

Evaluating Morphosyntactic Wellformedness of Text

Adithya Pratapa,* Antonios Anastasopoulos,* Shruti Rijhwani,
Aditi Chaudhary, David R. Mortensen, Graham Neubig and Yulia Tsvetkov

EMNLP 2021



Carnegie Mellon University
Language
Technologies
Institute



Evaluating Language Generation

- Rapid progress in language generation in recent years
- A need for an interpretable evaluation of output text quality
- Metrics should be usable across many human languages

We propose a linguistically grounded evaluation of wellformedness in multilingual settings,

L'AMBRE

Grammatical Wellformedness

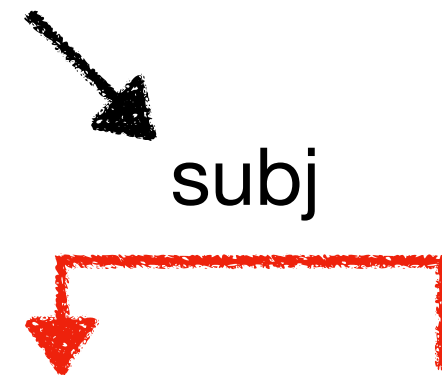
(*) Ich werd**en** lange Bücher lesen

Ich werde lange Bücher lesen

(*) Ich werde lang**en** Bücher lesen

Ich werde lange Bücher lesen

number agreement



PRON

AUX

ADJ

NOUN

VERB

Ich

werd**en**

lange

Bücher

lesen

I-NOM.1SG

will-1**PL**

long-ACC.PL

Book-ACC.PL

read-PTCP

Grammatical Wellformedness

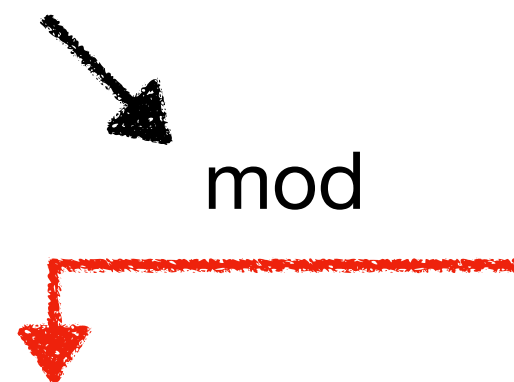
(*) Ich werden lange Bücher lesen

