**Rudolf Rosa**, David Mareček, Aleš Tamchyna
{rosa,marecek,tamchyna}@ufal.mff.cuni.cz

# Deepfix:

# Statistical Post-editing of Statistical Machine Translation Using Deep Syntactic Analysis

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

ACL SRW, Sofia, 6th August 2013

# Motivation

- Source text in English:

  *EU criticizes not only the Greek government.*

# Motivation

- Source text in English:

  *EU criticizes not only the Greek government*

- Google Translate to Czech (6[th] Aug 2013):

  *EU kritizuje nejen řecká vláda*

# Motivation

- Source text in English:

  *EU criticizes not only the Greek government*

- Google Translate to Czech (6$^{th}$ Aug 2013):

  *EU kritizuje nejen **řecká vláda**nominative (subject)*

  - *Not only **the Greek government** criticizes EU*

# Motivation

- Source text in English:

*EU criticizes not only the Greek government*

- Google Translate to Czech (6[th] Aug 2013):

*EU kritizuje nejen* **řeck*á* vlád*a***_nominative (subject)_

  - *Not only* **the Greek government** *criticizes EU*

- Post-editation by Deepfix:

*EU kritizuje nejen* **řeck*ou* vlád*u***_accusative (object)_

  - *EU criticizes not only* **the Greek government**

# Outline

1. Problem definition

2. Sentence analysis

3. Sentence post-editing

4. Results

# Outline

1. Problem definition

   ➔ Errors in valency in SMT outputs

2. Sentence analysis


3. Sentence post-editing


4. Results

# Outline

1. Problem definition

   → Errors in valency in SMT outputs

2. Sentence analysis (DEEP)

   → Deep dependency parsing

3. Sentence post-editing

4. Results

# Outline

1. Problem definition

   ➔ Errors in valency in SMT outputs

2. Sentence analysis (DEEP)

   ➔ Deep dependency parsing

3. Sentence post-editing (FIX)

   ➔ Statistical model of valency

4. Results

# Outline

1. **Problem definition**

   ➔ Errors in valency in SMT outputs

2. **Sentence analysis (DEEP)**

   ➔ Deep dependency parsing

3. **Sentence post-editing (FIX)**

   ➔ Statistical model of valency

4. **Results**

   ➔ Automatic & manual evaluation of Deepfix

# Subject – object dichotomy

- English: **position** (left/right constituent)
  - Subject *criticize* Object
- Czech: **morphological case** (nominative/other); word order relatively free
  - Subject$_{nominative}$ *kritizovat* Object$_{accusative}$
  - Object$_{accusative}$ *kritizovat* Subject$_{nominative}$
  - Subject$_{nominative}$ Object$_{accusative}$ *kritizovat*
  - Object$_{accusative}$ Subject$_{nominative}$ *kritizovat*

# Valency of *criticize* (*kritizovat*)

- example sentence

  - *EU*subject *criticizes not only the Greek government*object

  - *EU*nominative *kritizuje nejen řeckou vládu*accusative

# Valency of *criticize* (*kritizovat*)
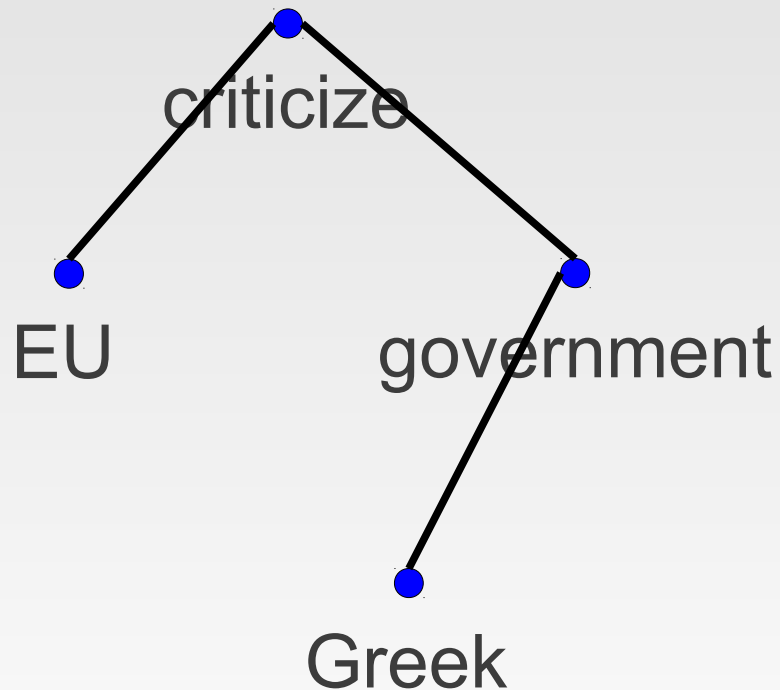
- example sentence

  - *EU*<sub>subject</sub> <u>*criticizes*</u> *not only the Greek government*<sub>object</sub>

  - *EU*<sub>nominative</sub> <u>*kritizuje*</u> *nejen řeckou vládu*<sub>accusative</sub>

- a valency frame of a verb

  - subject        *criticize*     object
  - nominative     *kritizovat*    accusative

