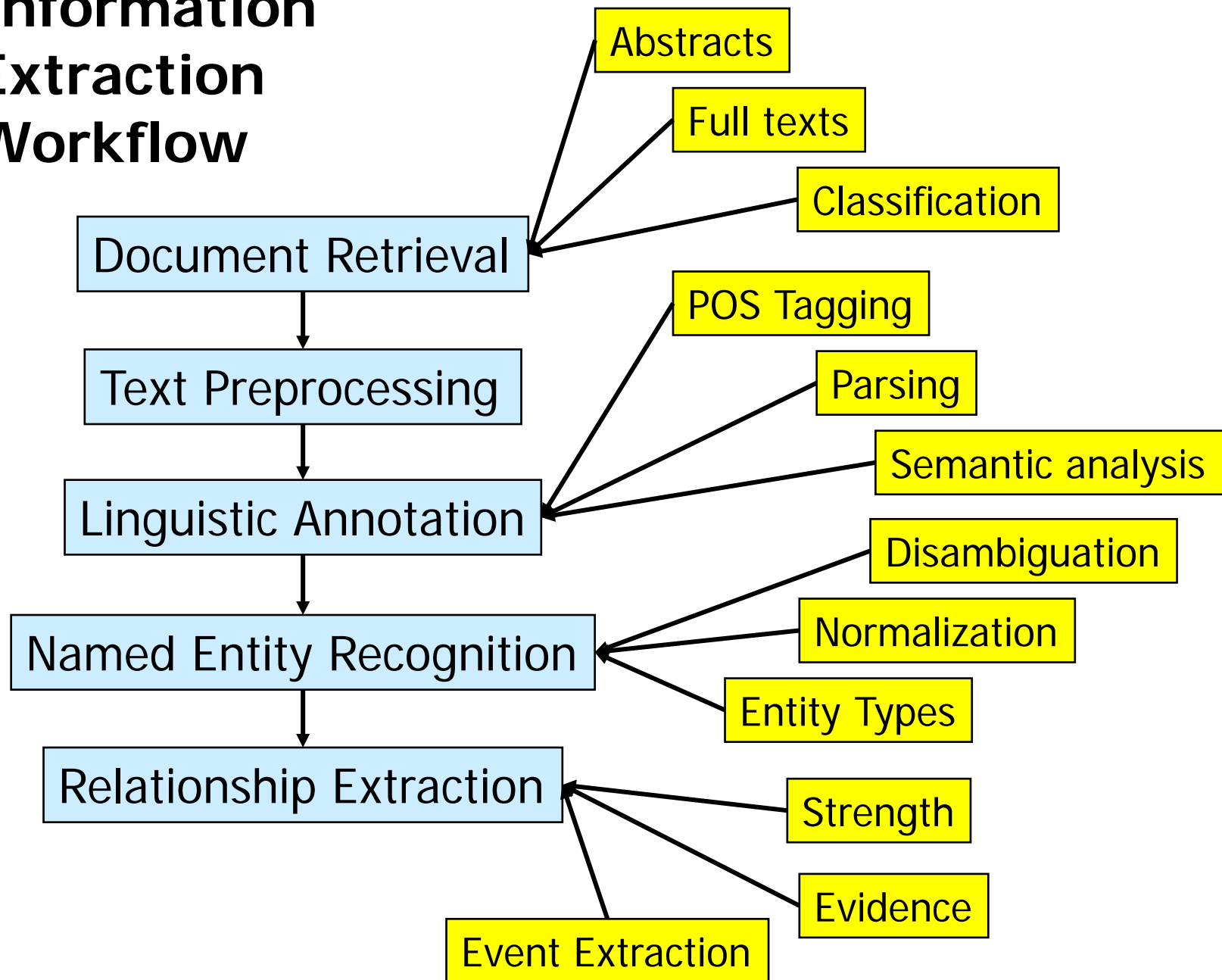
# Normalize Entities

Z-100 is an *arabinomannan* extracted from Mycobacterium tuberculosis that has various immunomodulatory activities, [Tax: 9606] [induc]tion of **interleukin 12**, **interferon gamma** (**IFN-gamma**) and beta-chemokine[s]. The effects of *Z-100* on [human immunodeficiency] virus type 1 (HIV-1) replication in human **monocyte-derived** [macrophages (MDMs)] are investigated in this paper. In **MDMs**, *Z-100* markedly suppressed the [re]plication of not only macrophage-tropic (M-tropic) HIV-1 strain (HIV-1JR-CSF), but also HIV-1 pseudotypes that possessed amphotropic Moloney murine leukemia virus or vesicular stomatitis virus G envelopes. *Z-100* was found to inhibit HIV-1 expression, even when added 24 h after infection. In addition, it substantially inhibited the expression of the pNL43lucDeltaenv vector (in which the *env* gene is defective [and the env gene is repl]aced with the *firefly luciferase* gene) when this vector was tr[ansfected into the M]DMs. These findings suggest that *Z-100* inhibits virus replication, mainly at HIV-1 transcription. However, *Z-100* also downregulated expression of the cell surface receptors **CD4** and **CCR5** in **MDMs**, suggesting some inhibitory effect on HIV-1 entry. Further experiments revealed that *Z-100* induced **IFN-beta** production in these cells, resulting in induction of the 16-kDa **CCAAT/enhancer binding protein** (**C/EBP**) **beta transcription factor** that represses HIV-1 long terminal repeat transcription. These effects were alleviated by SB 203580, a specific inhibitor of **p38 mitogen-activated protein kinases** (**MAPK**), indicating that the **p38 MAPK** signalling pathway was involved in *Z-100*-induced repression of HIV-1 replication in **MDMs**. These findings suggest that *Z-100* might be a useful immunomodulator for control of HIV-1 infection.

Tax: 1773

Entrez: 3458

UMLS: C0001175

GO:0009986

# Find Relationships

Z-100 is an arabinomannan [induction] from Mycobacterium tuberculosis that has various immunomodulatory activities [the induction of] interleukin 12, interferon gamma (IFN-gamma) and beta-chemokines. The effects of Z-100 on human immunodeficiency virus type 1 (HIV-1) replication in human monocyte-derived macrophages (MDMs) are investigated in this paper. In MDMs, Z-100 markedly suppressed the replication of not only macrophage-tropic (M-tropic) HIV-1 strain (HIV-1JR-CSF), but also HIV-1 pseudotypes that possessed amphotropic Moloney murine leukemia virus or vesicular stomatitis virus G envelopes. Z-100 was [inhibit] inhibit HIV-1 expression, even when added 24 h after infection. In addition, specifically inhibited the expression of the pNL43lucDeltaenv vector (in which the env gene is defective and the nef gene is replaced with the firefly luciferase gene) when this vector was transfected directly into MDMs. These findings suggest that Z-100 inhibits virus replication, mainly at HIV-1 transcription. However, Z-100 