# Maschinelle Sprachverarbeitung
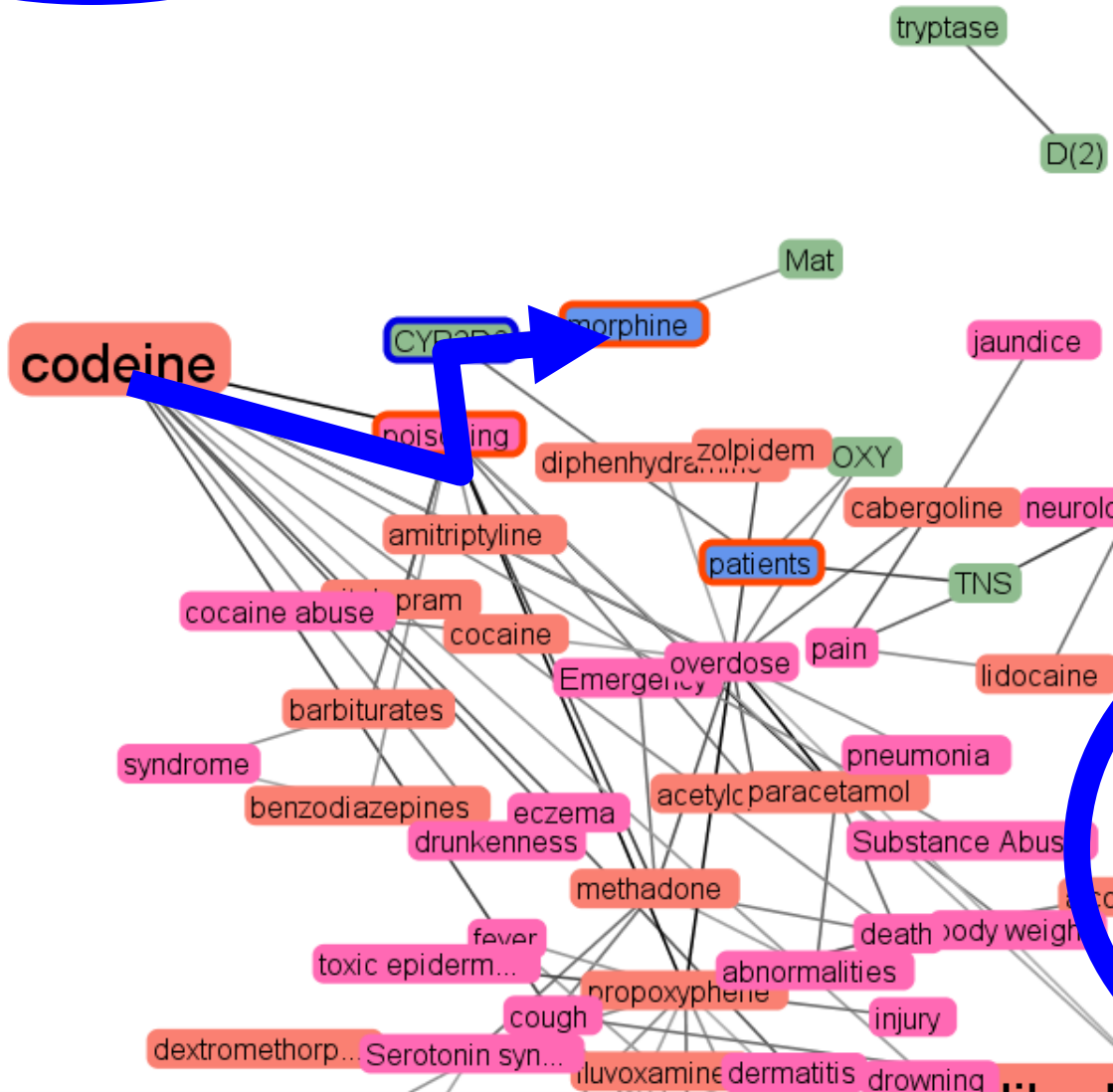
Ulf Leser

# 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" -> 170 abstracts
  - Aren't there better ways?

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

# 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.

# Find Relationships

**Z-100** is an *arabinomannan* [induction] from Mycobacterium tuberculosis that has various immunomodulatory activi... 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 ampho... **M**... ey murine leukemia virus or vesicular stomatitis virus G envelopes. **Z-100** was [inhibit] ...hib... HIV-1 expression, even when added 24 h after infection. In addition, ...lly 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.

