



Vincent Kríž

# Detecting Entity Relations in Texts

Intelligent library (INTLIB, TA02010182)

KEG Seminar, 2013-12-19

Institute of Formal and Applied Linguistics  
Faculty of Mathematics and Physics  
Charles University in Prague  
Czech Republic

kriz@ufal.mff.cuni.cz  
<http://ufal.mff.cuni.cz/~kriz>

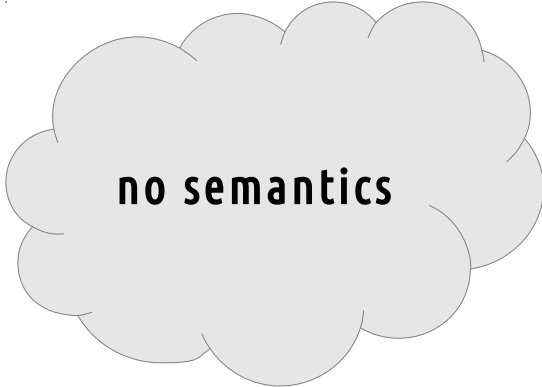
# Outline

- INTLIB
- Legislative documents
  - Structure
  - Semantics
    - JTagger
    - JStories
    - Rextractor
  - Improving syntax

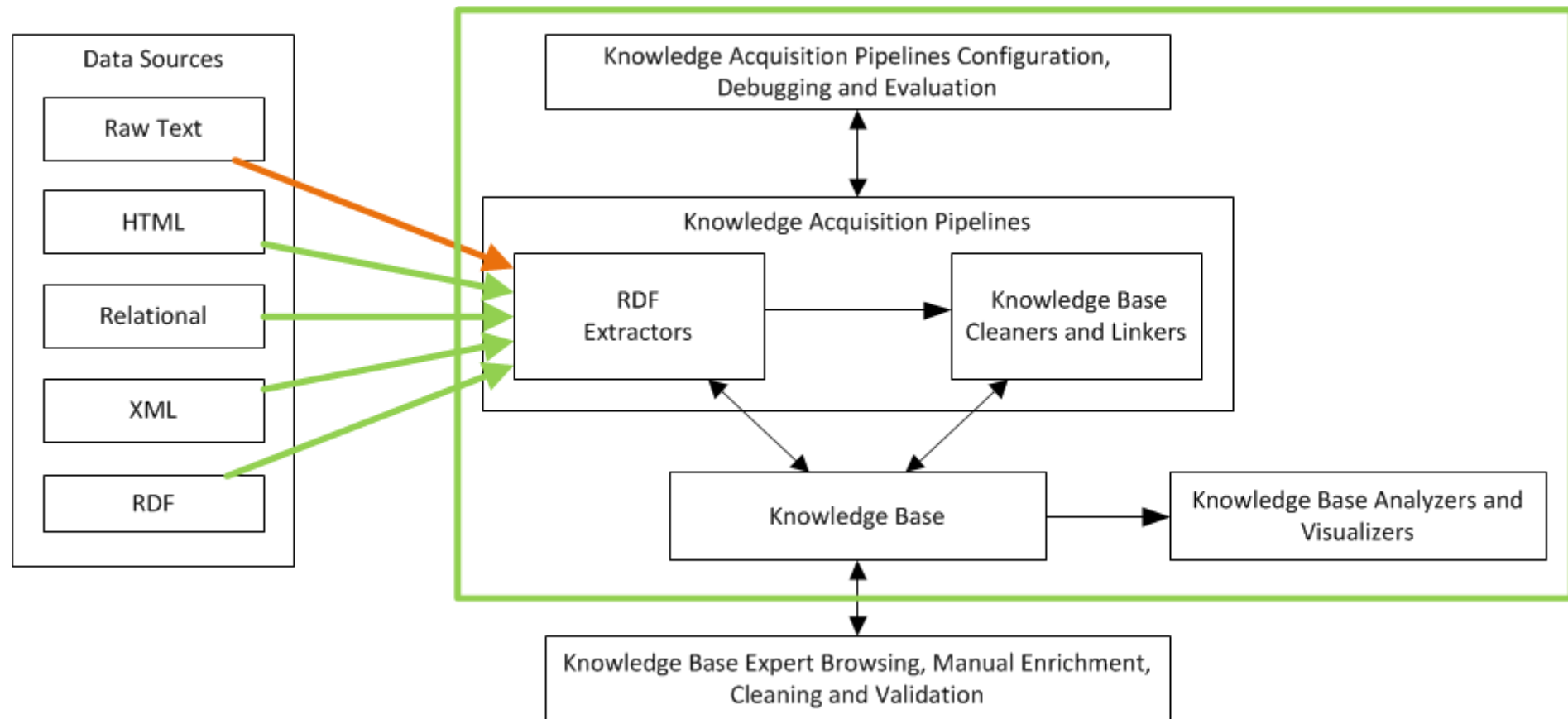
# Outline

- **INTLIB**
- Legislative documents
  - Structure
  - Semantics
    - JTagger
    - JStories
    - Rextractor