# Valency of *criticize* (*kritizovat*)

- example sentence

  - *EU*subject *criticizes not only the Greek government*object

  - *EU*nominative *kritizuje nejen řeckou vládu*accusative

- a valency frame of a verb

  - subject  *criticize*  object  (position)
  - nominative  *kritizovat*  accusative  (cases)

# Valency of *criticize* (*kritizovat*)

- example sentence

  - *EU*<sub>subject</sub> *criticizes not only the Greek government*<sub>object</sub>

  - *EU*<sub>nominative</sub> *kritizuje nejen řeckou vládu*<sub>accusative</sub>

- a valency frame of a verb

  - subject          *criticize*     object              (position)
  - nominative     *kritizovat*    accusative      (cases)

- decomposition into head-argument pairs
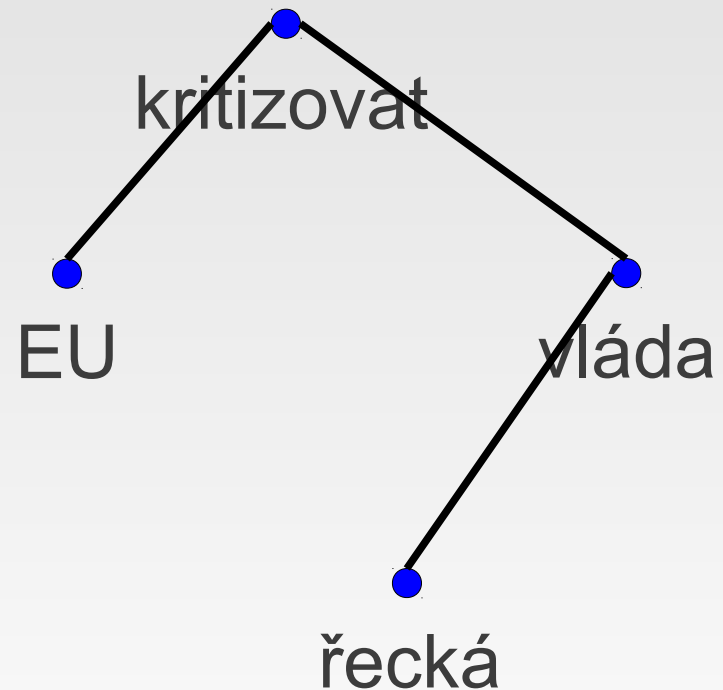
  - (*to criticize, government*) ~ (*kritizovat, vládu*)

# Sentence analysis (DEEP)

- tagging & lemmatization
  - combination of rule-based and statistical approach
- word-alignment
  - unsupervised methods (Giza++)
- dependency parsing
  - statistical, trained on manually created treebanks
  - parser adapted for parsing of SMT outputs
- induction of deep structure (tectogrammar)
  - rule-based