# 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.*
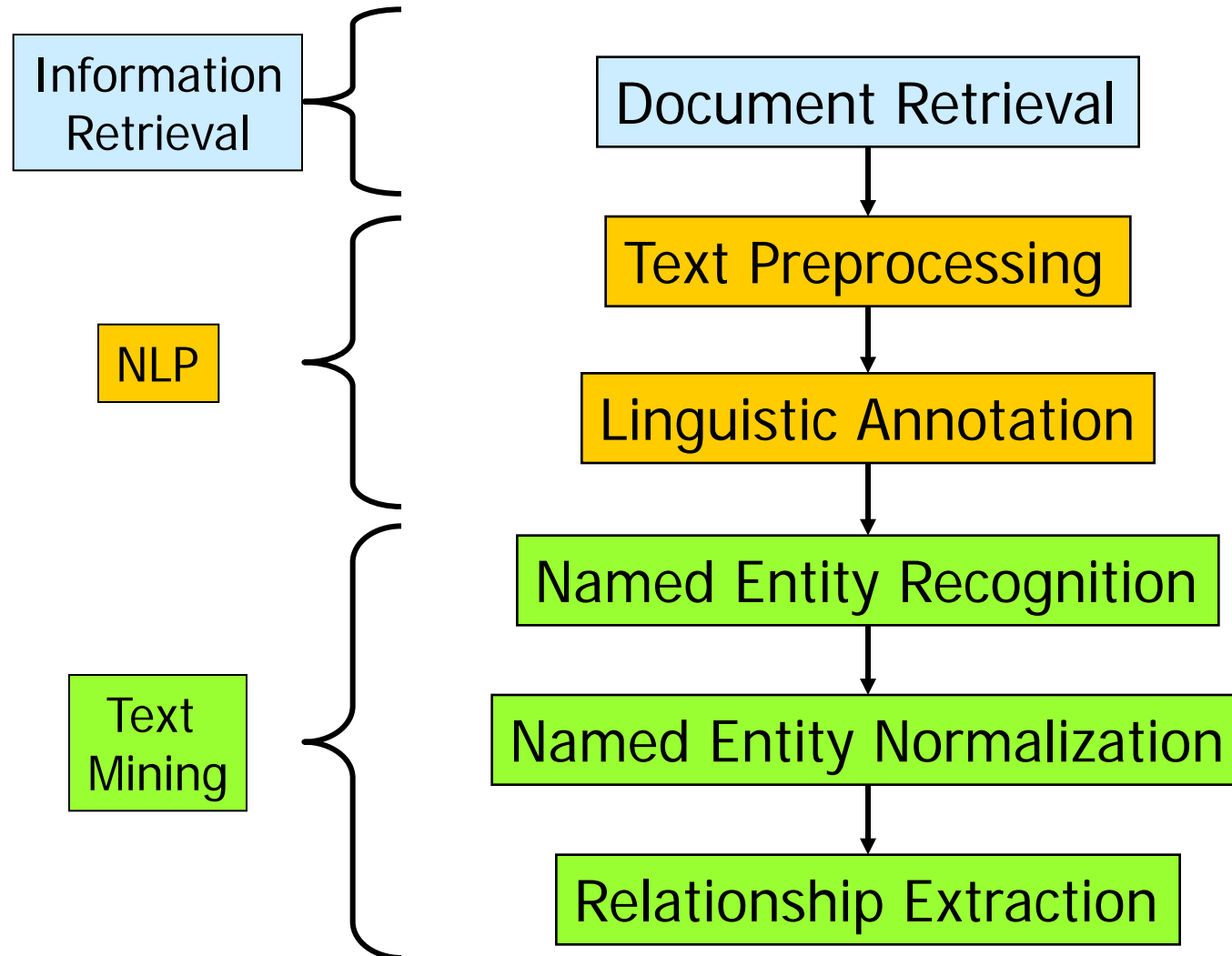
# Detecting Gene Names (Leser & Hakenberg, 2005)

*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.*

- Also: hedgehog, soul, the, white, …
- State-of-the-art methods reach ~85% in NEN
    - Plus 10% for less stringent boundary definitions
    - Large dicts, CRF, species classification, large background corpus, …
    - That's about as high as inter-annotator agreement
- Different performance for other classes (mutations, diseases, functional terms, cell lines, …)

# Typical IE-Workflow

# Applications in Business Intelligence

- **What problems** are most frequently reported by our customers? Which products, product lines, parts etc.?
  - Mails, knowledge bases, repair reports, call centers, …
- How does our customer satisfaction change?
  - Tone in communication?
  - Reports in Blogs, Twitter, …?
- Can we improve customer self service?
  - "Entity Search"
  - Precise routing and prioritization of requests

- See, e.g., Lang, A. and Reinwald, B. (2008). "Nutzung unstrukturierter Daten für Business Intelligence." *Datenbank Spektrum* **25**.

# Some Recent Students Work

- Can we predict the results of elections using Twitter?
    - Tweet classification, sentiment detection
- What aspects of mobile apps are good / bad?
    - Aspect extraction, topic modelling,  sentiment detection
- Can we find texts talking about the biology of stem cells?
    - Text classification, q-gram models
- Can we predict the success of a drug based on papers?
    - Named entity recognition, time series analysis, classification
- Can we semantically cluster tables from the web / papers?
    - Word similarity, text clustering
- Can active learning help for finding gene relationships?

# Modul Maschinelle Sprachverarbeitung

- Vorlesung           2 SWS
- Übung               2 SWS
- Slides are English

- Contact
    Ulf Leser
    <span style="color:blue">Raum:          IV.401</span>
    Tel:            (030) 2093 – 3902
    eMail:          leser (..) informatik . hu-berlin . de

# Literatur

- **Highly recommended**
  - Manning, C.D., Schütze, H. (1999). „Foundations of Statistical Natural Language Processing", MIT Press.
- **Other**
  - Weiss, Indurkhya, Zhang: „Fundamentals of Predictive Text Mining", Springer, 2016
  - Aggarwal, Zhai: „ Mining Text Data", Springer, 2012
  - Manning, C. D., Raghavan, P. and Schütze, H. (2008). "Introduction to Information Retrieval", Cambridge University Press.
  - Lemnitzer, L. and Zinsmeister, H. (2010). "Korpuslinguistik - Eine Einführung", narr Studienbücher.
  - Original papers

# Web

# Questions?

# Questions

- Diplominformatiker?
- IBI / Wirtschaftsinformatik?
- Bachelor?
- Semester?

- Special expectations, experiences, questions?