Ich werde lange Bücher lesen

(*) Ich werde langen Bücher lesen

Ich werde lange Bücher lesen

case agreement



PRON AUX ADJ NOUN VERB

Ich werde langen Bücher lesen

I-NOM.1SG will-1SG long-DAT.PL Book-ACC.PL read-PTCP

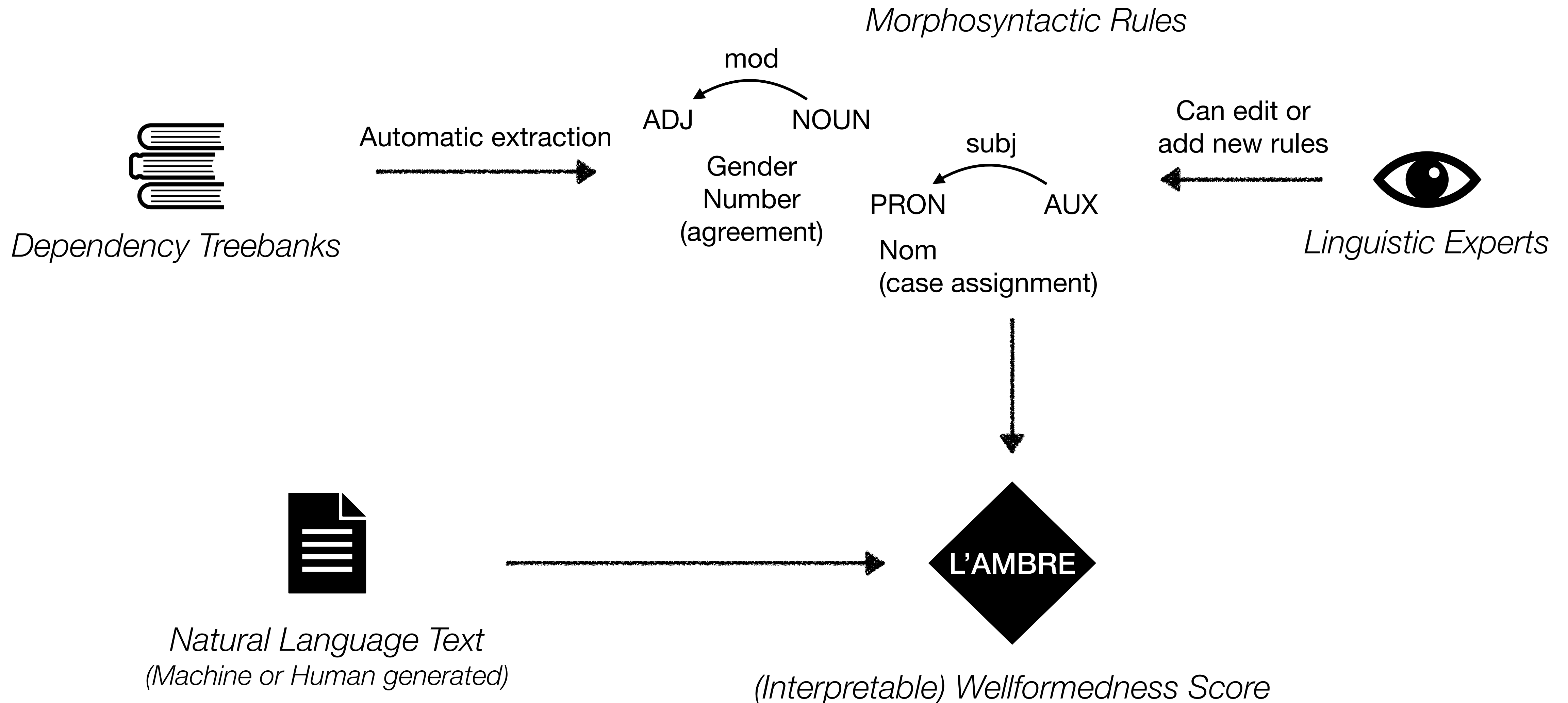
Wellformedness \subset Fluency

- Wellformedness is usually necessary for fluency but not sufficient (Sakaguchi et al., 2016)
- We focus on wellformedness for two reasons,
 1. Easier to measure via a simple comparison against a checklist
 2. Wellformedness rules can be extracted automatically and accurately

New Metric: L'AMBRE

- Linguistically **A**ware **M**orphosyntax-**B**ased **R**ule **E**valuation
- Evaluates natural language text by checking against language-specific morphosyntactic rules
- Metric is interpretable by nature
- Morphosyntactic rules are automatically extracted from dependency treebanks, can be easily extended to new languages

L'AMBRE Overview



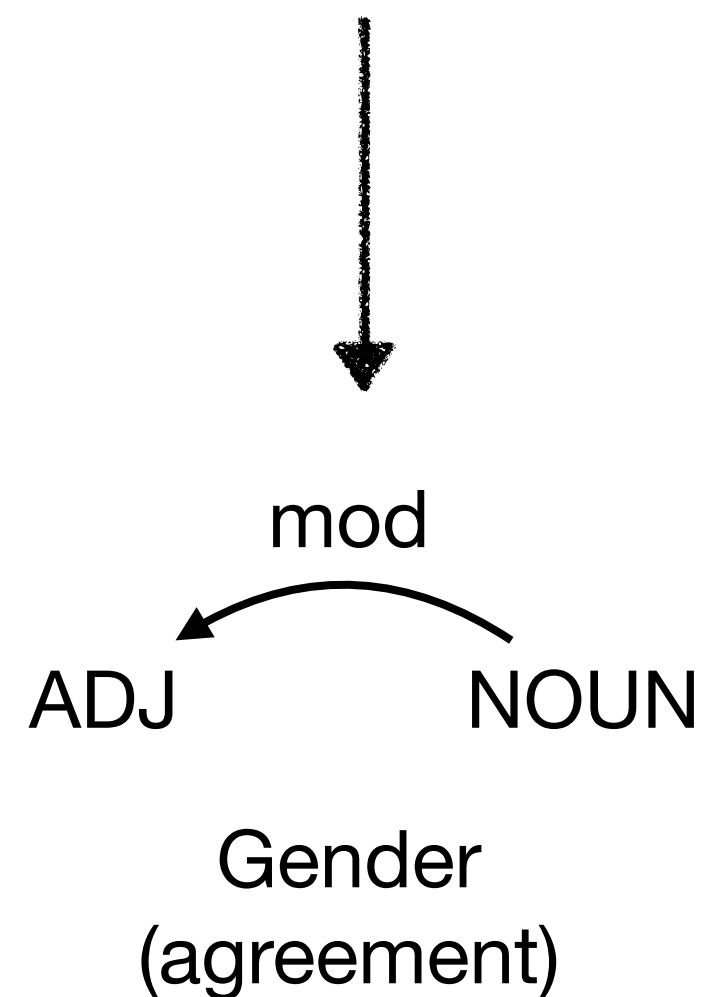
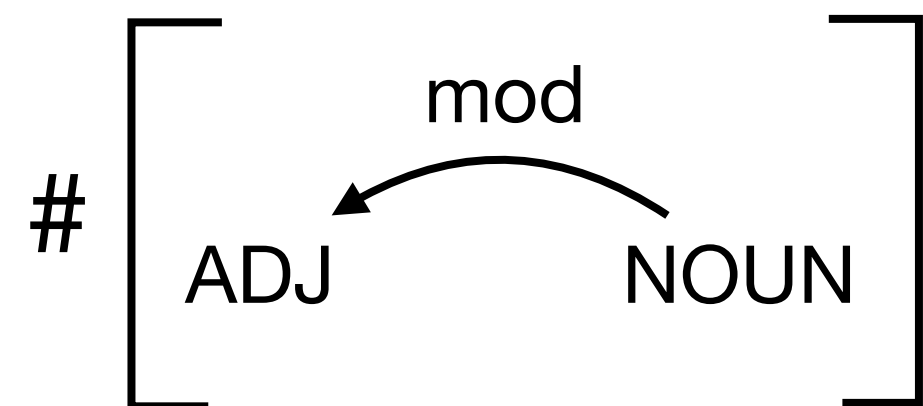
Morphosyntactic Rules