  - Improving syntax

# Motivation

- large collections of documents
  - efficient browsing & querying
  - typical approaches
    - full-text search
    - metadata search
- 
- A thought bubble with a grey fill and a black outline, containing the text "no semantics". It is connected to the "full-text search" and "metadata search" items by three small circles of increasing size.
- semantics interpretation of documents →  
suitable DB & query language →  
user-friendly browsing & querying

# INTLIB



# NLP Group

- **Documents**

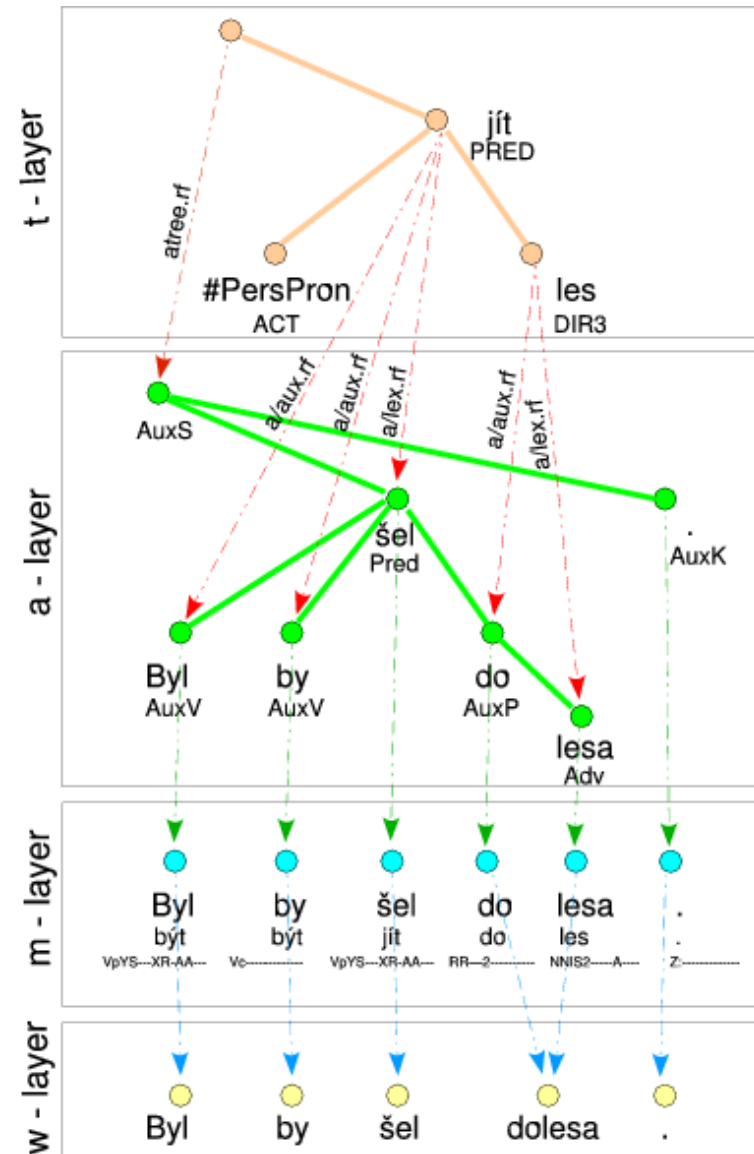
- semi-structured documents from some domain
- legislative documents, project/medical documentation