also downregulated expression of the cell surface receptors CD4 and CCR5 in MDMs, suggesting some inhibitory effect on HIV [induces] eriments revealed that Z-100 induced IFN-beta production in these ce tion of the 16-kDa CCAAT/enhancer binding protein (C/EBP) beta transcription factor that represses HIV-1 long terminal repeat transcription. These effects were alleviated by SB 203580, a specific inhibitor of p38 mitogen-activated protein kinases (MAPK), indicating that the p38 MAPK signalling pathway was involved in Z-100-induced repression of HIV-1 replication in MDMs. These findings suggest that Z-100 might be a useful immunomodulator for control of HIV-1 infection.

# Information Extraction Workflow

# Why Text-Mine?

- Curation
  Support construction of high quality knowledge bases

- Search
  Let users find specific information faster

- Biomedical Research
  Provide background information for specific types of biomedical analysis (network / systems biology)

# Some WBI Projects

- EUMed: Genotype-phenotype relations
  - Mutations, genes, diseases, drugs, mutation-disease, DDI, ...
- CellFinder: Genes characteristic for a given cell
  - Genes, cells, cell lines, transcription, location, function
- OncoPath: Regulatory relationships between TF
  - Genes, species, methods, regulation
- Virtual Liver: Metabolic reactions
  - Chemicals, quantities, units, reactions, methods
- ColoNet: Mammalian clock genes and the clock network

# Topics Today

- Named Entity Recognition
- Applications in Curation
- Application in Search

# Detecting Gene Names

*The human T cell leukemia lymphotropic virus type 1 Tax protein represses MyoD-dependent transcription by inhibiting MyoD-binding to the KIX domain of p300.*

*The human T cell leukemia lymphotropic virus type 1 Tax protein represses MyoD-dependent transcription by inhibiting MyoD-binding to the KIX domain of p300.*

# State-of-the-Art Solutions

- Example: GNAT
  - Large dictionaries, 2nd order CRF
  - Species disambiguation using Linnaeus (NER again)
  - Background texts describing genes and their function
  - Single sense per discourse: Further appearances in the same text
  - No resolution of type and splice variants
- Evaluated performance: ~85% F1-measure
  - Only measured on abstracts, full texts are different
  - Correctness is in the eye of the beholder
  - Different results on different corpora
  - Different quality per species

# What Bio-TM People Study (Neves 2014)



WBI Corpus Repository:
http://corpora.informatik.h

# Experiences: NER+NEN

- No true applications without normalization
  - ChemSpot2 includes OPSIIN
- Every entity type requires a mid-size project
  - Training data, annotator recruitment, domain-specific features, etc.
- In general: The more concrete, the easier
  - Species names versus gene function
- Evaluation usually is a nightmare
  - "The human FGD-3' protein is soluble in water at ..."
  - w/o "protein", w/o splice variants, w/o taxa, w/o water, ...