# Feedback MaschSprach 2015/2016

| Konzeption | Viel neues | Lernziele | Materialien | Rhetorik | Beispiele | Klare Struktur | Verbindungen | Klar verständlich | Übung hilft | Verständliche Antworten | Kritische auseinandersetz | Gute Athmosphäre | Tempo | Schwierigkeit | Arbeitsaufwand | Dozent | Vorlesung | Abweichung vom Optimum | Abweichung pro Frage | Gefehlt | Warum kommen? | Studiengang | Fachsemester |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 6 | 5 | 6 | 5 | 5 | 6 | 6 | 5 | 6 | 6 | 4 | 5 | 5 | 3 | 3 | 3 | 2 | 2 | 10 | 0,56 | 1 | 4 | M | 4 |
| 6 | 5 | 6 | 5 | 5 | 6 | 6 | 6 | 5 | 6 | 5 | 5 | 6 | 2 | 2 | 3 | 1 | 1 | 8 | 0,44 | 0 | 4 | M | 1 |
| 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 3 | 3 | 3 | 1 | 1 | 0 | 0,00 | 0 | 4 | M | 2 |
| 6 | 6 | 6 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 5 | 6 | 6 | 3 | 3 | 3 | 1 | 1 | 2 | 0,11 | 0 | 2,3,4, | M | 3 |
| 6 | 6 |  | 6 | 6 | 6 | 6 | 6 | 6 |  |  | 6 | 6 | 3 | 3 | 3 | 1 | 1 | 0 | 0,00 | 0 | 4 | M | 6 |
| 6 | 6 | 6 | 5 | 5 | 6 | 6 | 6 | 6 | 6 | 5 | 4 | 6 | 3 | 3 | 2 | 1 | 1 | 6 | 0,33 | 1 | 2,3,4 | M | 1 |
| 6 | 6 | 6 | 5 | 6 | 6 | 6 | 6 | 6 | 5 | 4 | 6 | 6 | 3 | 4 | 3 | 1 | 1 | 5 | 0,28 | 2 | 2,4 | M | 3 |
| 6 | 3 | 6 | 4 | 6 | 4 | 6 | 6 | 6 | 5 | 6 | 6 | 6 | 2 | 2 | 3 | 1 | 1 | 10 | 0,56 | 0 | 2,4 | M | 4 |
| 6 | 6 | 5 | 6 | 6 | 5 | 6 | 6 | 6 | 6 | 5 | 6 | 6 | 3 | 3 | 3 | 1 | 1 | 3 | 0,17 | 3 | 4 | M | 5 |
| 6 | 5 | 6 | 4,0 | 5 | 6 | 6 | 6 | 6 | 5 | 3 | 2 | 6 | 2 | 2 | 3 | 2 | 2 | 16 | 0,89 | 2 | 4 | BA | 6 |
| 6 | 5 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 5 | 6 | 5 | 6 | 3 | 3 | 3 | 1 | 1 | 4 | 0,22 | 0 | 3,4 | M | 1 |
| 6 | 5 | 5 | 6 | 6 | 6 | 4 | 5 | 6 | 4 | 5 | 6 | 6 | 3 | 3 | 3 | 2 | 2 | 10 | 0,56 | 0 | 4 | BF | 3 |
| 6 | 6 | 6 | 5 | 6 | 6 | 5 | 6 | 6 | 6 | 5 | 6 | 6 | 3 | 3 | 3 | 1 | 1 | 3 | 0,17 | 0 | 4 | M | 1 |
| 4 | 5 | 4 | 4 | 6 | 5 | 4 | 5 | 6 | 6 | 4 | 6 | 4 | 3 | 3 | 3 | 2 | 2 | 17 | 0,94 | 0 | 3,4 | M | 3 |
| 6 | 6 | 4 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 5 | 6 | 3 | 3 | 3 | 1 | 1 | 4 | 0,22 | 2 | 3 | M | 9 |
| 6 | 6 |  | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 5 | 6 | 6 | 3 | 3 | 3 | 1 | 1 | 1 | 0,06 | 0 | 3,4 | M | 9 |
| 6 | 5 | 6 | 6 | 6 | 5 | 6 | 6 | 6 | 6 | 4 | 5 | 6 | 3 | 3 | 3 | 1 | 1 | 5 | 0,28 | 1 | 3,4 | M | 3 |
| 6 | 6 | 5 | 6 | 6 | 6 | 6 | 6 | 6 | 6 | 3 | 5 | 5 | 3 | 3 | 3 | 1 | 1 | 6 | 0,35 | 1 | 4 | M | 5 |
| 6 | 5 | 6 | 5 | 5 | 5 | 6 | 6 | 6 | 5 | 4 | 5 | 5 | 3 | 3 | 3 | 1 | 2 | 10 | 0,56 | 3 | 2,4 | M | 2 |
| 5,89 | 5,42 | 5,53 | 5,26 | 5,74 | 5,74 | 5,68 | 5,84 | 5,95 | 5,56 | 4,72 | 5,32 | 5,74 | 2,84 | 2,89 | 2,95 | 1,21 | 1,26 |  |  | 0,8 |  |  | 3,7 |
| 6,00 | 6,00 | 6,00 | 6,00 | 6,00 | 6,00 | 6,00 | 6,00 | 6,00 | 6,00 | 6,00 | 6,00 | 6,00 | 3,00 | 3,00 | 3,00 | 1,00 | 1,00 |  |  |  |  |  |  |
| 0,11 | 0,58 | 0,47 | 0,74 | 0,26 | 0,26 | 0,32 | 0,16 | 0,05 | 0,44 | 1,28 | 0,68 | 0,26 | 0,16 | 0,11 | 0,05 | -0,21 | -0,26 |  |  |  |  |  |  |
| 0,11 | 0,58 | 0,47 | 0,74 | 0,26 | 0,26 | 0,32 | 0,16 | 0,05 | 0,44 | 1,28 | 0,68 | 0,26 | 0,16 | 0,11 | 0,05 | 0,21 | 0,26 | 6,40 |  |  |  |  |  |

