



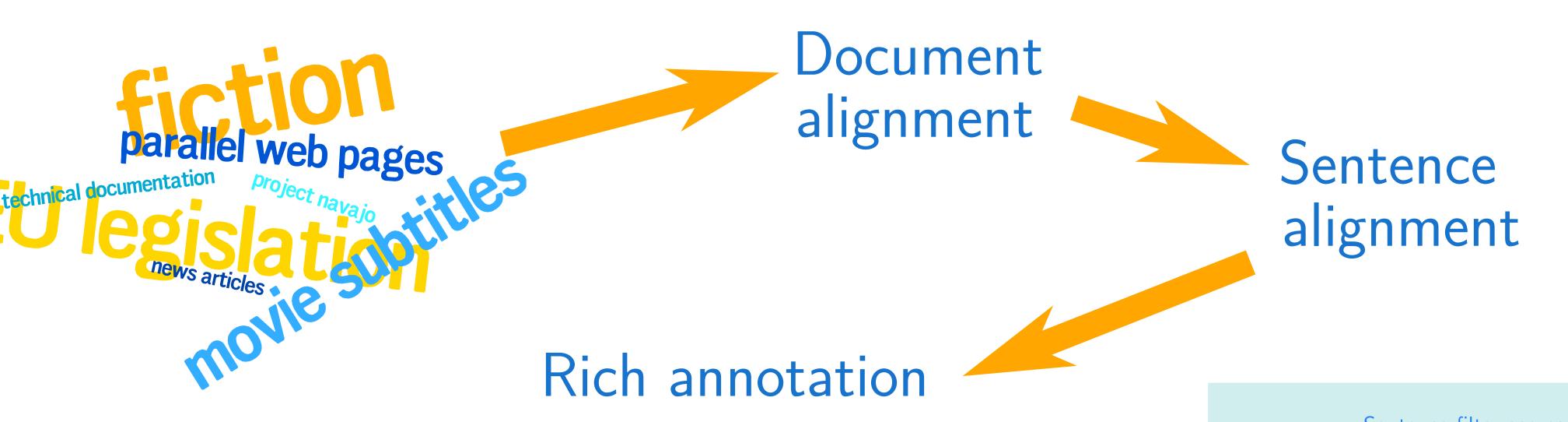


Ondřej Bojar, Zdeněk Žabokrtský, Ondřej Dušek, Petra Galuščáková, Martin Majliš, David Mareček, Jiří Maršík, Michal Novák, Martin Popel, Aleš Tamchyna *surname*@ufal.mff.cuni.cz, *except* {mnovak,odusek}@ufal.mff.cuni.cz, jiri.marsik89@gmail.com Charles University in Prague, Faculty of Mathematics and Physics, Institute of Formal and Applied Linguistics

http://ufal.mff.cuni.cz/czeng

CzEng 1.0

- updated release of our Czech-English parallel corpus
- freely available for non-commercial research and educational purposes
- 15 million sentence pairs (doubled since last release)
- ca. 200 million words per language
- improved automatic rich annotation part-of-speech tags, morphology surface syntax dependency trees (a-)



- deep syntax dependency trees (t-)
- co-reference links
- carefully filtered to reduce the amount of non-matching sentence pairs
- automatic sentence and word alignment

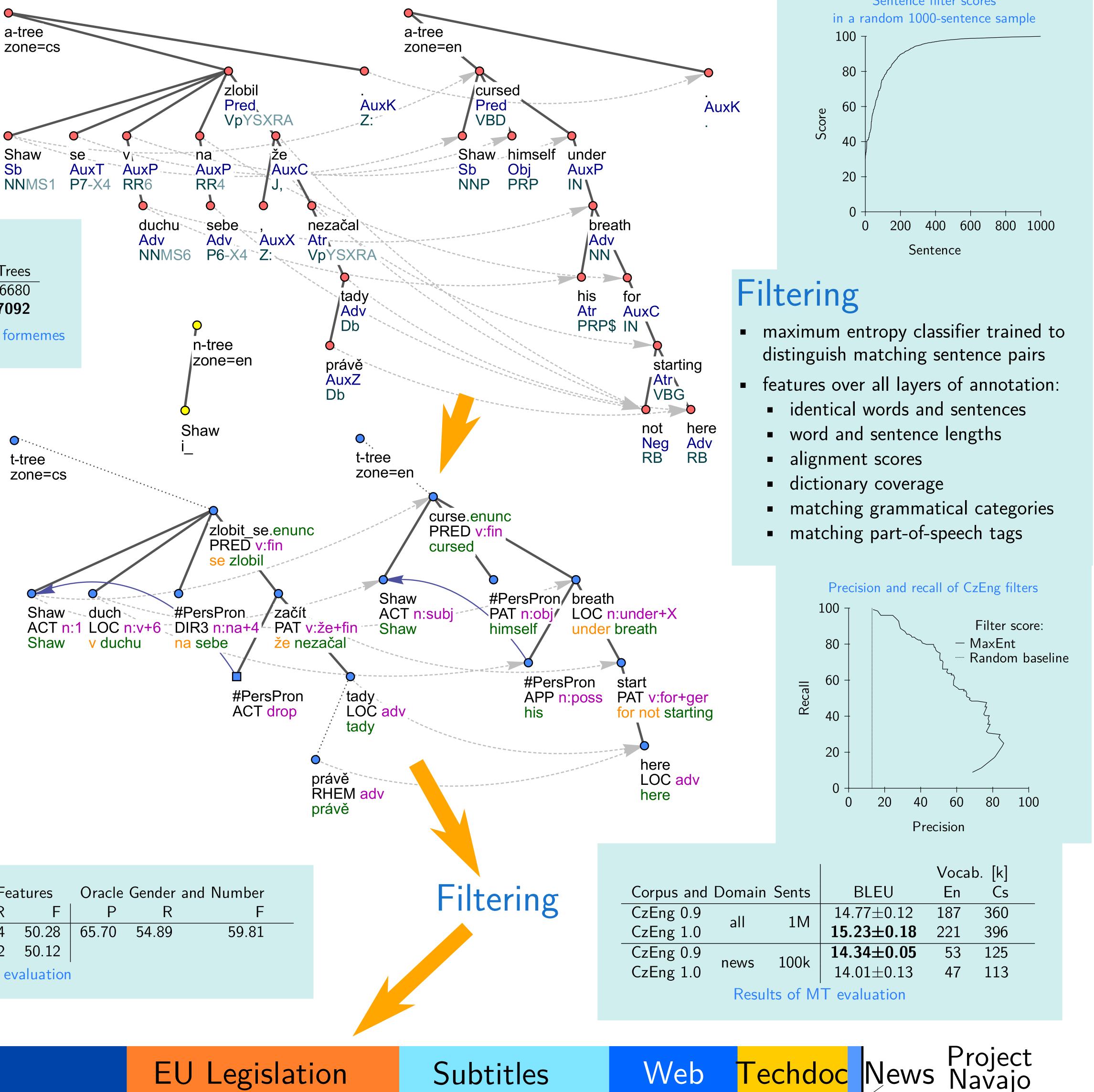
	Formeme Detection on					
	Automatic Trees	Manual Trees				
Baseline	1.5981	1.6680				
Improved	1.6873	1.7092				