- **Extractor**

- NLP techniques

- **Data » Linked data**

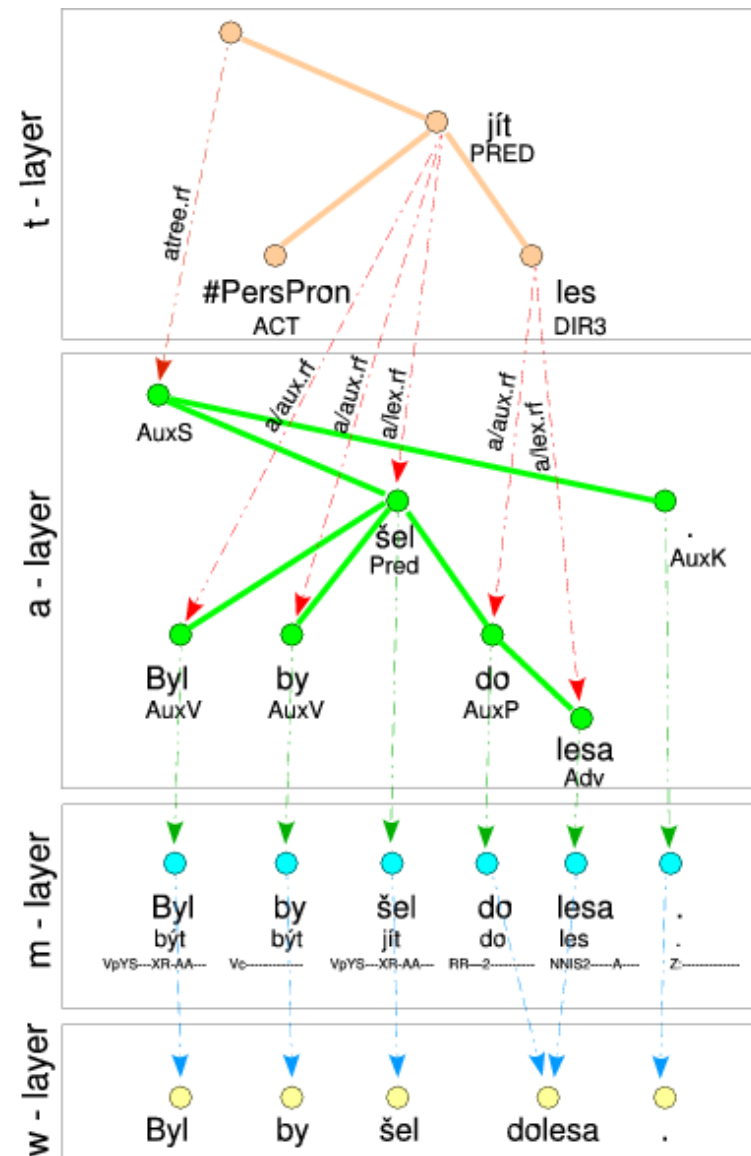
- automatically interconnected with other related data and with the original documents



# NLP Group

## • Tools

- segmentation & tokenization
- lemmatization & morphology
- syntactic parsing
- deep syntactic parsing