# Was ich ändern will

- NER und RE etwas unsystematisch
- Scalable IE rausnehmen und was dazu?
  - Deep learning and word embeddings
  - Question Answering, Opinion Mining, Topic modelling?
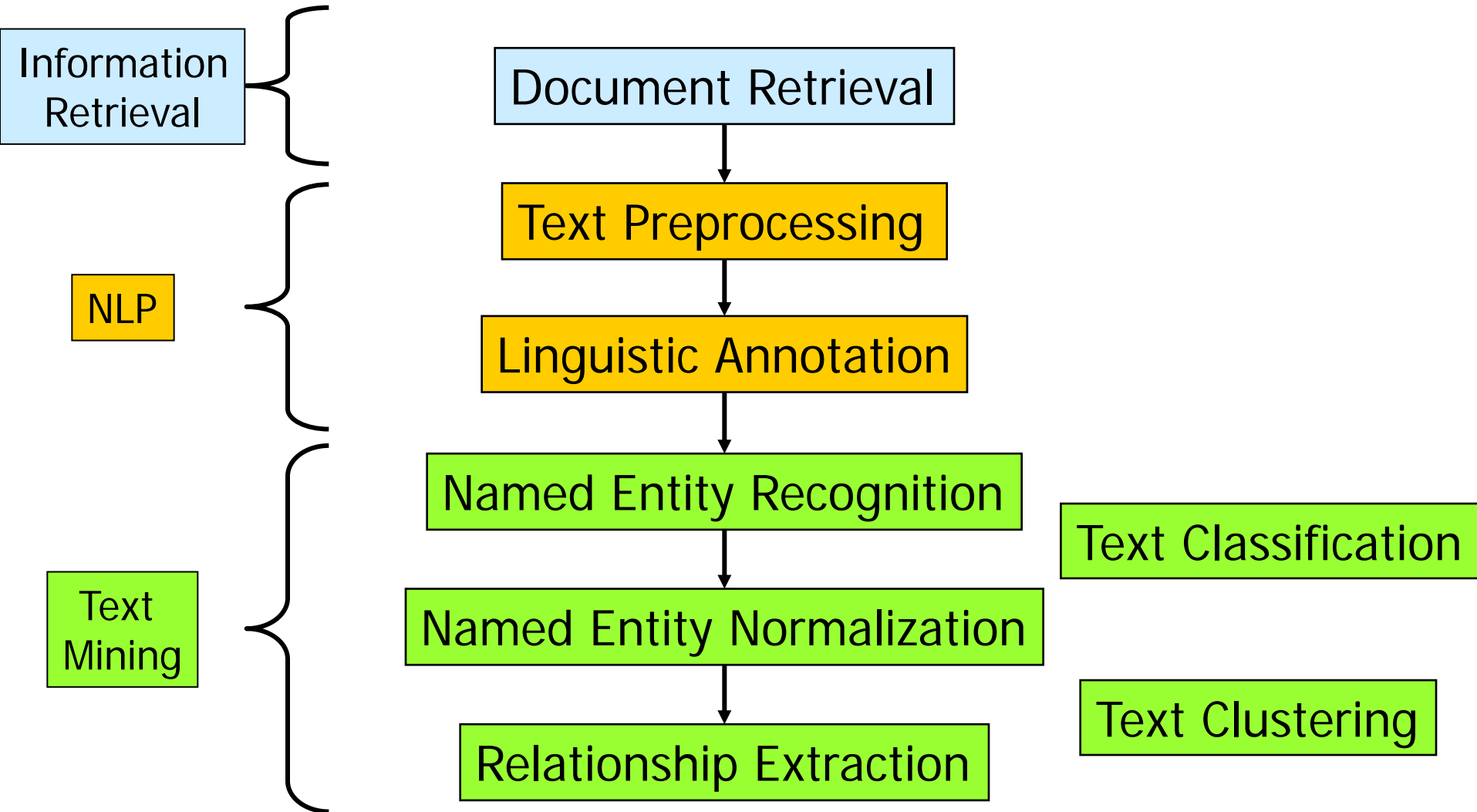- Vorstellende Gruppen vorher festlegen?

# Content of this Lecture

- A very short primer on Information Retrieval
- A very short primer on Natural Language Processing
- A very short primer on Text Mining

# Phases in Text Mining

Information Retrieval

Document Retrieval

↓

Text Preprocessing

↓

NLP

Linguistic Annotation

↓

Named Entity Recognition

↓

Text Mining

Named Entity Normalization

Text Classification

↓

Relationship Extraction

Text Clustering

# Information Retrieval (aka "Search")

- Find all documents which contain the following words
- „Leading the user to those documents that will best enable him/her to satisfy his/her need for information"[Robertson 1981]
  - A user wants to know something
  - The user needs to tell the machine what he wants to know: query
  - Posing exact queries is difficult: room for interpretation
  - Machine interprets query to compute the (hopefully) best answer
  - Goodness of answer depends on original intention of user, not on the query (relevance)

# Difference to Database Queries

- Queries: Formal language versus natural language
- Exactly defined result versus loosely described relevance
- Result set versus ranked list of results
- DB: Posing the right query is completely left to the user
- IR: Understanding the query is a problem of the software