# Deep syntactic dependency trees
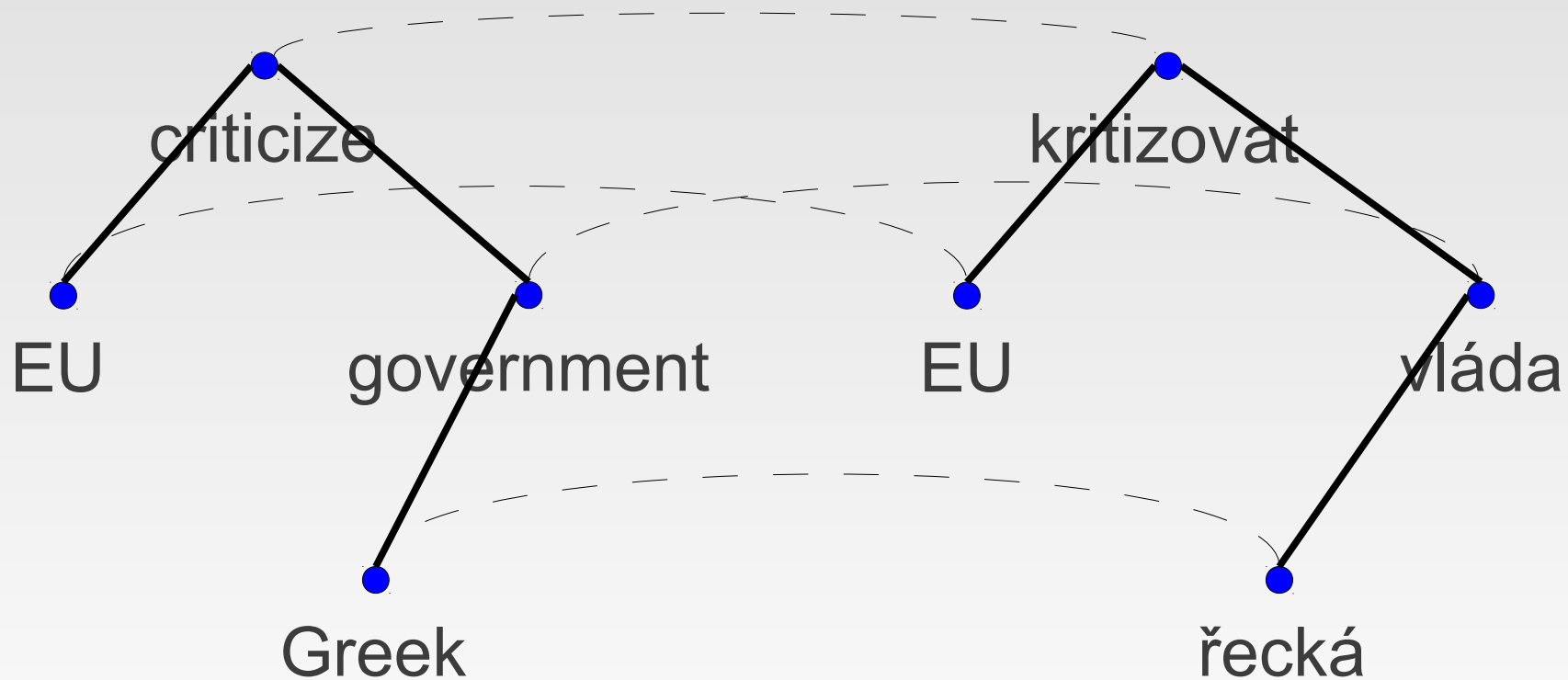
*EU criticizes*
*the Greek government*

*EU kritizuje*
*řecká vláda*

# Deep syntactic dependency trees



*EU criticizes
the Greek government*

*EU kritizuje
řecká vláda*

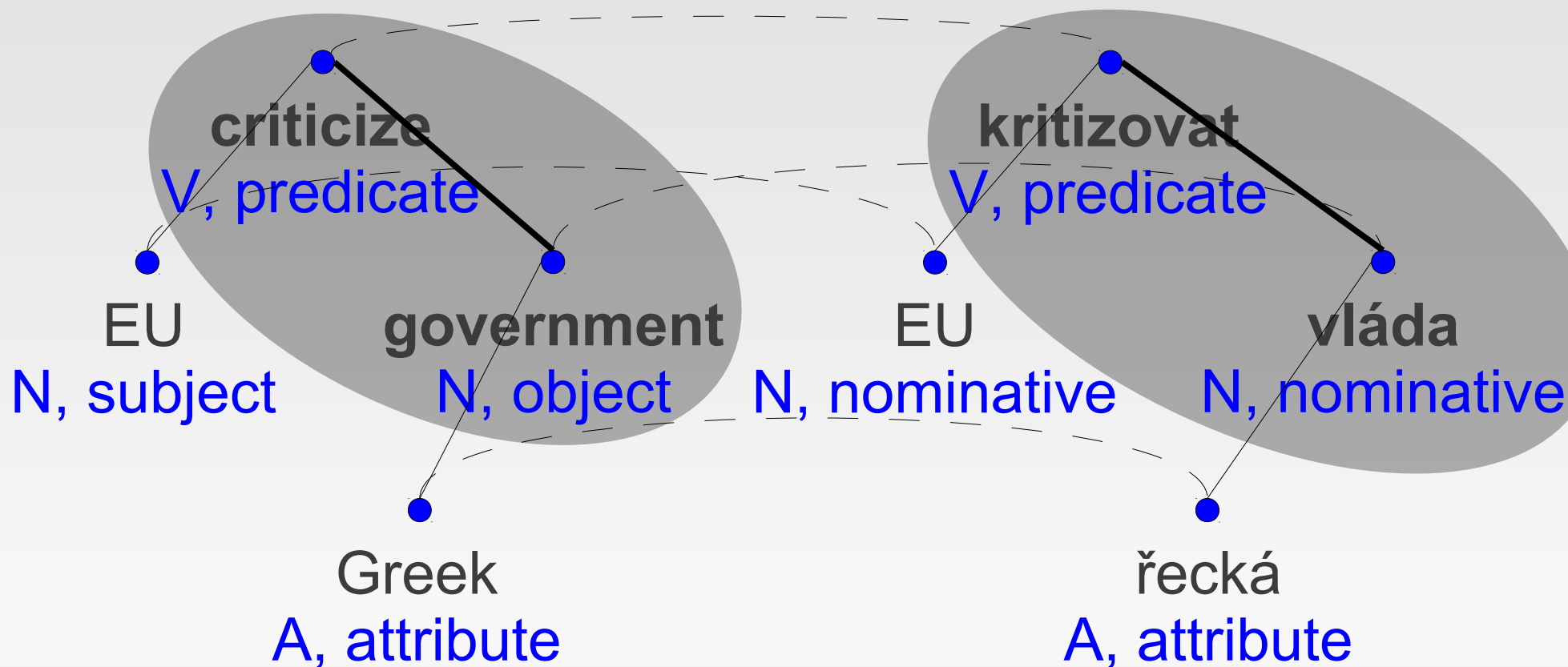# Deep syntactic dependency trees

EU criticizes
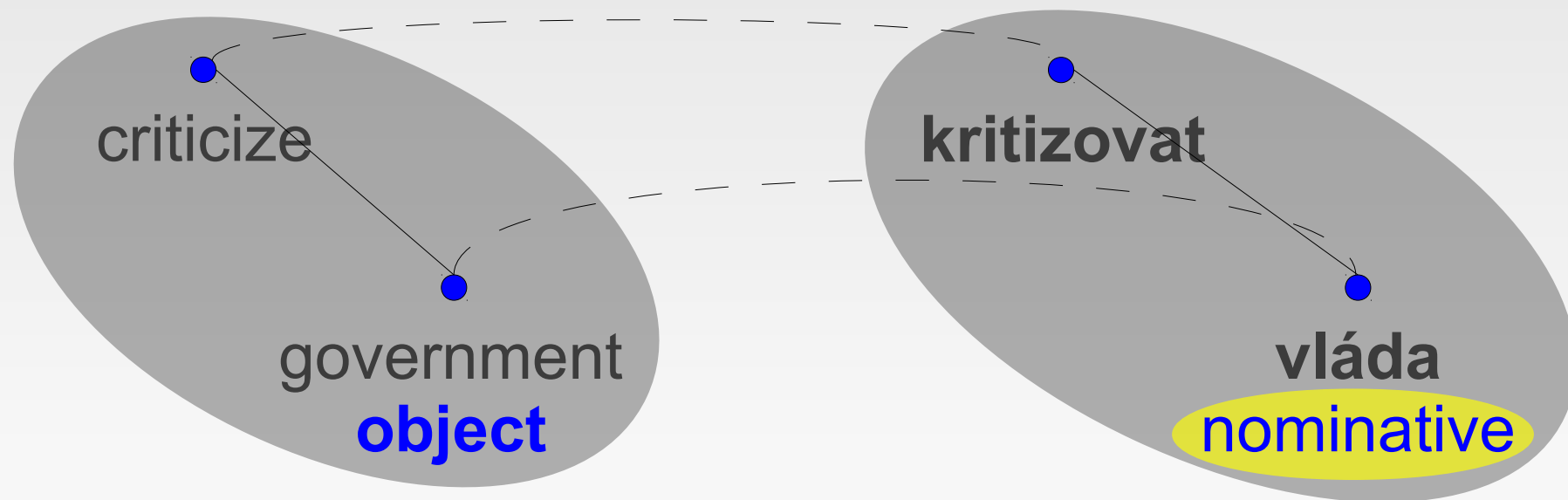the Greek government

EU kritizuje
řecká vláda



criticize