# Topics Today

- Named Entity Recognition
- Applications in Curation
- Application in Search

# Gene-Regulatory Relationships

- Many published "PPI" actually are derived from co-expression, in-silico predictions, phylogenetic inference, …
- Essential information: Strength of experimental evidence?
- Example: Gene regulation and transcription factors



Source: http://apt.bea.ki.se

# Showing that a TF X directly regulates gene Y

TF X must bind to the promoter region of gene Y



Presence of TF X must change expression of gene Y

Destroying the promoter where X binds must change expression of gene Y

# Towards a Human TF Network (Thomas et al. 2015)

- Human:  ~1200 TFs
- Only ~500 TF-TF relationships are publically available
  - Transfac, ORegAnno, TRRD
- TM-REG: Enhance human TF core network
  - Identification of potential TF-TF relationships using text mining
  - Manual verification by biological experts
  - Distinguish hypothetical from proven relationships

# Workflow

# Top-Scores

- Ranking is crucial to use human time effectively



**Fig. 4.** Precision of our workflow for the $n$ most confidently classified and manually curated sentences. Pairs already contained in a regulatory database are ignored (see Table 3).

# Results
## (Top 2000 sentences, known relationships filtered)

| Experimental Evidence | Databases | Curation | All* |
|---|---|---|---|
| Solely one evidence | | | |
| -E1 : proof of binding on a DNA-Region | 352 | 39 | 381 |
| -E2: change in expression upon activation of the TF | 0 | 45 | 45 |
| -E3: impact of binding site on expression | 6 | 16 | 22 |
| Exactly two evidences | | | |
| -E1 and E2 | 1 | 22 | 24 |
| -E2 and E3 | 1 | 17 | 18 |
| -E1 and E3 | 108 | 43 | 151 |
| All three evidences | | | |
| -E1, E2, E3 | 37 | 128 | 170 |
| Total | 505 | 310 | 815 |

**Table 4.** Identified regulatory relationships from the review in comparison to relationships found in common databases and the according experimental evidence.

# Specific Pathway: Regulation in Liver



- Black: regulations contained in existing RegDBs;
- Red: added by curation of the top 2,500 unspecific sentences;
- Orange: regulations found by manual curation of all 1,435 sentences with co-occurring liver specific TFs.

# Biological Usefulness



- With TM_REG
- Current network
- Randomized

# Experiences: Curation

- Key: Combine text mining and human wisdom
  - No way that current text mining techniques can check evidences at such detail automatically
- Confidence scores of classifiers do make sense
- High throughput is possible
  - Clear problem setting
  - High performance text mining
  - Appropriate curation tools
- Important difference
  - Fast curation: Get certain data from all of PubMed
  - Full curation: Get everything out of this paper

Manageable

Hard & costly

# Topics Today

- Named Entity Recognition
- Application in Curation
- Application in Search

# GeneView

- PubMed does not allow searching synonyms of entities
  - "Find mentionings of BRCA1 (170 synonyms)"
- PubMed does not resolve homonyms
  - "Find mentionings of human BRCA1"
- PubMed cannot rank by "semantic" content
  - "Find mutations in BRCA1"
- PubMed does not allow searching for relationships
  - "Find proteins interacting with BRCA1"
- Geneview does
  - Free web interface
  - Free database: Annotated PubMed

# GeneView

# Biomedical IE at Scale

| Entity ty | Genes | GNAT | of articles |
|---|---|---|---|
| Cell-type | Chemicals | ChemSpot | 5,622 |
| Chemical | | | 9,851,536 |
| Disease | Species | Linneaus | 74,583 |
| Drugs | | | 6,246,067 |
| Enzyme | SNPs | MutationFinder | 590,301 |
| Genes | | | 2,959,439 |
| Histone-mo | Histone-Mod | SETH | 7,673 |
| SNP | Cell type | Banner | 192,544 |
| Species | | | 9,119,134 |
| Tissue | Disease | (unpublished) | 222 |
| Overall | | | 13,463,850 |
| | … | … | |

# Experiences: Search

- User interface more important than 2% more F-measure
- Errors are inevitable; sometimes accepted, sometimes not
- Keeping it up-to-date (tools & data) is a challenge
- Users use context you cannot capture (author, journal, …)
- Much work, few papers
- Scalability is a real issue

# Acknowledgements

- Philippe Thomas
- Mariana Lara Neves
- Tim Rocktäschel
- Torsten Huber
- Michael Weidlich

- Astrid Rheinländer
- Anja Kunkel
- Martin Beckmann
- Marc Bux
- Jörgen Brandt

- Christine Sers
- Nils Blüthgen
- Bertram KLinger
- Andreas Kurz
- Angela Relogio