- We focus on three types of rules,
 - Agreement
 - Case assignment
 - Verb form choice
- Extract directly from SUD treebanks (Chaudhary et al. 2020)

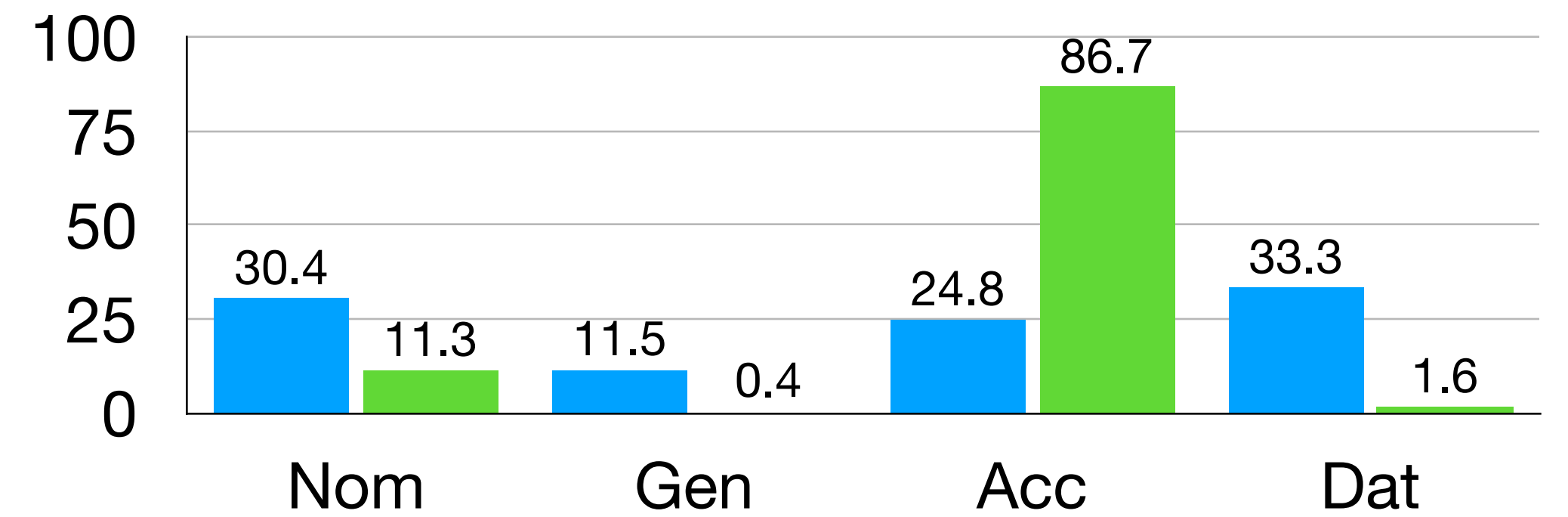
Extracting Rules

Agreement

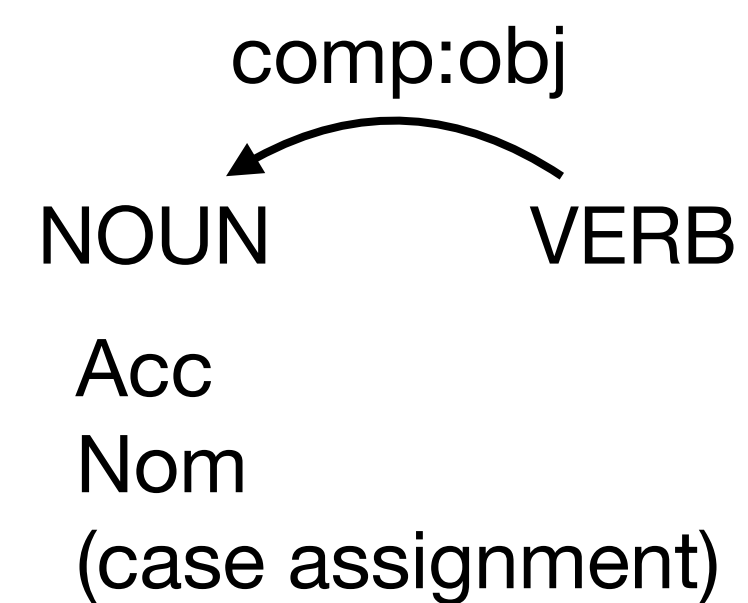
[Gender_{ADJ} = Gender_{NOUN}] ≥ 0.9



Case Assignment



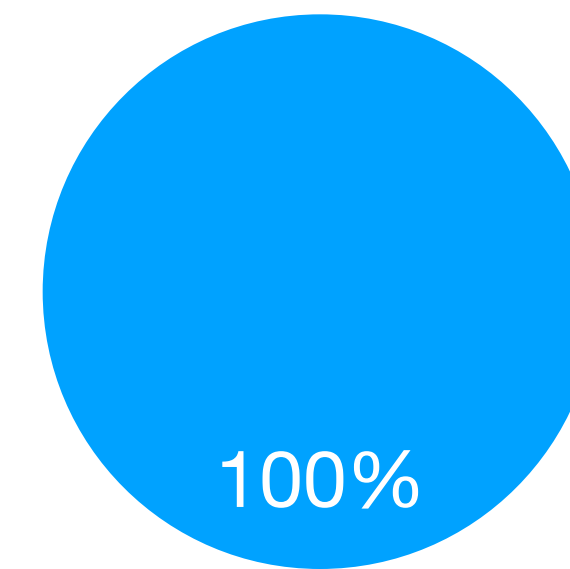
■ Global (NOUN)
■ Local (comp:obj-NOUN-VERB)