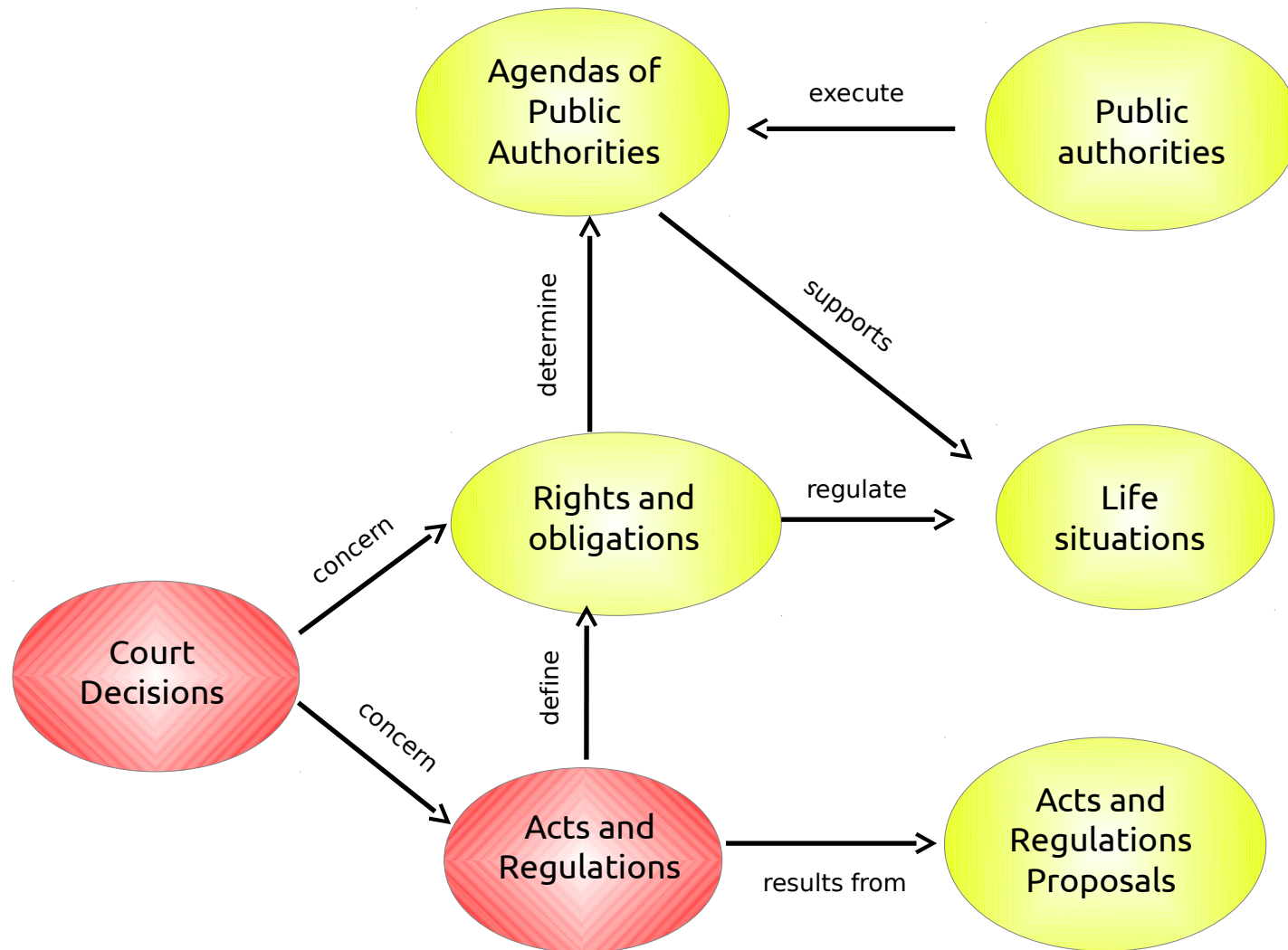


# Outline

- INTLIB
- Legislative documents
  - Structure
  - Semantics
    - Entities
      - JTagger
    - Relations
      - Rextractor
      - JStories
  - Improving syntax



# Legislation domain



# Outline

- INTLIB
- Legislative documents
  - Structure
  - Semantics
    - JTagger
    - JStories
    - Rextractor
  - Improving syntax