V, predicate

kritizovat
V, predicate

EU
N, subject

government
N, object

EU
N, nominative

vláda
N, nominative

Greek
A, attribute

řecká
A, attribute

# (head, arg) pair identification

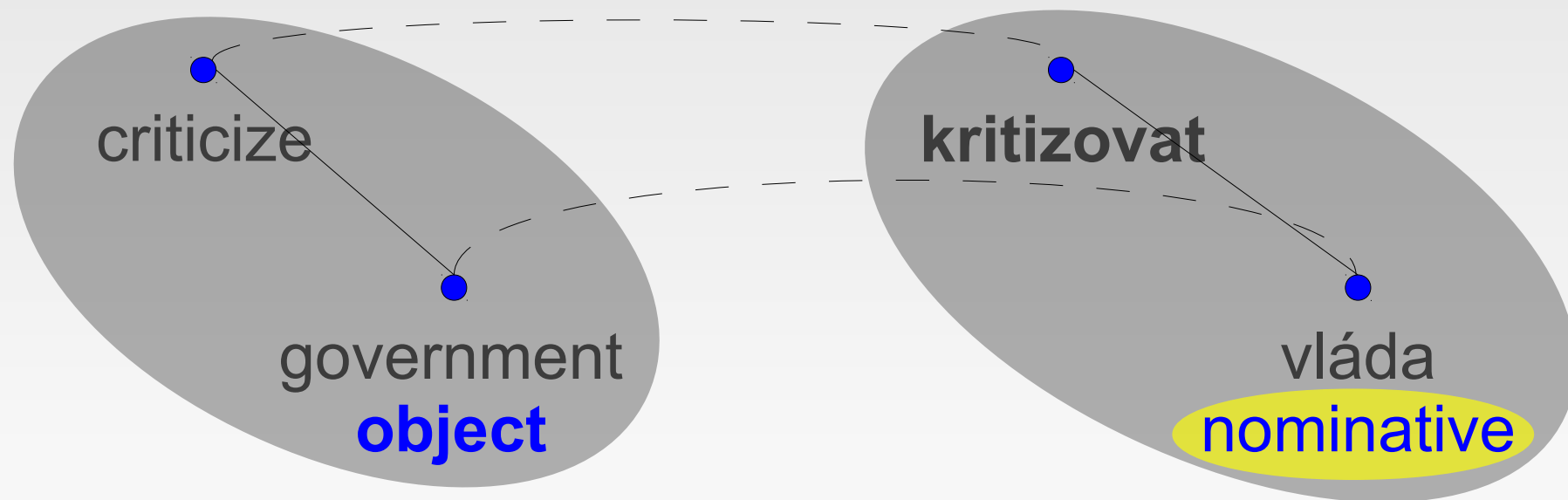# Valency models (FIX)

- $P(\text{arg}_{case} \mid \text{head}_{lemma}, \text{English\_arg}_{case})$

- $P(\text{arg}_{case} \mid \text{head}_{lemma}, \text{English\_arg}_{case}, \text{arg}_{lemma})$