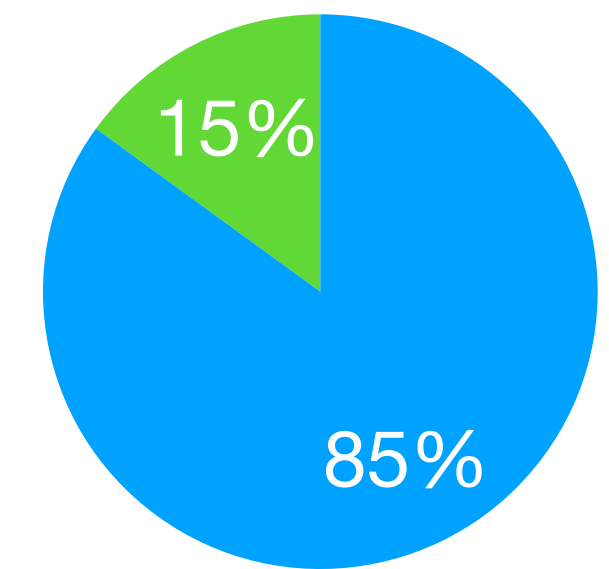
Human Evaluation of Rules

- We ask human experts to label each rule as one of
 - almost always true
 - sometimes true
 - need not be true
- Our rules are of high quality

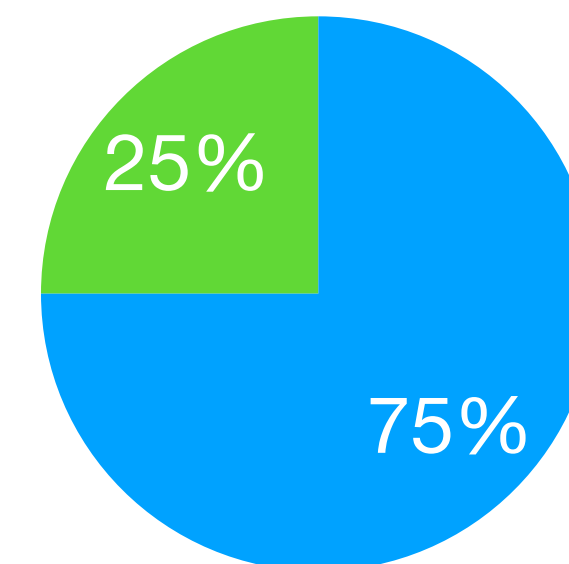
● Almost Always True
● Sometimes True
● Need not be true



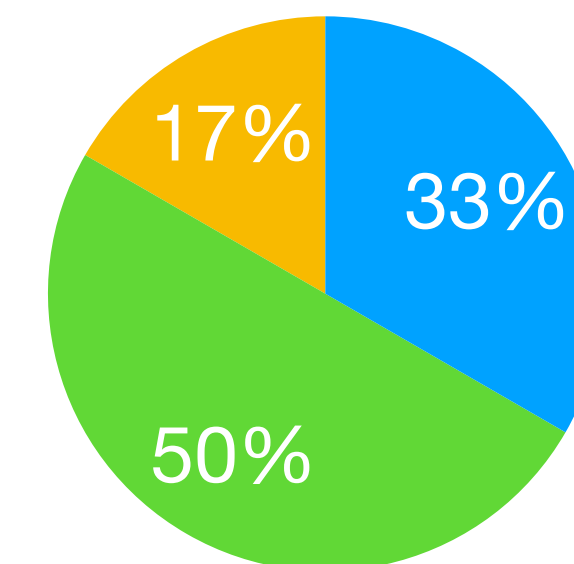
Greek Agreement



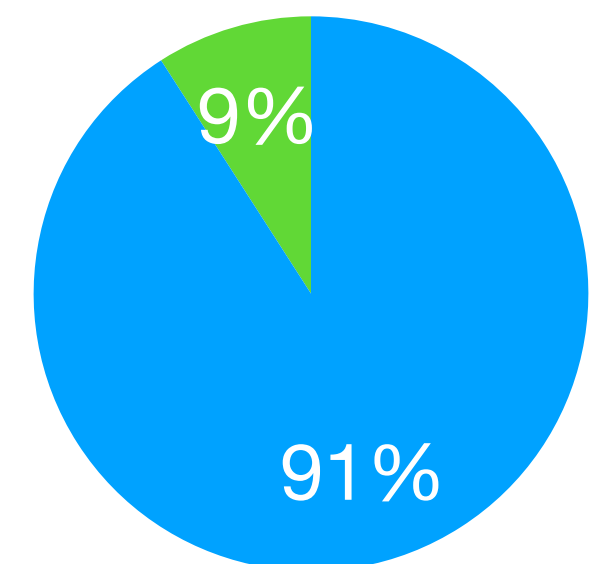
Russian Agreement



Greek Argument Structure



Russian Argument Structure



Italian Argument Structure

Checking Rules

- **L'AMBRE**: parse text and match against the morphosyntactic rules
- But...
 - Noisy (or grammatical incorrect) text can throw off the standard parsers (Hashemi and Hwa, 2016 and others)