# Legislation domain - structure

## What we have done

- ontology of legislative documents
- metadata and structure of acts, regulations and decrees represented as Linked Data
  - metadata about each version of each act, regulation and decree since 1945
  - structured content of versions of all acts, regulations and decrees valid in 2011, 2012

# Legislation domain - structure

HLAVA I  
ÚVODNÍ USTANOVENÍ

§ 1  
Předmět úpravy

Tato vyhláška zpracovává příslušné předpisy Evropské unie a upravuje:

- a) způsob vymezení hydrogeologických rajonů, vymezení útvarů podzemních vod,
- b) způsob hodnocení stavu podzemních vod a
- c) náležitosti programů zjišťování a hodnocení stavu podzemních vod.

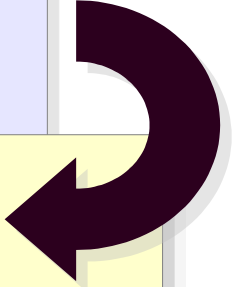
# Legislation domain - structure

HLAVA I  
ÚVODNÍ USTANOVENÍ

§ 1  
Předmět úpravy

Tato  
a  
b  
c

```
<head id="11" label="HLAVA I">
  <title>ÚVODNÍ USTANOVENÍ</title>
  <section id="12" label="§ 1">
    <title>Předmět úpravy</title>
    <text>Tato vyhláška zpracovává příslušné
      předpisy Evropské unie a upravuje:</text>
    <section id="13" label="a)">
      <text>způsob vymezení hydrogeologických rajonů,
        vymezení útvarů podzemních vod,</text>
    </section>
    <section id="14" label="b)">
      <text>způsob hodnocení stavu podzemních vod a</text>
    </section>
    <section id="15" label="c)">
      <text>náležitosti programů zjišťování a
        hodnocení stavu podzemních vod.</text>
    </section>
  </section>
</head>
```



# Outline

- INTLIB
- Legislative documents
  - Structure
  - Semantics
    - JTagger
    - Jstories
    - Rextractor
  - Improving syntax

# Legislation domain - semantics

Extracting **concepts** and **relationships** between them from documents

- court decisions
  - **Entities:** references, institutions, acts, dates
  - **Relations:** whole *story* of a case
- acts, regulations, ...
  - **Entities:** subjects, things, locations, ...
  - **Relations:** rights, obligations, ...

# Outline

- INTLIB
- Legislative documents
  - Structure
  - Semantics
    - JTagger
    - JStories
    - Rextractor
  - Improving syntax



# JTagger

- **Entities**

- References on
  - court decisions
  - acts
- Effectiveness of Act
- Institutions

- **Relations**

- Publisher
  - Institution → Decision
- Abbreviation

# JTagger

- Annotation in **Brat** (<http://brat.nlplab.org>)

The Constitutional Court states, first, that the identical legal issue addressed the position taken by the Plenum of the Constitutional Court on 28th April 2009 file no. Pl. US-st 27/09 (ST 27/53 SbNU 885; 136/2009 Coll.). Here said ... because from that date a unilateral increase rent allowed by § 3, paragraph 2 of Act No. 107/2006 Coll. Unilateral Increase of Rent and Amending Act No. 40/1964 Coll., the Civil Code, as amended.

Annotations in the image:

- Institution** (yellow box) above "The Constitutional Court" and "Plenum of the Constitutional Court".
- publisher** (yellow box) above the arrow connecting the two "Institution" entities.
- Decision** (green box) above "file no. Pl. US-st 27/09".
- Act** (red box) above "136/2009 Coll.", "§ 3, paragraph 2 of Act No. 107/2006 Coll.", and "Act No. 40/1964 Coll.".
- Effectiveness** (yellow box) above "as amended".