- estimated from CzEng 1.0 (15M parallel stcs)



criticize

kritizovat
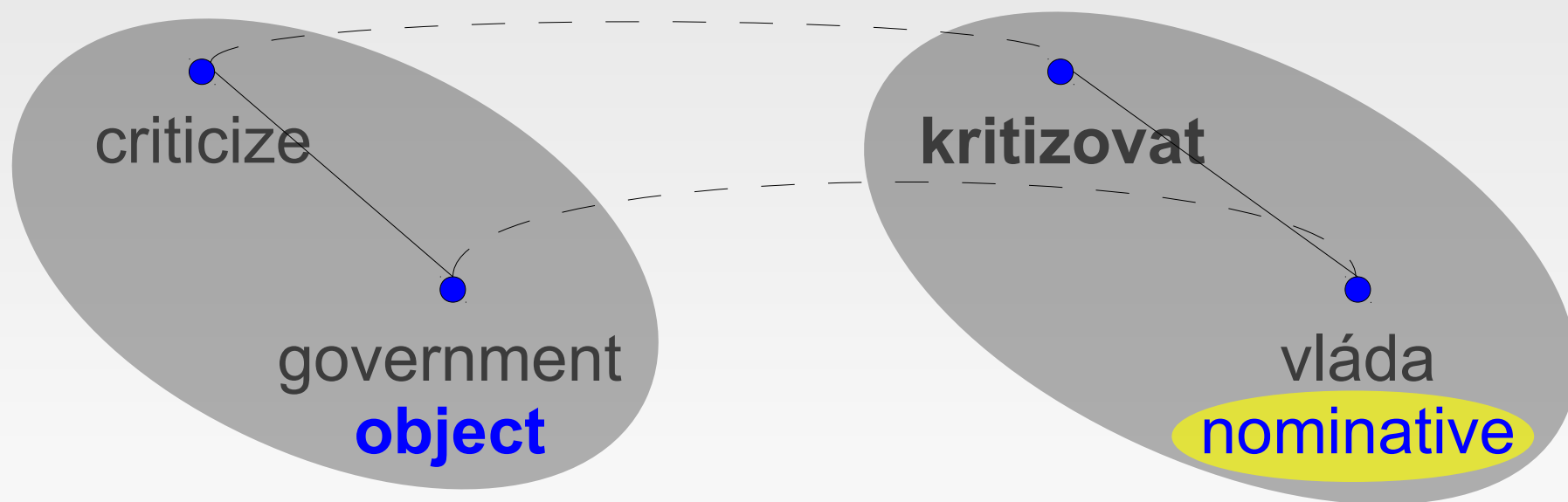
government
**object**

vláda
nominative

# Argument case probabilities

- P(nominative | *kritizovat*, object) = 0.03
- P(accusative | *kritizovat*, object) = 0.80



criticize

government
**object**

**kritizovat**

vláda
nominative

# Argument case probabilities

- P(nominative | *kritizovat*, object) = 0.03
- P(accusative | *kritizovat*, object) = 0.80
- threshold: 0.55

# Argument case correction

- P(nominative | *kritizovat*, object) = 0.03
- P(**accusative** | *kritizovat*, object) = **0.80**
- threshold: **0.55**

criticize

government
**object**

**kritizovat**

vláda
**accusative**

# Sentence correction

- Statitical machine translation output:

*EU kritizuje nejen řeck$a_{nominative}$ vlád$a_{nominative}$*

  - *Not only **the Greek government** criticizes EU*

# Sentence correction

- Statitical machine translation output:

  *EU kritizuje nejen* **řeck**$\textcolor{red}{\textbf{á}}$**<sub>nominative</sub>** **vlád**$\textcolor{red}{\textbf{a}}$**<sub>nominative</sub>**

  - *Not only* **the Greek government** *criticizes EU*

- Valency model correction:

  *EU kritizuje nejen* **řeck**$\textcolor{red}{\textbf{á}}$**<sub>nominative</sub>** **vlád**$\textcolor{green}{\textbf{u}}$**<sub>accusative</sub>**

# Sentence correction

- Statitical machine translation output:

  *EU kritizuje nejen* **řecká**<sub>nominative</sub> **vláda**<sub>nominative</sub>