# Natural Language Processing

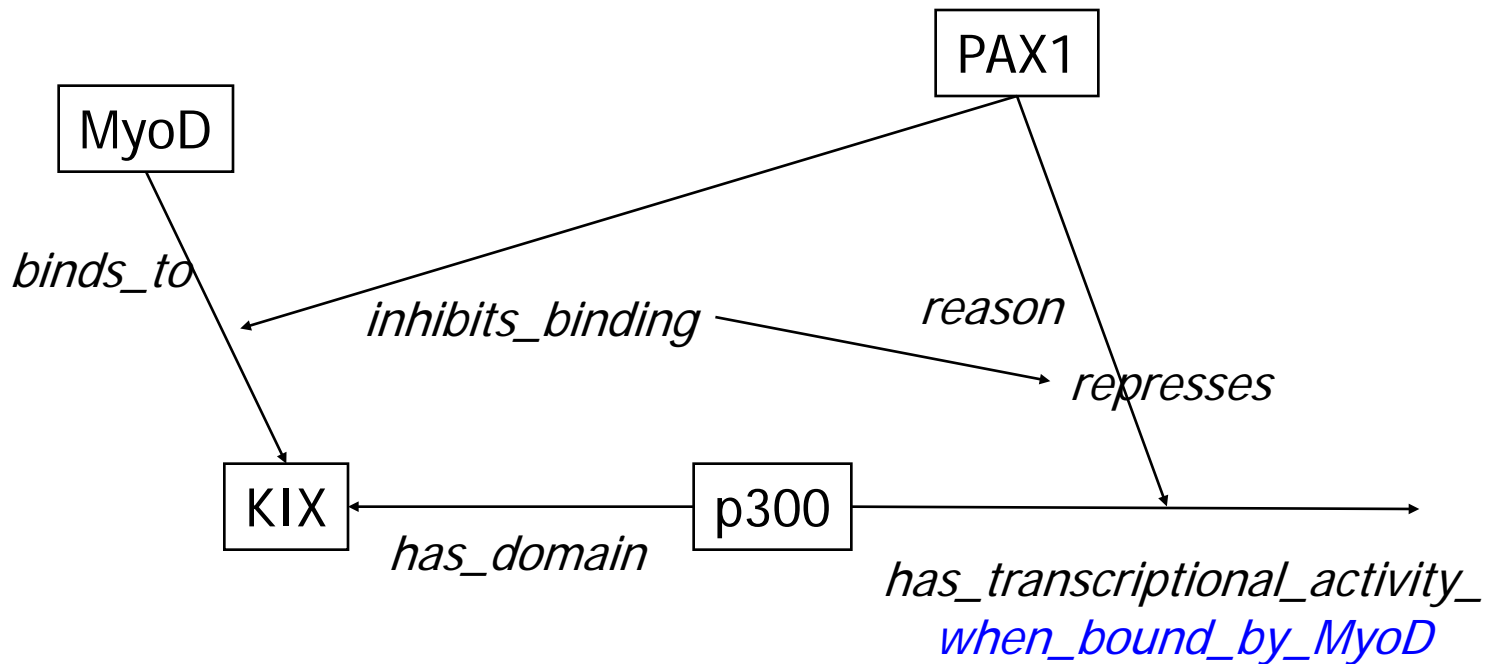- Making natural language text accessible to a computer
  - Find semantic units: words, tokens, phrases, clauses, sentences
    - "Implementing the C4.5 algorithm with languages such as DOT.NET, Java etc. is not as simple as one might think ..."
    - "The $\alpha(3)$-helicase-5' mRNA is ..."
  - Find grammatical role of words
  - Find grammatical structure of sentences
  - Find syntactic relationships between entities
  - Draw semantic inferences from a text
  - ...

- "Understand" the text

# Understanding Text is Difficult (even for us)

*„The PAX1 protein represses MyoD-dependent transcription by inhibiting MyoD-binding to the KIX domain of p300."*

# Part-Of-Speech Tagging

- Part-of-Speech (POS) is the grammatical class of a word
  - Adverb, verb, adjective, ...
  - Verb: Tense, number, ...
  - Noun: Gender, case, number, ...

| | | | |
|---|---|---|---|
| **DT** | Determiner | **SYM** | Symbol (math or scientific) |
| **EX** | Existential *there* | **UH** | Interjection |
| **FW** | Foreign word | **VB** | Verb, base form |
| **NN** | Noun, singular or mass | **VBD** | Verb, past tense |
| **NNS** | Noun, plural | **VBG** | Verb, gerund/present participle |
| **NNP** | Proper noun, singular | | |
| **NNPS** | Proper noun, plural | **VBZ** | Verb, 3rd person/singular, present |
| **RP** | Particle | | |

- POS tagging: Given a text, assign each word its POS
  - *"Does/VBZ flight/NN LH750/NNP serve/VB dinner/NN ?"*
  - Caveat: There are different tag sets
- POS tags are very useful for many tasks
  - NER: names of entities should be nouns
- Method: Maximum Entropy, Hidden Markov Models

# Text Mining

- Text Mining = "Data Mining on text"
- Text Mining = "Statistical NLP minus parsing" [Altmann/Schütze]
- Typical tasks
  - Document classification (route emails to the right operator)
  - Document clustering (group search results by topics)
  - Information extraction (find all celebrities and their partners)

# Clinical Entity Recognition for ICD-9 Code Prediction in Clinical Discharge Summaries

Diploma Thesis Presentation

Jonathan Bräuer

02.10.2017

- Clinical data is often stored in textual form
- Reports contain valuable information
  - Diseases, symptoms, treatments, drugs, dosage, family history, lab measurements, images/radiology, progression, ...
  - Many of these not available in structured form
- Especially important: Disease (symptoms, phenotypes)
  - For accounting
  - For decision support

- **DATE OF ADMISSION:** MM/DD/YYYY

- **DATE OF DISCHARGE:** MM/DD/YYYY

- **DISCHARGE DIAGNOSES:**

- 1. Vasovagal syncope, status post fall.
  2. Traumatic arthritis, right knee.
  3. Hypertension.
  6. History of chronic obstructive pulmonary disease.

- **BRIEF HISTORY:** The patient is an (XX)-year-old female with history of previous stroke; hypertension; COPD, stable; renal carcinoma; presenting after a fall and possible syncope. While walking, she accidentally fell to her knees and did hit her head on the ground, near her left eye. Her fall was not observed, but the patient does not profess any loss of consciousness, recalling the entire event. The patient does have a history of previous falls, one of which resulted in a hip fracture. She has had physical therapy and recovered completely from that...

- **DIAGNOSTIC STUDIES:** All x-rays including left foot, right knee, left shoulder and cervical spine showed no acute fractures. The left shoulder did show old healed left humeral head and neck fracture with baseline anterior dislocation. ...

- **HOSPITAL COURSE:**

- 1. Fall: The patient was admitted and ruled out for syncopal episode. Echocardiogram was normal, and when the patient was able, ...