# JTagger

- **Data sets**
  - Corpus of manually annotated court decisions
    - The Supreme Court (150)
    - The Constitutional Court (150)

	SC			CC		
	# of docs	# of tokens	# of entities	# of docs	# of tokens	# of entities
Training set	135	332,535	8,487	135	312,191	7,910
Test set	15	36,999	943	15	34,701	879
Total	150	369,534	9,430	150	346,892	8,789

# JTagger

- **Data sets**
  - Corpus of manually corrected court decisions
    - The Supreme Court (93)
    - The Constitutional Court (91)

	SC			CC		
	# of docs	# of tokens	# of entities	# of docs	# of tokens	# of entities
<b>Total</b>	93	120,856	6,047	91	100,464	4,945

# JTagger

- **Machine Learning experiments**
  - Hidden Markov models (HMM)
  - Perceptron Algorithm with Uneven Margins (PAUM)

# JTagger

- **Hidden Markov models (HMM)**

- pattern recognition - speech, handwriting, gesture recognition, **part-of-speech tagging**, ...

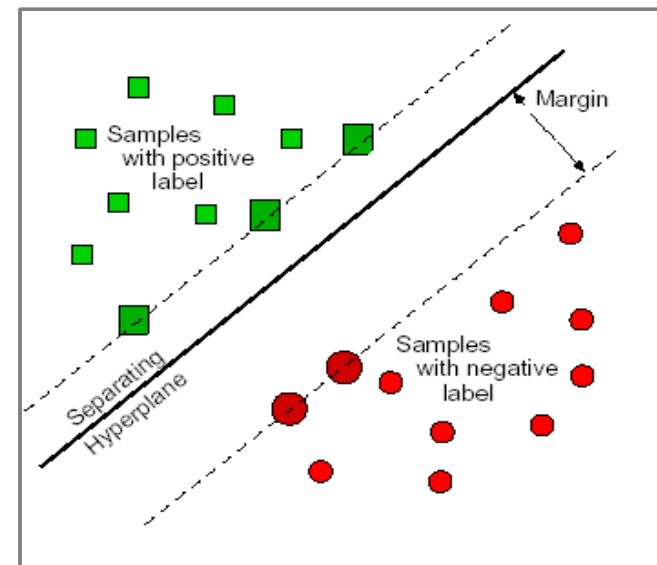
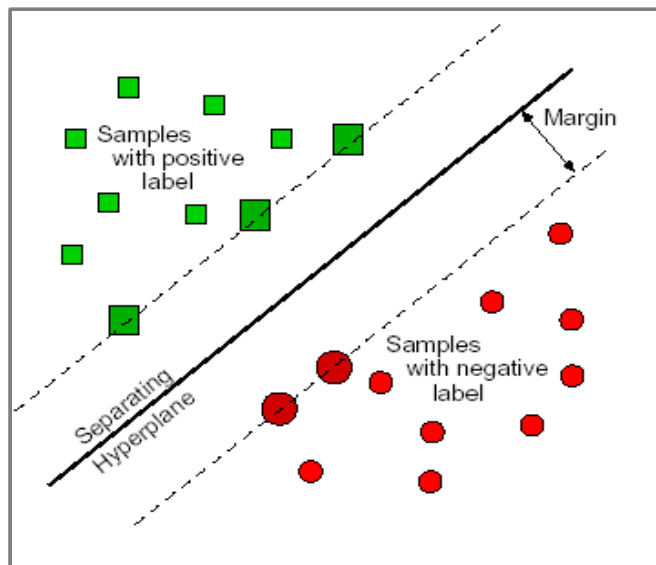
the	Plenum	of	the	Constitutional	Court	on	28th	April
DT	NNP	IN	DT	NNP	NNP	IN	JJ	NNP
NONE	NONE	NONE	INTS	INST	INST	NONE	NONE	NONE

- **noisy channel**



# JTagger

- **Perceptron Algorithm with Uneven Margins (PAUM)**
  - Implementation in the GATE framework
  - Perceptron and Support Vector Machines
- PAUM doesn't position the separator right between the points, but over one side



# JTagger

- Models
  - **HMM**
  - **PM Small**
    - trigrams of word forms
  - **PM**
    - 5-grams of word forms
  - **PM POS**
    - 5-grams of lemmas and POS-tags
  - **PM POS EXT**
    - PM POS + orthography features