Making Parsers Robust

- Utilize UniMorph to add morphology related noise to standard treebanks
- Retrain Stanza parsers on noisy treebanks
- Improved results in feature tagging and dependency parsing on noisy text

Treebank Sentence

(οικισμός + N;ACC;SG)

Στο μικρό **οικισμό** της Λίνδου

In-the small settlement of Lindos



Στο μικρό **οικισμούς** της Λίνδου

(οικισμός + N;ACC;PL)

Noisy Sentence

Validating L'AMBRE: Grammatical Error Correction (GEC)

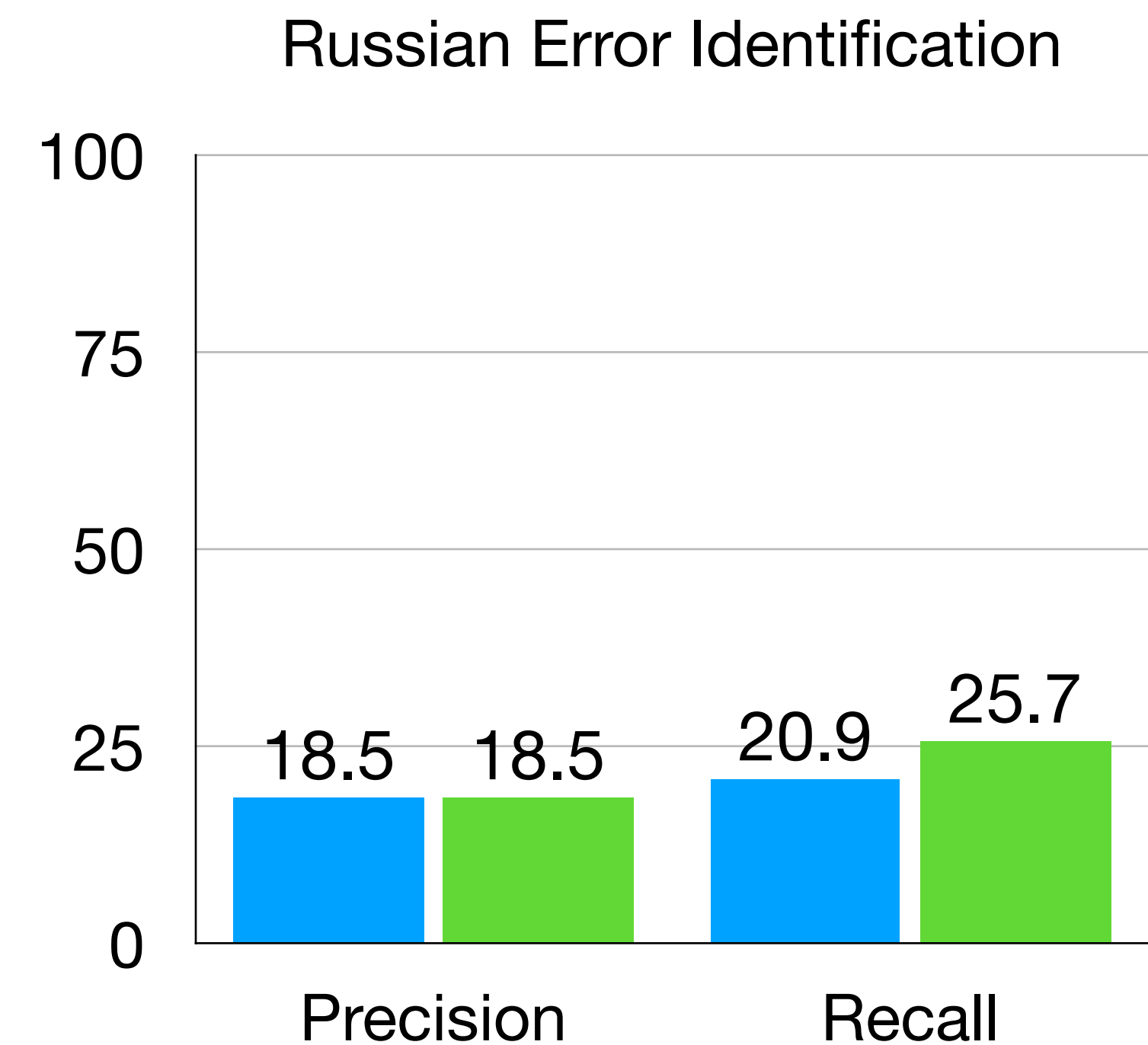
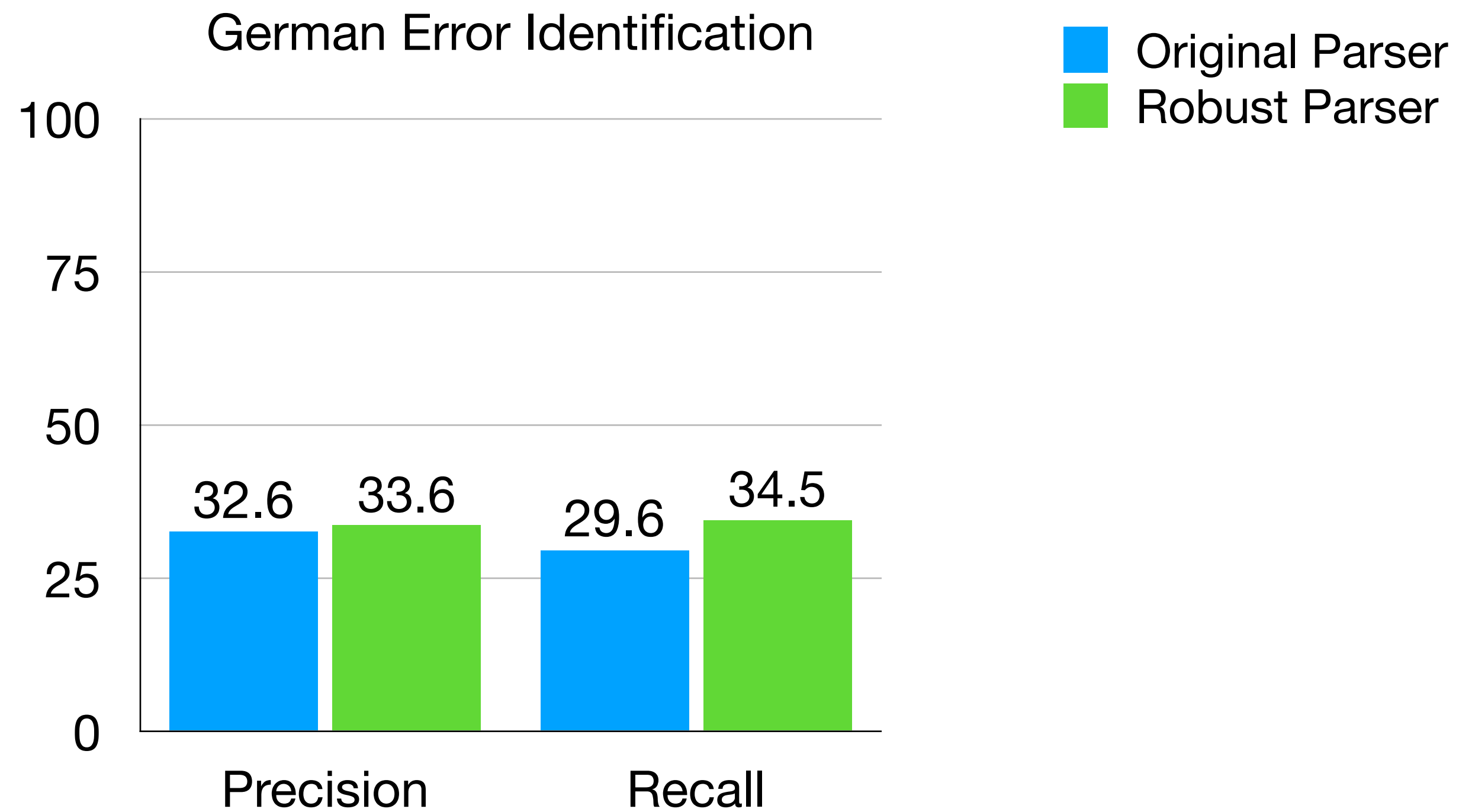
- Does L'AMBRE accurately capture grammaticality?
- Validate on Russian and German GEC
- We focus on error identification (not correction)