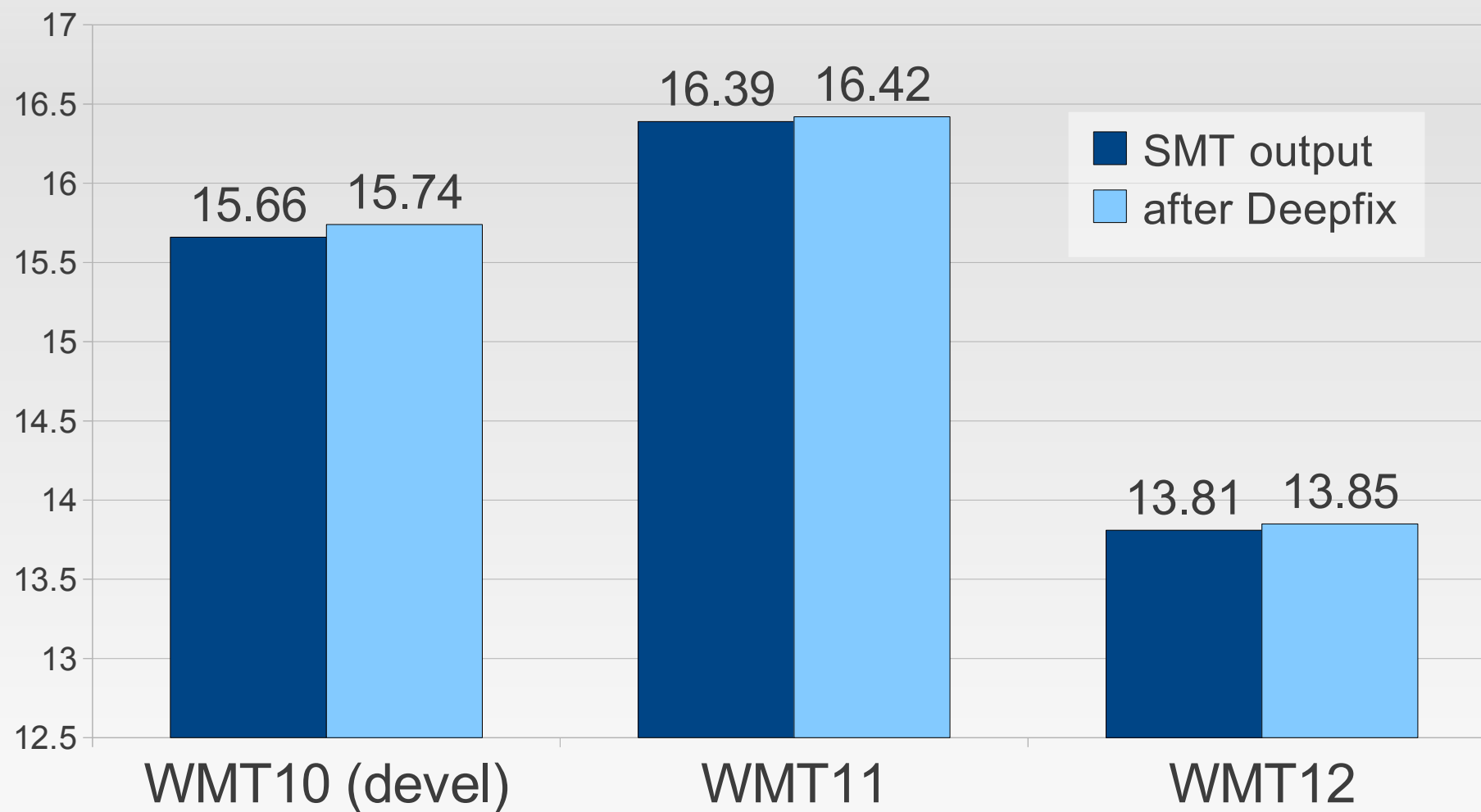  - *Not only* **the Greek government** *criticizes EU*

- Valency model correction:

  *EU kritizuje nejen* **řecká**<sub>nominative</sub> **vládu**<sub>accusative</sub>
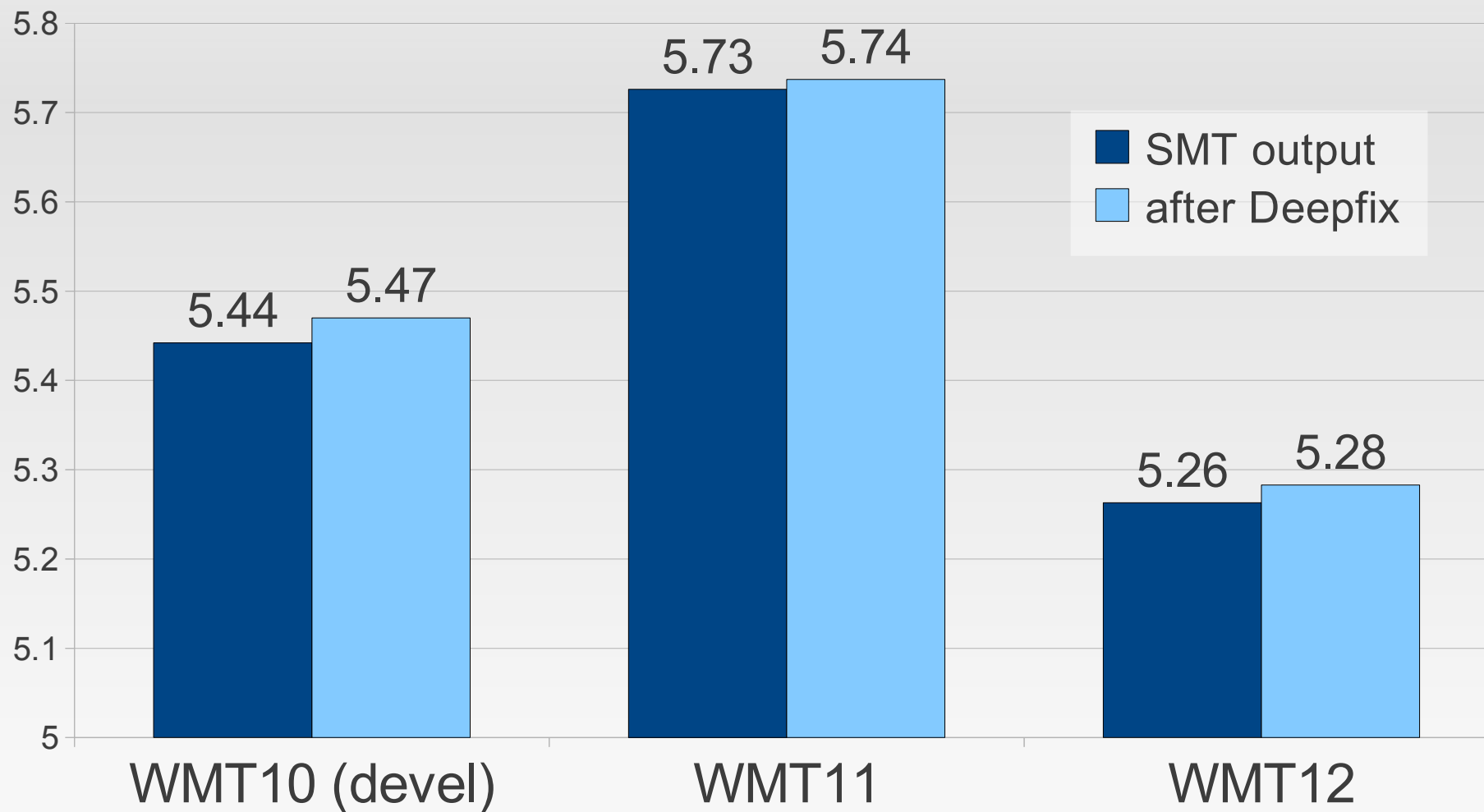
- Agreement enforcement:

