

# 10 – Extracting structured information from text

## IA161 Advanced Techniques of Natural Language Processing

Zuzana Nevěřilová, Vojtěch Kovář

NLP Centre, FI MU, Brno

November 20, 2017

1 What?

2 Why?

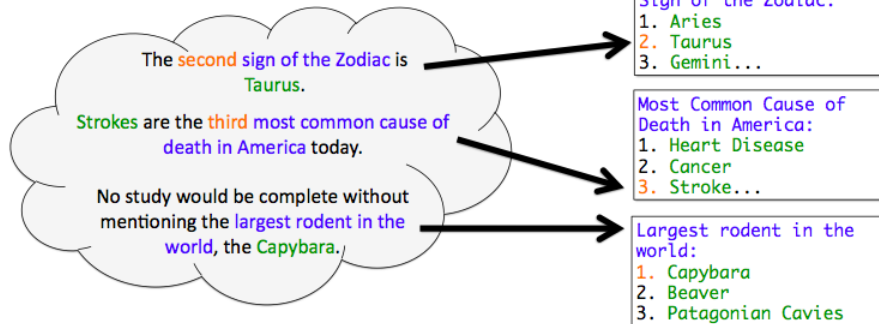
3 How?

4 Who?

## Unstructured Web Text



## Structured Sequences



# Information Extraction Goals

Fed Chairman  
Ben Bernanke  
said the U.S.  
economy...  
The euro rose to  
\$1.2008,  
compared to  
\$1.1942  
on Tuesday.



# Information Extraction Goals

- Types of Factual Information
  - ▶ keywords
  - ▶ entities
  - ▶ relations
  - ▶ events
- Extracted from
  - ▶ Different text types: news articles, emails, novels, output from speech recognizer
  - ▶ Different domains or the general domain

# Information Extraction Applications

- Direct applications for specific users:
  - ▶ financial analysts
  - ▶ media analysts
  - ▶ PR workers
- Use in subsequent computer applications
  - ▶ information systems
  - ▶ question answering
  - ▶ automatic reasoning
  - ▶ automatic summarization
  - ▶ ...
- Disambiguate and shorten the information
- Find informational redundancy, aggregate information from several sources









# Successful Information Extraction Systems

Google

museums in prague

Web Maps Images Videos Shopping More Search tools

Museums frequently mentioned on the web

							
National Museum, Prague	National Technical Museum	Prague Jewish Museum	Museum of Communism, Czech Republic	Prague National Gallery	Antonin Dvořák Museum	Museum Kampa	Museum Decorative Arts in Prague

## Prague Museums - Visitor Information - My Czech Republic

[www.myczechrepublic.com](http://www.myczechrepublic.com) > [Prague Guide](#) > [Museums & Galleries](#) ▼

Museums in Prague. National Museum. National Technical Museum and other

# Successful Information Extraction Systems

- x.ai – automatic personal assistant Amy
  - ▶ agrees automatically on meeting times
  - ▶ recognizes/asks for contact details
  - ▶ operates over Google calendar
- Extracting protein interaction from research texts
- Summarizing and filtering stock market news
- Extracting information about conflicts from news
- Smaller systems for more specialized tasks

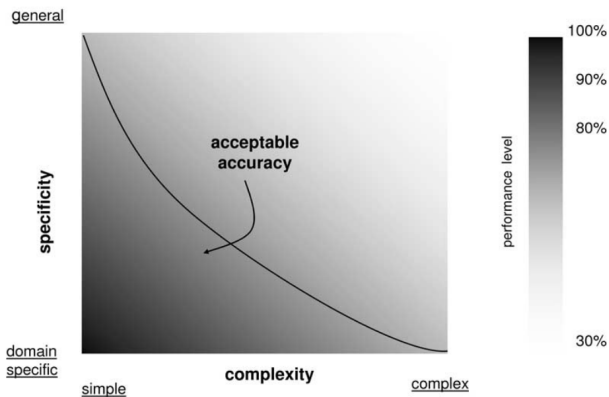


# Information Extraction Evaluation

- Message Understanding Conference + Text REtrieval Conference
- series of conferences starting in 80s and 90s
- shared tasks + competition among systems
- helped standardization in the field

# Information Extraction Approaches

- Specific domain / Complex information
  - ▶ precise, narrow requests from small homogeneous corpora
  - ▶ weighting/ordering/refining results
- General domain / Simple snippets of information
  - ▶ vague request from huge data
  - ▶ aggregation of the response




# Information Extraction Components

named entity recognition (NE)	finds and classifies names, places, dates, keywords etc.	rocket, Tuesday, Dr Head, Dr Big Head, We Build Rockets Inc.
coreference resolution (CO)	finds identity relations between entities	It = rocket, Dr Head = Dr Big Head
relation extraction (RE)	add description to entities, finds relation between entities (based on CO)	rocket = red shiny, rocket – brainchild – Dr Head, Dr Head – works for – We Build Rockets Inc.
event extraction (EE)	fits RE into event scenarios	rocket launching event

The *shiny red rocket* was fired on *Tuesday*. It is the *brainchild* of *Dr Big Head*. *Dr Head* is a staff scientist at *We Build Rockets Inc.*

# Information Extraction Components

named entity recognition (NE)	discussed in detail in lecture 04	Z. Nevěřilová, 20/11/2017, A219, IA161
coreference resolution (CO)	discussed in detail in lecture 08	jej = IA161
relation extraction (RE)	discussed later in this lecture	A219 = computer room, IA161 – being taught – 20/11/2017
event extraction (EE)	event recognition, “filling the gaps”	course: name [IA161], date [20/11/2017], lecture room [A219], teacher [Z. Nevěřilová]

  
domain dependent  
tied to scenarios of interest

Výuka předmětu IA161 se koná v pondělí v počítačové učebně A219.  
20. listopadu 2017 jej učí Zuzana Nevěřilová.