Inter-lingual pointwise mutual information on formemes

Formemes

v:fin v:for+ger n:obj n:na+4

- description of morpho-syntactic form in t-tree nodes, including:
 - syntactic position and usage, morphological case
 - prepositions and conjunctions
- improved version aimed at greater cross-language consistency



Subtitles

3.1 M

Sentence filter scores

Techdoc News

0.2 M

-0.03 M

0.5 M

0.6 M

1.6 M

Web

1.9 M

Co-reference

- automatic co-reference links in t-trees for both languages
- using rules and perceptron ranking
- grammatical co-reference
 - relative and reflexive pronouns

Sent

pairs

verbs of control

Total

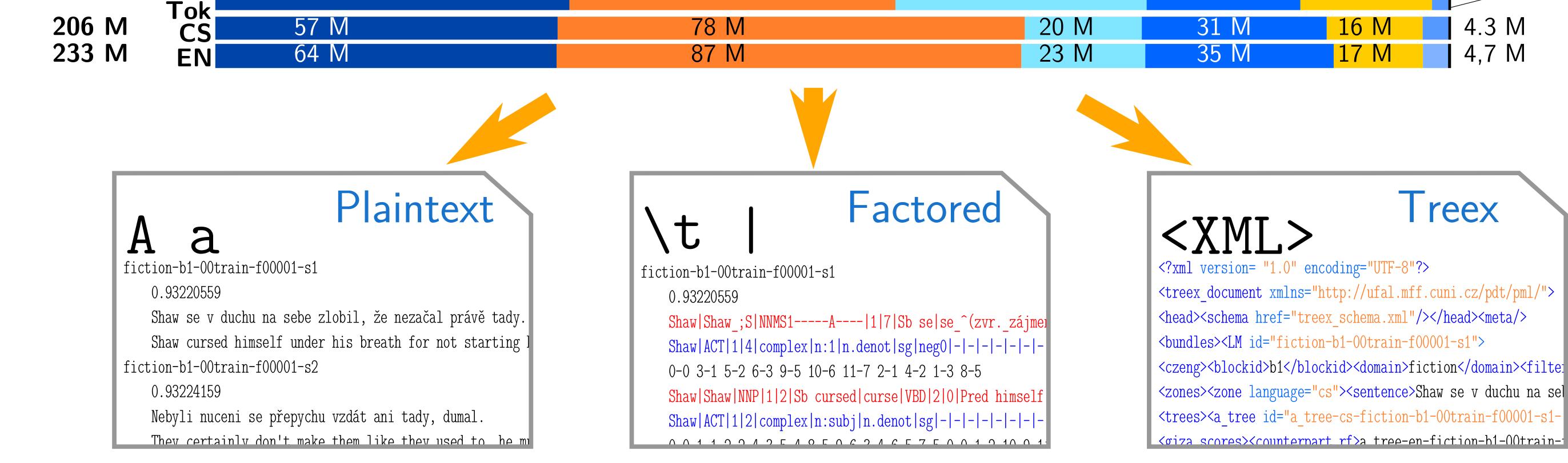
15.1 M

textual pronominal co-reference

	Gold Standard Features			Automatic Features			Oracle Gender and Number			
Language	P	R	F	P	R	F	P	R	F	
Czech	77.06	77.58	77.32	55.23	46.14	50.28	65.70	54.89	59.81	
English	45.52	58.69	51.27	44.53	57.32	50.12				
Co-reference resolution evaluation										

Fiction

4.3 M



EU Legislation

4.0 M

This work was supported by the project EuroMatrixPlus (FP7-ICT-2007-3-231720 of the EU and 7E09003+7E11051 of the EU and 7E09003+7E11051 of the Czech Republic), Czech Science Foundation grants P406/10/P259 and 201/09/H057, GAUK 4226/2011, and the FAUST project (FP7-ICT-2009-4-247762 of the EU and 7E11041 of the Czech Republic). This work has been using language resources developed and/or stored and/or distributed by the LINDAT-Clarin project of the Ministry of Education of the Czech Republic (project LM2010013). This poster was presented at LREC 2012.