- 2. Status post fall with trauma: The patient was unable to walk normally secondary to traumatic injury of her knee, causing significant pain and swelling. Although a scan showed no acute fractures, ...
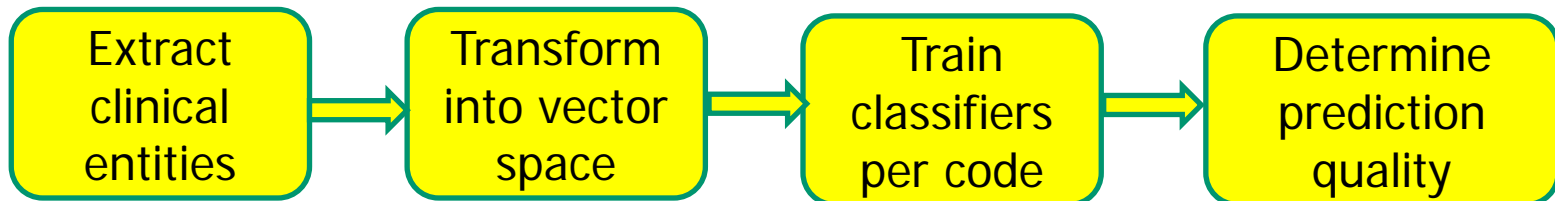
# Goals and Methods

- Predict discharge diagnosis based on clinical texts
- Approach 1: Recognize diseases in text (NER-based approach)

```
Extract          Map to          Compare to
clinical    →    ICD-9-CM   →    assigned
entities                         codes
```
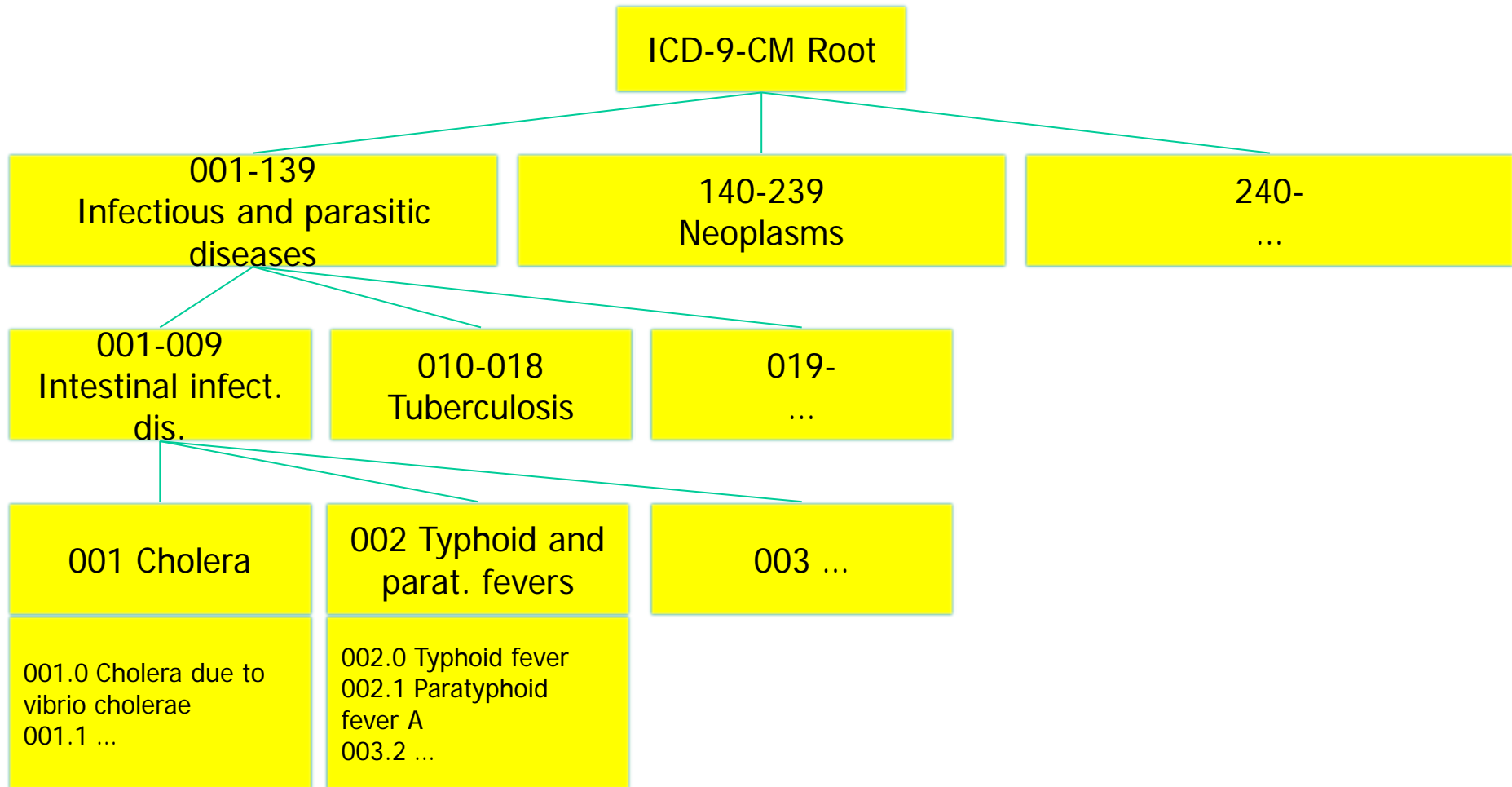
- Approach 2: Predict disease based on (entire, partial) text (classification-based approach)

```
Extract          Transform       Train           Determine
clinical    →    into vector →   classifiers →   prediction
entities         space           per code        quality
```
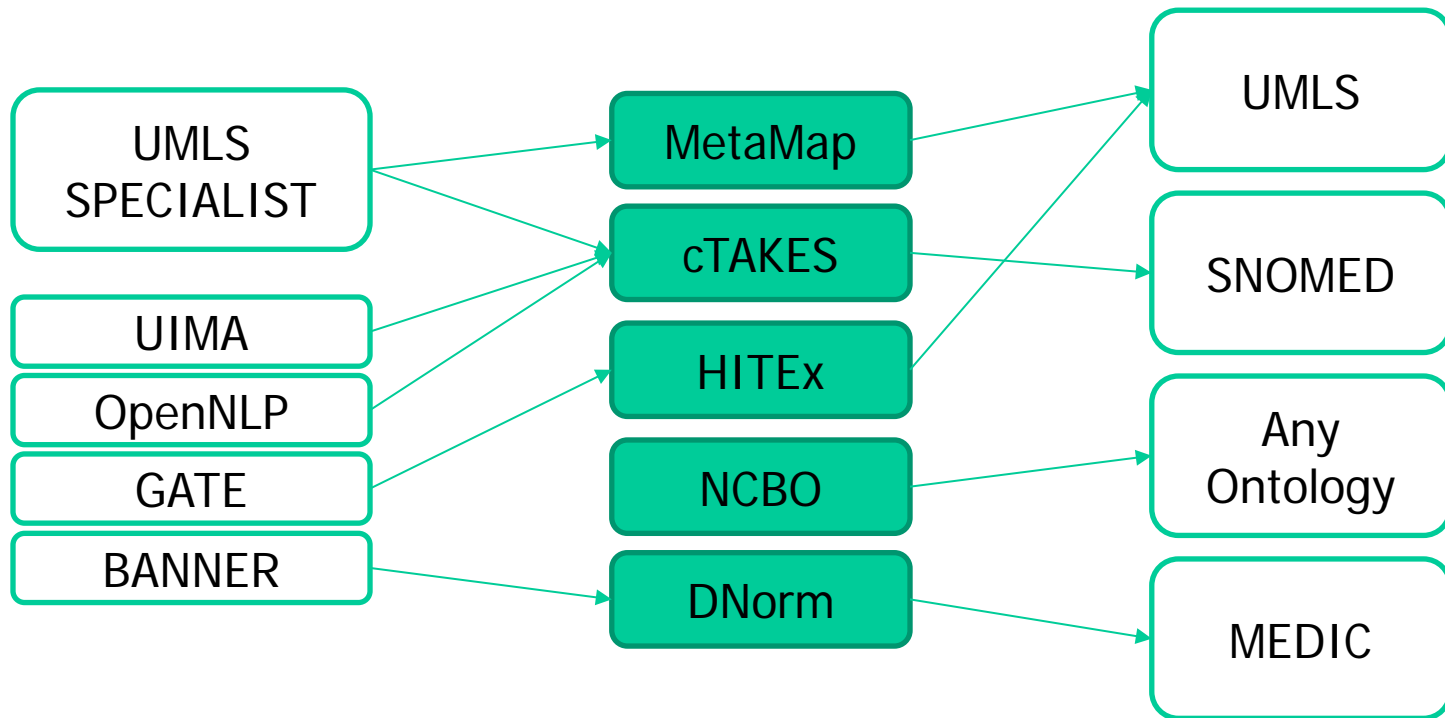
# Disease Names: ICD-9