  *EU kritizuje nejen* **řeckou**<sub>accusative</sub> **vládu**<sub>accusative</sub>

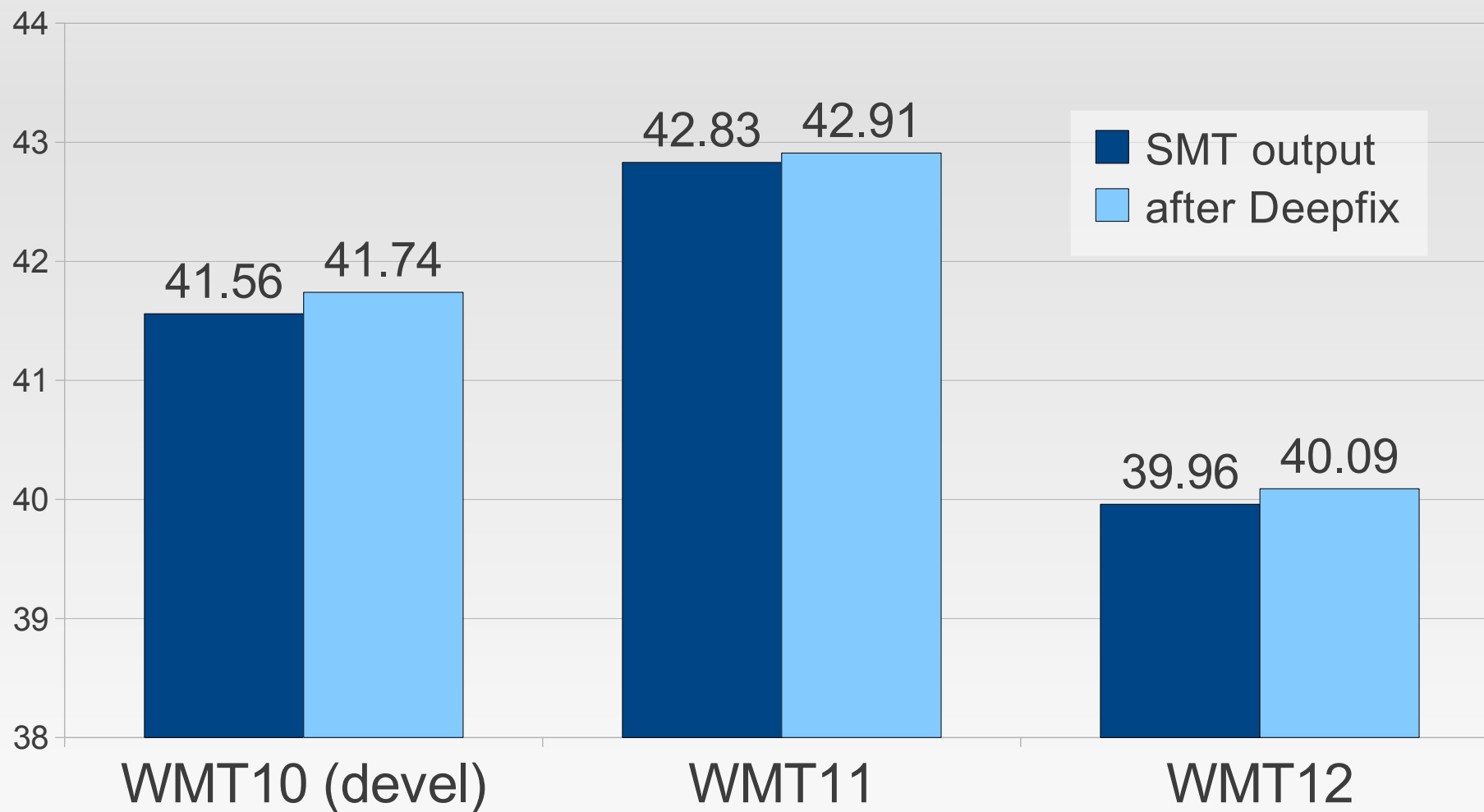  - *EU criticizes not only* **the Greek government**
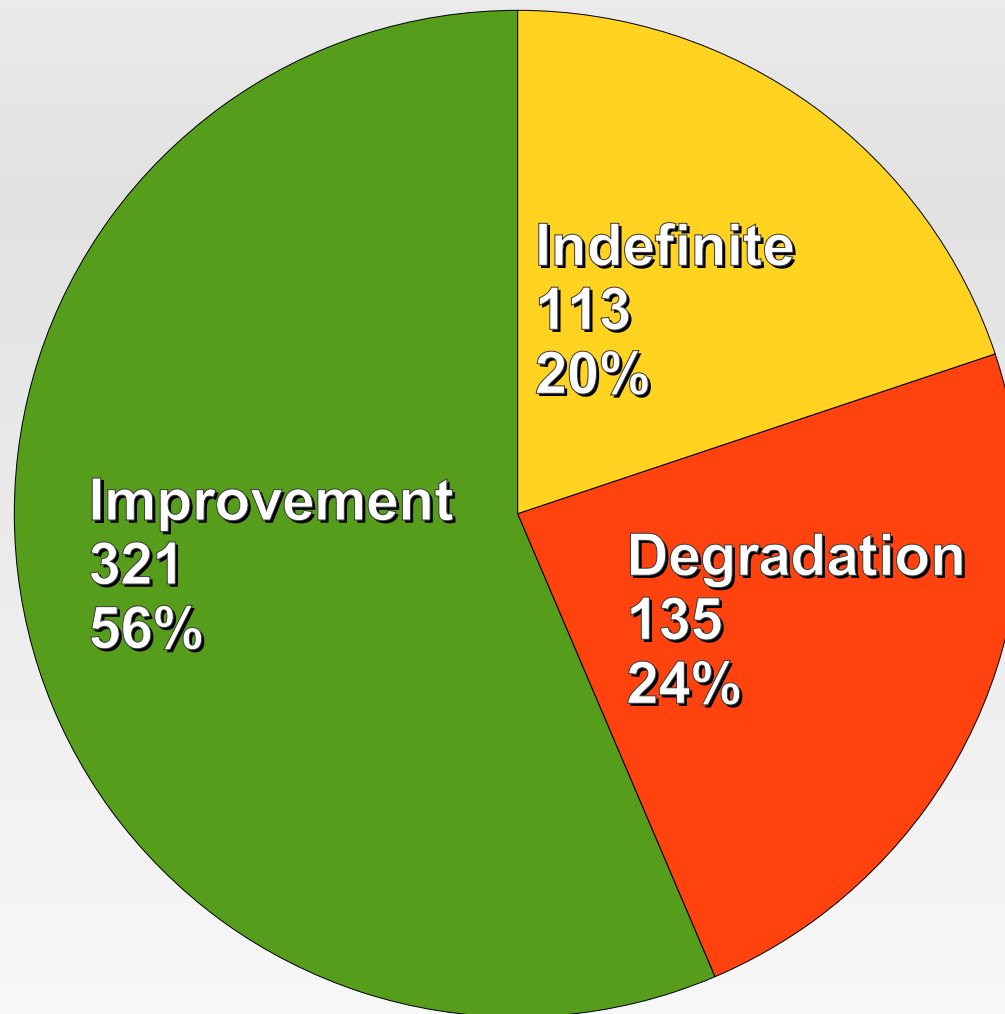
# Automatic evaluation (BLEU)

# Automatic evaluation (NIST)

# Automatic evaluation (1-PER)

# Manual evaluation (changed stcs)

# Conclusion

- address valency errors
  - statistical post-editing of SMT
- identify head-argument pairs (DEEP)
  - deep syntactic analysis
- find the best case for the arguments (FIX)
  - statistical valency model
- obtain slight improvement of translation quality
  - indicated by automatic evaluation
  - confirmed by manual evaluation

# Future work

- explore existing valency lexicons

- more intricate modelling

  - combine more models

  - machine learning (now thresholds semi-manual, and overfitted to development data)

- further adapt underlying NLP tools (tagger)

- extend to other language pairs

# Thank you for your attention

Rudolf Rosa, David Mareček, Aleš Tamchyna
{rosa,marecek,tamchyna}@ufal.mff.cuni.cz

**Deepfix:**
**Statistical Post-editing**
**of Statistical Machine Translation**
**Using Deep Syntactic Analysis**

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

For this presentation and other information, please visit:

http://ufal.mff.cuni.cz/~rosa/