# JTagger

## Strict F1 on entities

	Entity	HMM	PM pos ext	PM pos	PM	PM small
SC	A	0,75±0,02	0,91±0,02	0,91±0,03	0,89±0,03	0,88±0,03
	D	0,82±0,08	0,97±0,02	0,96±0,02	0,95±0,03	0,94±0,02
	E	0,89±0,04	0,90±0,05	0,89±0,05	0,88±0,08	0,82±0,1
	I	0,92±0,03	0,96±0,02	0,96±0,02	0,95±0,02	0,96±0,02
CC	A	0,63±0,05	0,87±0,02	0,86±0,02	0,84±0,03	0,78±0,03
	D	0,83±0,05	0,95±0,03	0,95±0,03	0,93±0,03	0,92±0,03
	E	0,96±0,03	0,96±0,03	0,96±0,03	0,96±0,03	0,96±0,03
	I	0,91±0,02	0,93±0,02	0,93±0,02	0,92±0,01	0,92±0,01

# JTagger

## Lenient F1 on entities

	Entity	HMM	PM pos ext	PM pos	PM	PM small
SC	A	0,93±0,02	0,96±0,01	0,96±0,01	0,95±0,01	0,95±0,02
	D	0,91±0,03	0,98±0,01	0,97±0,02	0,96±0,02	0,95±0,02
	E	0,94±0,04	0,91±0,05	0,90±0,05	0,90±0,06	0,83±0,1
	I	0,97±0,01	0,98±0,00	0,98±0,01	0,97±0,01	0,97±0,01
CC	A	0,89±0,02	0,94±0,01	0,94±0,01	0,94±0,01	0,93±0,02
	D	0,93±0,03	0,97±0,02	0,97±0,02	0,96±0,02	0,95±0,03
	E	0,96±0,03	0,96±0,03	0,96±0,03	0,96±0,03	0,96±0,03
	I	0,97±0,01	0,98±0,01	0,98±0,01	0,97±0,01	0,97±0,01

# JTagger

- On-line DEMO
  - <http://ufal.mff.cuni.cz/jtagger>
- Open data
  - JTagger as a component of ODCleanStore
    - <http://sourceforge.net/projects/odcleanstore/>
  - daily, fully automatic
  - processing and publication of the new decisions