			ihr	aktuelles	
Vielleicht	verdienen	Ärzte	ihre	aktuelle	Gehalt
<i>Maybe</i>	<i>to earn</i>	<i>doctors</i>	<i>their</i>	<i>current</i>	<i>salary</i>

Example from German GEC dataset

Results

- Precision and Recall of detecting morphology related errors
- Robust parsers improve on recall



Results

- Comparison to other metrics for error identification
 - grammaticality-based metric (Napoles et al. 2016; GBM)
 - LM perplexity
- Observations
 - Perplexity performs the best but not interpretable
 - L'AMBRE does well at morphology related errors, and GBM does well with other fluency related rules

Courtney Napoles, Keisuke Sakaguchi, and Joel Tetreault. *There's no comparison: Referenceless evaluation metrics in grammatical error correction*. EMNLP 2016.

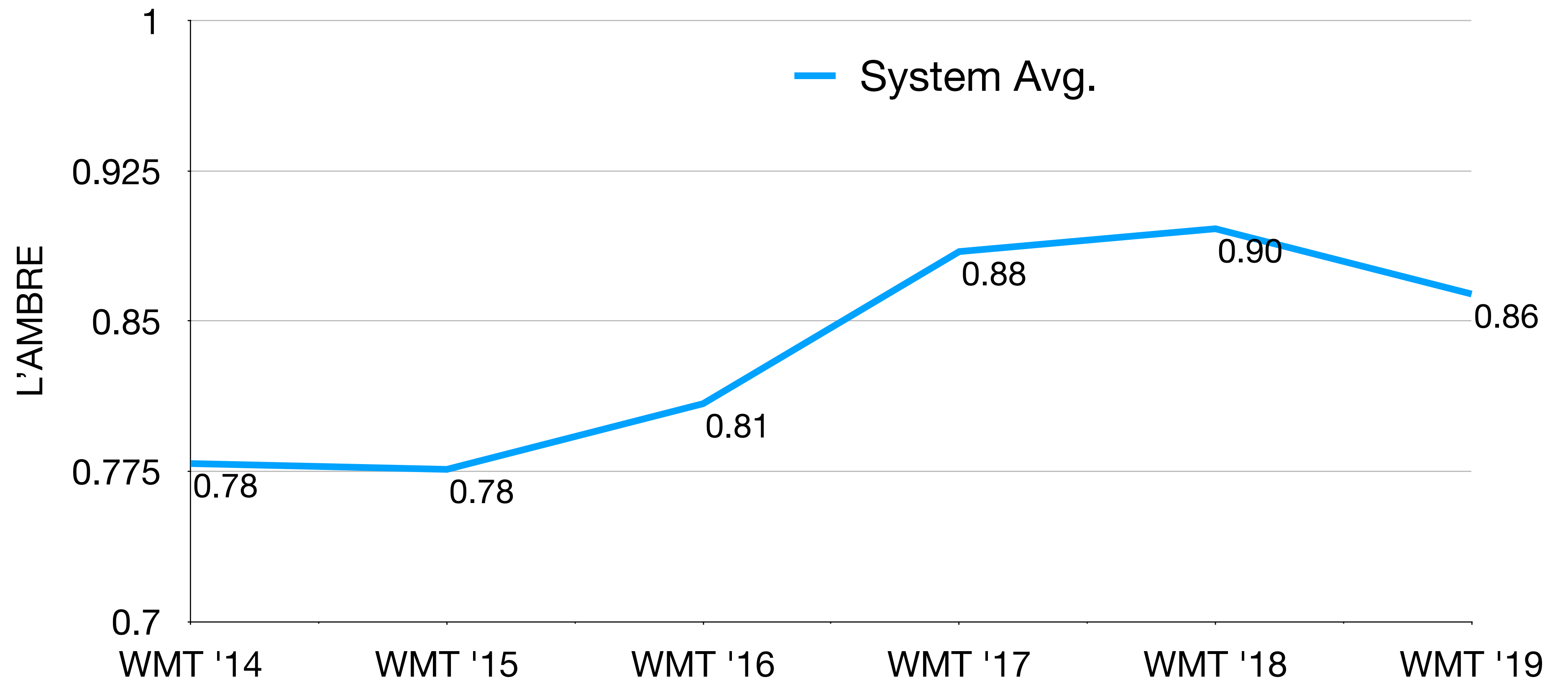
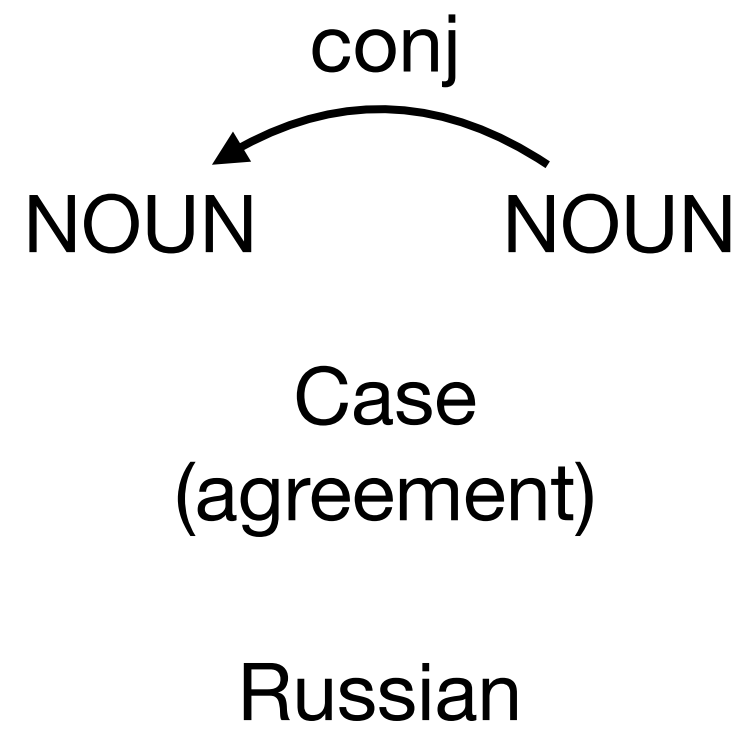
NMT Case Study

English → X	WMT '18	WMT '19
Czech	0.84	0.80
German	-0.06	0.16
Estonian	0.68	-
Finnish	0.86	0.85
Russian	0.86	0.57
Turkish	0.58	-