# Medical NER Tools Evaluated

# Number of Extracted Concepts (Per Document)



Legend: ■ Total ■ Disease Mentions ICD-9 Mapped

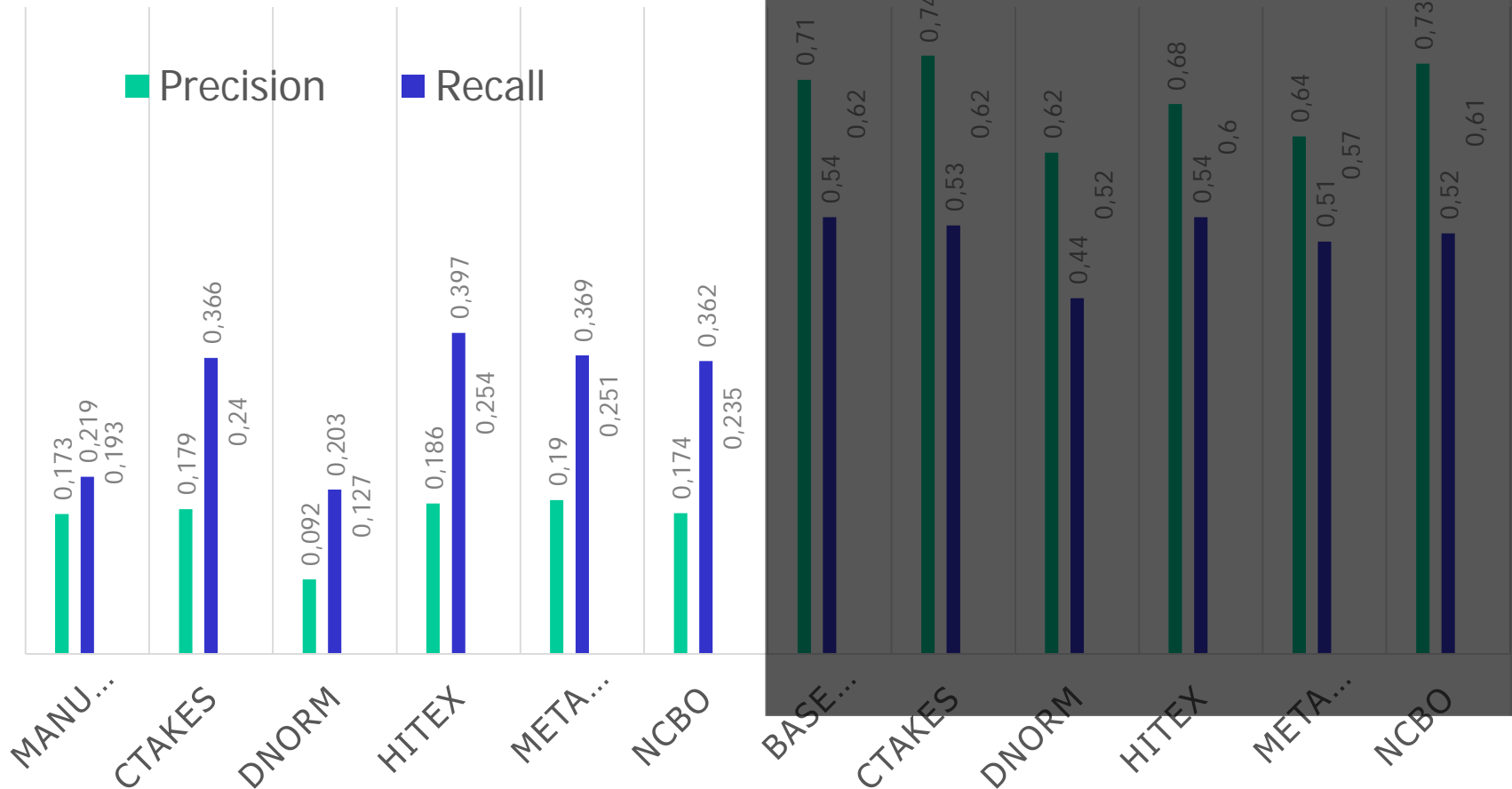| | Total | Disease Mentions | ICD-9 Mapped |
|---|---|---|---|
| CTAKES | 635,80 | 191,40 | 33,30 |
| DNORM | 18,50 | 18,50 | 4,80 |
| HITEX | 217,20 | 61,00 | 35,10 |
| METAMAP | 294,90 | 119,70 | 32,90 |
| NCBO | 642,30 | 152,70 | 36,70 |
| MANUAL | | 59,50 | 10,90 |