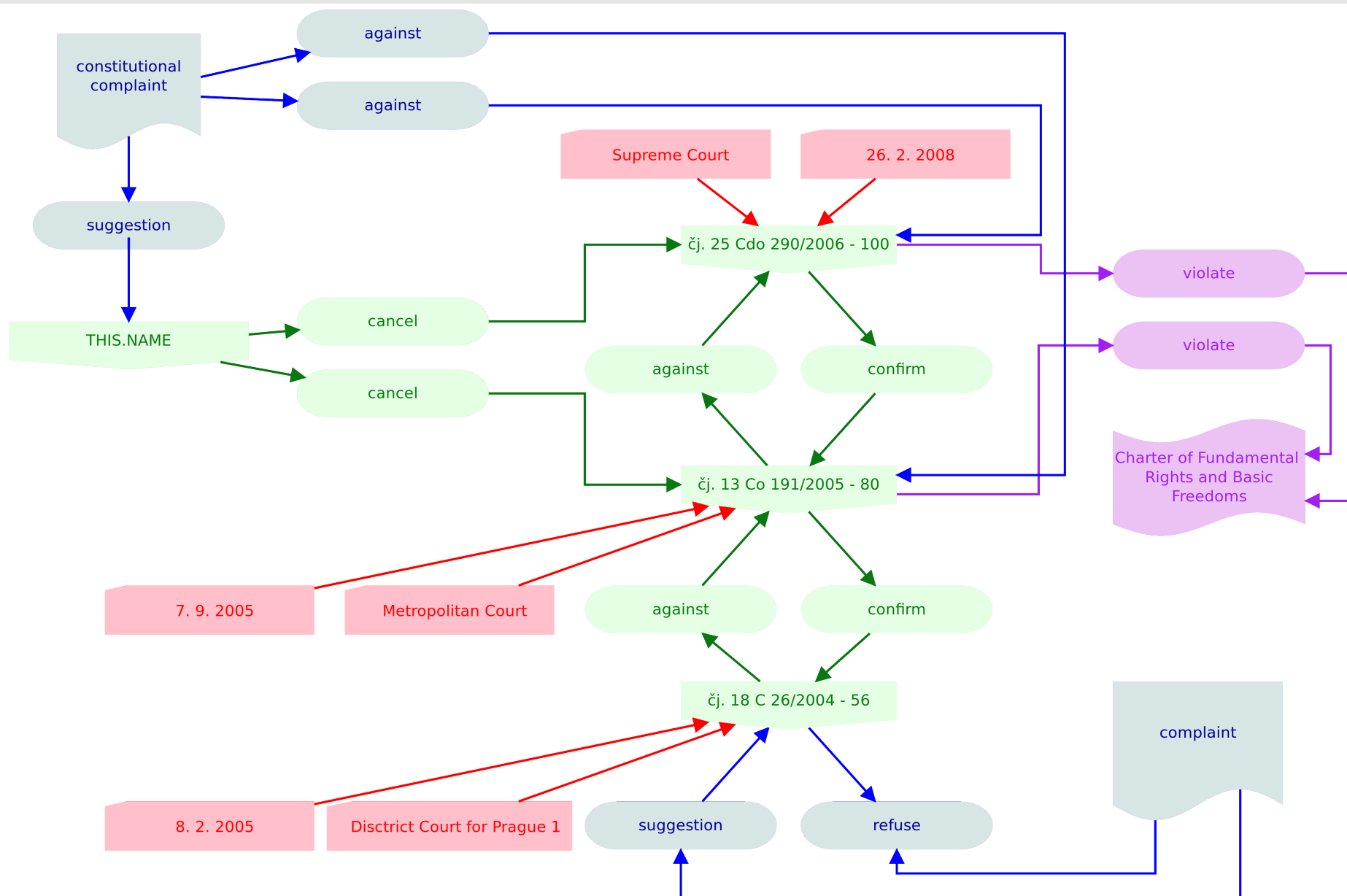
# Outline

- INTLIB
- Legislative documents
  - Structure
  - Semantics
    - JTagger
    - JStories
    - Rextractor
  - Improving syntax

# JStories

- Legal case retrospection
  - a story
  - begins when an accuser files a complaint to a court
  - ends when a final decision is rendered

# JStories




# Outline

- INTLIB
- Legislative documents
  - Structure
  - Semantics
    - JTagger
    - JStories
    - Rextractor
  - Improving syntax

# RExtractor

- General framework for information extraction from texts
- Entities detection and identifying relationships between them
- Pipeline:
  - syntactic analysis of the text
  - entities detection (in syntactic trees)
  - PML-TQ for relations

[http://prezi.com/jvrkcl1ui-ug/?utm\\_campaign=share&utm\\_medium=copy&rc=ex0share](http://prezi.com/jvrkcl1ui-ug/?utm_campaign=share&utm_medium=copy&rc=ex0share)



pipeline  
visualization  
@prezi.com




# Outline

- INTLIB
- Legislative documents
  - Structure
  - Semantics
    - JTagger
    - JStories
    - Rextractor
  - Improving syntax

# Improving syntactic parsing

- Parsers trained on PDT data
  - How they work on legal texts?
  - Too many nodes make automatic parsing almost impossible
- Manual syntactic analysis of Czech legal texts

[http://prezi.com/m4r9r1nhzi7n/?utm\\_campaign=share&utm\\_medium=copy&rc=ex0share](http://prezi.com/m4r9r1nhzi7n/?utm_campaign=share&utm_medium=copy&rc=ex0share)



pipeline  
visualization  
@prezi.com

# Conclusion

- INTLIB
- Legislative documents
  - Structure
  - Semantics
    - JTagger
    - JStories
    - Rextractor
  - Improving syntax