# Relation Extraction

- noun phrase recognition
- verb group recognition
- adjective phrase recognition
- adverbial phrase recognition
- partial parsing
- event recognition: actors = noun phrases, action = verb, place = adverbial phrase, time = adverbial phrase
- rule-based or statistical

within a given task, the set of relations is *fixed*

Best MUC results:  $\approx 75\text{--}80\%$  (humans  $\approx 90\%$ )

# Scenario Templates

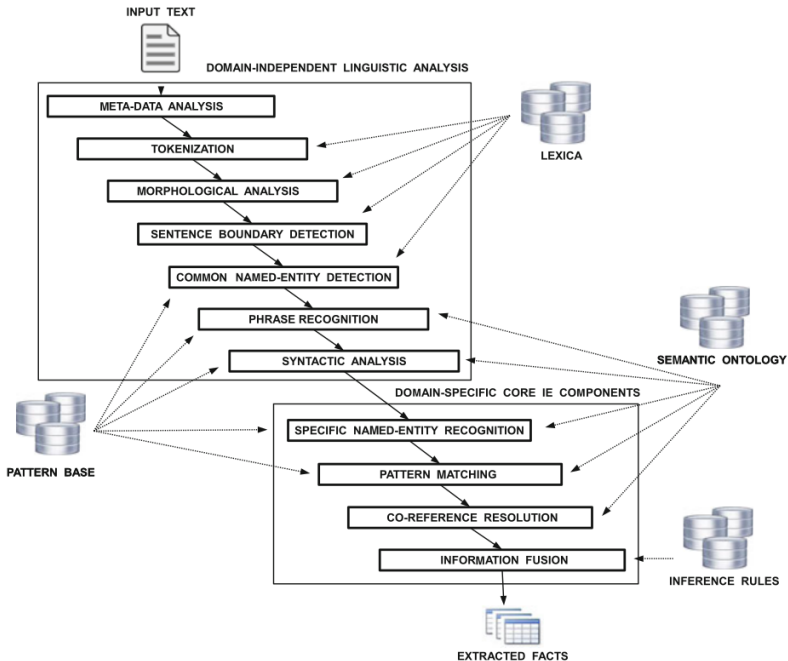
prototypical outputs

- precision–recall trade-off
- other evaluation metric: slot error rate

$$S = \frac{\textit{incorrect} + \textit{missing}}{\textit{key}},$$

where *incorrect* is the number of incorrectly assigned slots,  
*missing* is the number of missing slots,  
and *key* is the total number of slots.

Best MUC results:  $\approx 60\%$  (humans  $\approx 80\%$ )



# Accuracy

- Still not very consistent evaluation metrics
- General texts
  - ▶ “fill in the gaps” task (as in MUCs): around 60 %
  - ▶ EFa – precision of phrase detection and classification: 70 %
  - ▶ far from reliable and usable analysis
  - ▶ OIE reports over 80 % *precision*
- Specialized systems
  - ▶ simpler task, e.g. only dates, places, ...
  - ▶ e.g. Amy, the automated personal assistant
  - ▶ much better, human level accuracy



# Information extraction: Summary

- extracting structured information from text
- named entity detection + coreference resolution + relation extraction
- event recognition = domain specific, task specific
- successful in very specialized tasks, not very usable in general tasks

## Trends:

- social media
- cross-lingual extraction
- open (general) domain

# Information Extraction Systems

- Open Information Extraction (OIE), or TextRunner
  - ▶ <http://openie.allenai.org>
  - ▶ 100 million web pages
  - ▶ 500 million assertions
- GATE – general architecture for text engineering
  - ▶ <http://gate.ac.uk>
  - ▶ huge system for language annotation and all levels of automatic processing
  - ▶ contains a customizable information extraction component
- EFa – Extraction of Facts
  - ▶ <http://nlp.fi.muni.cz/projects/set/efa>
  - ▶ in NLP centre at FI
  - ▶ analysis of running text
  - ▶ syntactic analysis
  - ▶ phrase detection
  - ▶ semantic classification of phrases

# References I



Banko, M., Cafarella, M. J., Soderland, S., Broadhead, M., and Etzioni, O. (2007).

Open information extraction for the web.

*IJCAI*, 7:2670–2676.



Chang, C.-H., Kayed, M., Girgis, M. R., and Shaala, K. F. (2006).

A survey of web information extraction systems.

*Knowledge and Data Engineering, IEEE Transactions on*, 18(10):1411–1428.



Cunningham, H. (2005).

Information Extraction, Automatic.

*Encyclopedia of Language and Linguistics, 2nd Edition*.

## References II



Fader, A., Soderland, S., and Etzioni, O. (2011).  
Identifying relations for open information extraction.  
*In Proceedings of the Conference on Empirical Methods in Natural Language Processing, EMNLP '11*, pages 1535–1545, Stroudsburg, PA, USA. Association for Computational Linguistics.



Mitkov, R. (2005).  
*The Oxford handbook of computational linguistics*.  
Oxford University Press.



Piskorski, J. and Yangarber, R. (2013).  
*Information Extraction: Past, Present and Future*, pages 23–49.  
Springer Berlin Heidelberg, Berlin, Heidelberg.