# Issues (Typical)

- **Hierarchical classification** – which level of ICD-9?
  - Higher levels: More training data, few classes, high accuracy
    But: Little value
  - Lower levels: Little training data, many classes, low accuracy
    But: High value

- **Mapping** between ontologies
  - Concepts with different syntax & synonyms
  - Concepts at different granularities
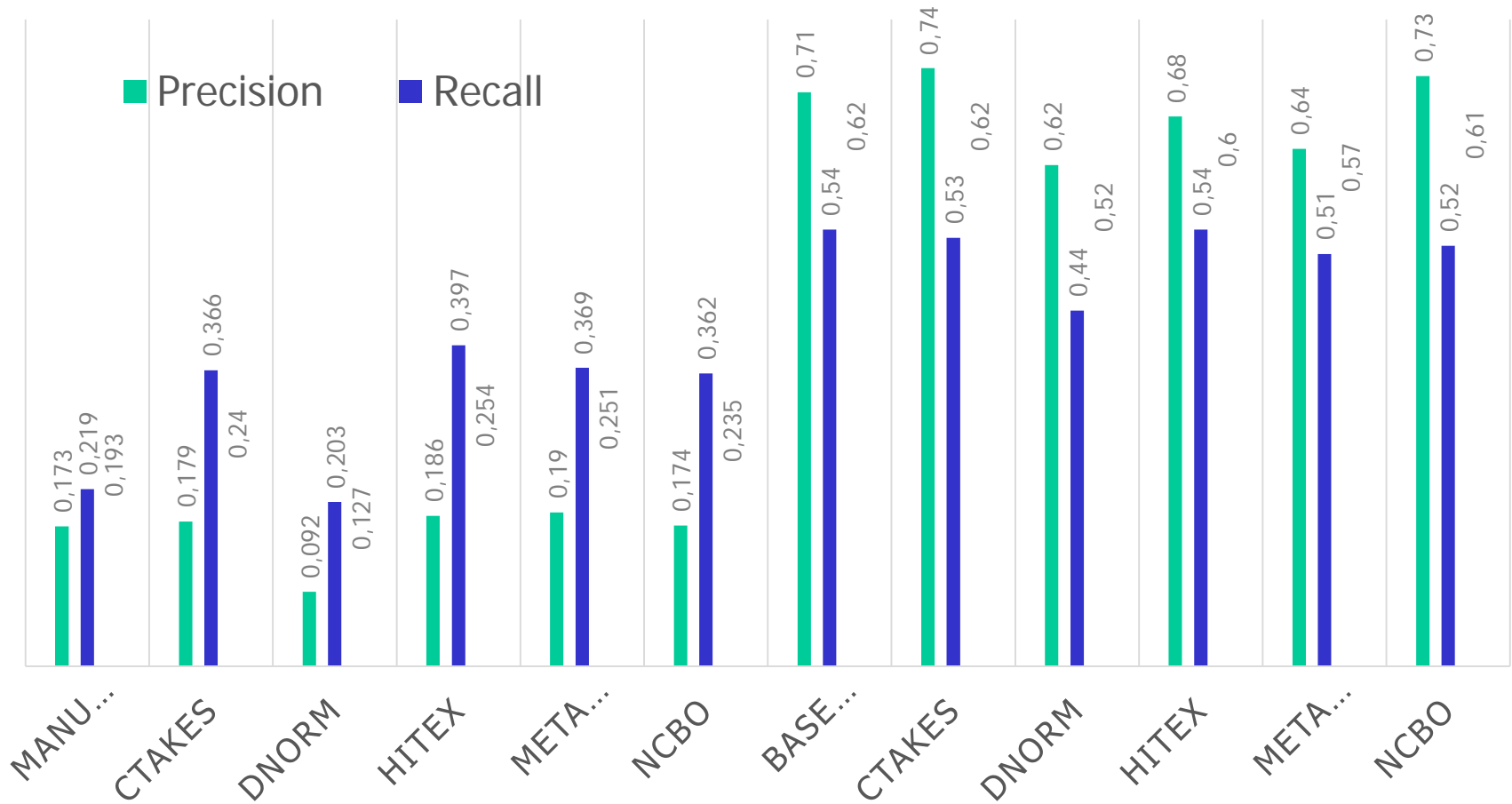  - Conflicting subsumption relationships
  - Diverging coverage
  - …

```
  DNorm        MetaMap        HITEx          cTAKES

  Other    →   UMLS      →   SNOMED-CT   →   ICD-9-
  Ontologies   (CUI)         (Code)          CM
```

# Results / Evaluation

- 50 k discharge summaries
- 7 k classes (diagnosis codes)

# Results / Evaluation

# Error Analysis

- False positive code assignments
  - Mapping errors
  - Contextual errors
  - Negation / temporal status
- False negative code assignments
  - Obvious codes not tagged in gold standard (hypertension)
  - Heavy use of abbreviations and acronyms
  - Missing sections
  - Missing mention

# What we will not cover

- Linguistic analysis beyond POS / parsing
- Spoken language
- Machine translation
- Cross-language search / analysis
- User interfaces
- Special classification problems: Sentiment analysis, question answering, Watson
- Topic modelling
- ...