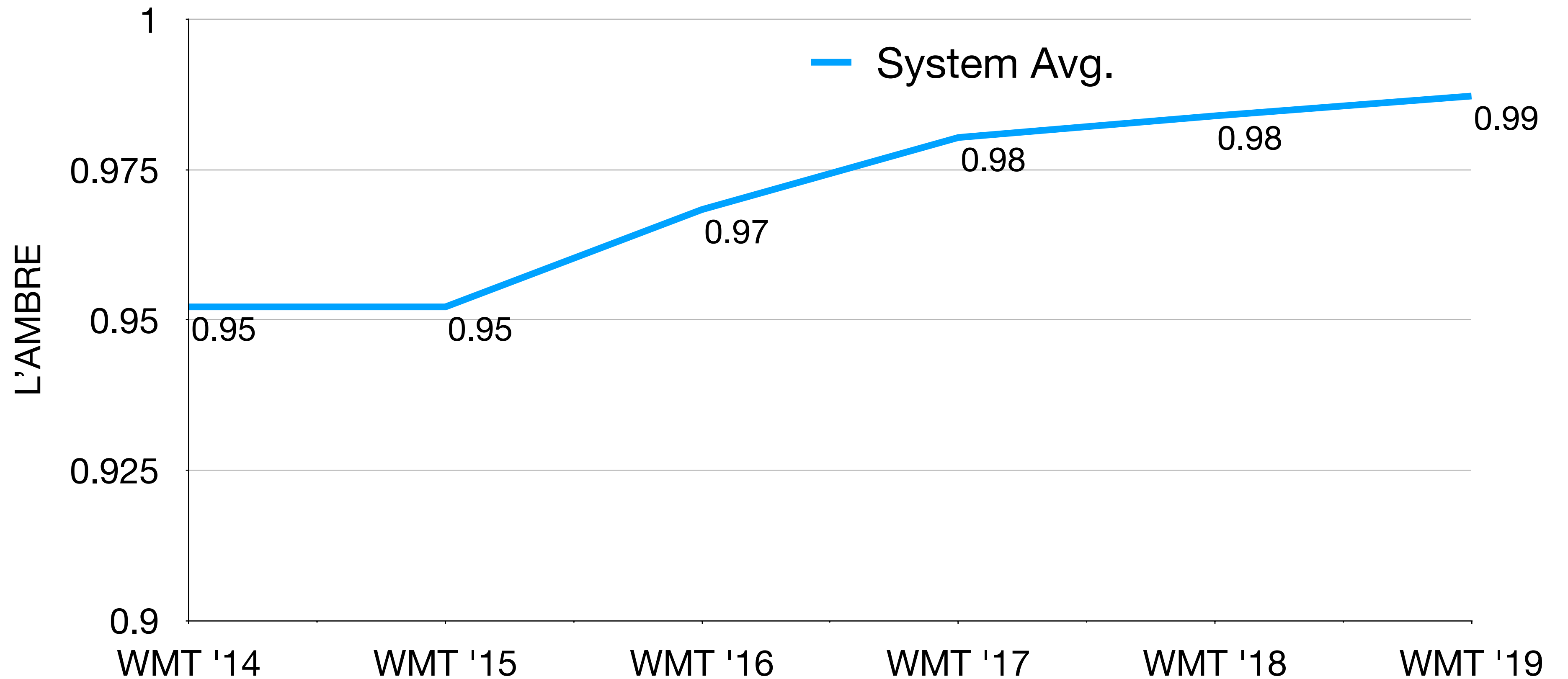
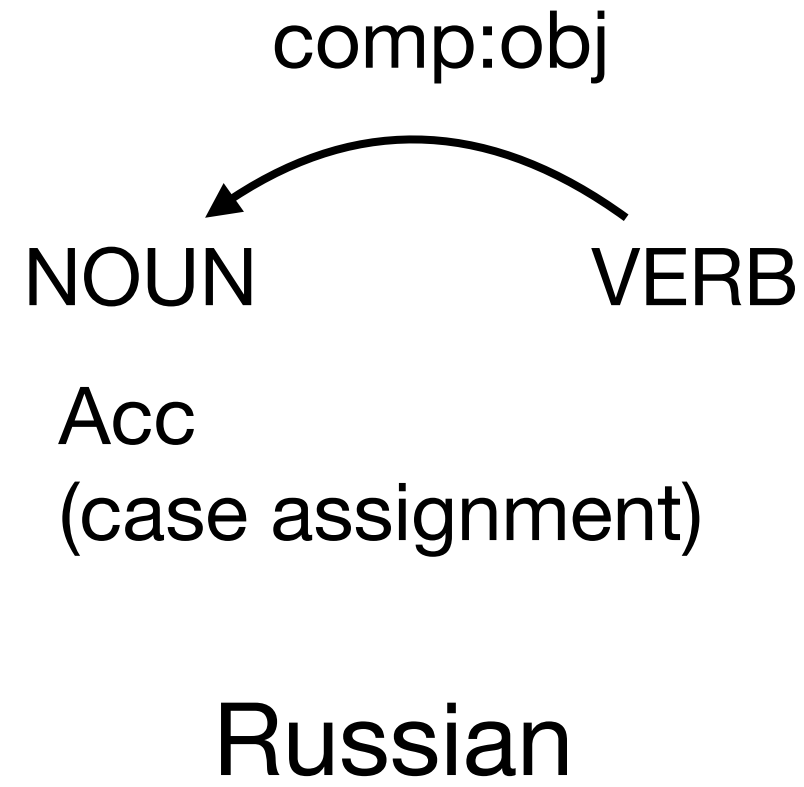
Pearson's r Correlations to human z-scores

- Analyze systems from WMT shared tasks
- L'AMBRE correlates well with human scores
- L'AMBRE also correlates well with perplexity
- L'AMBRE can provide fine-grained analysis

Fine-grained Diachronic Analysis

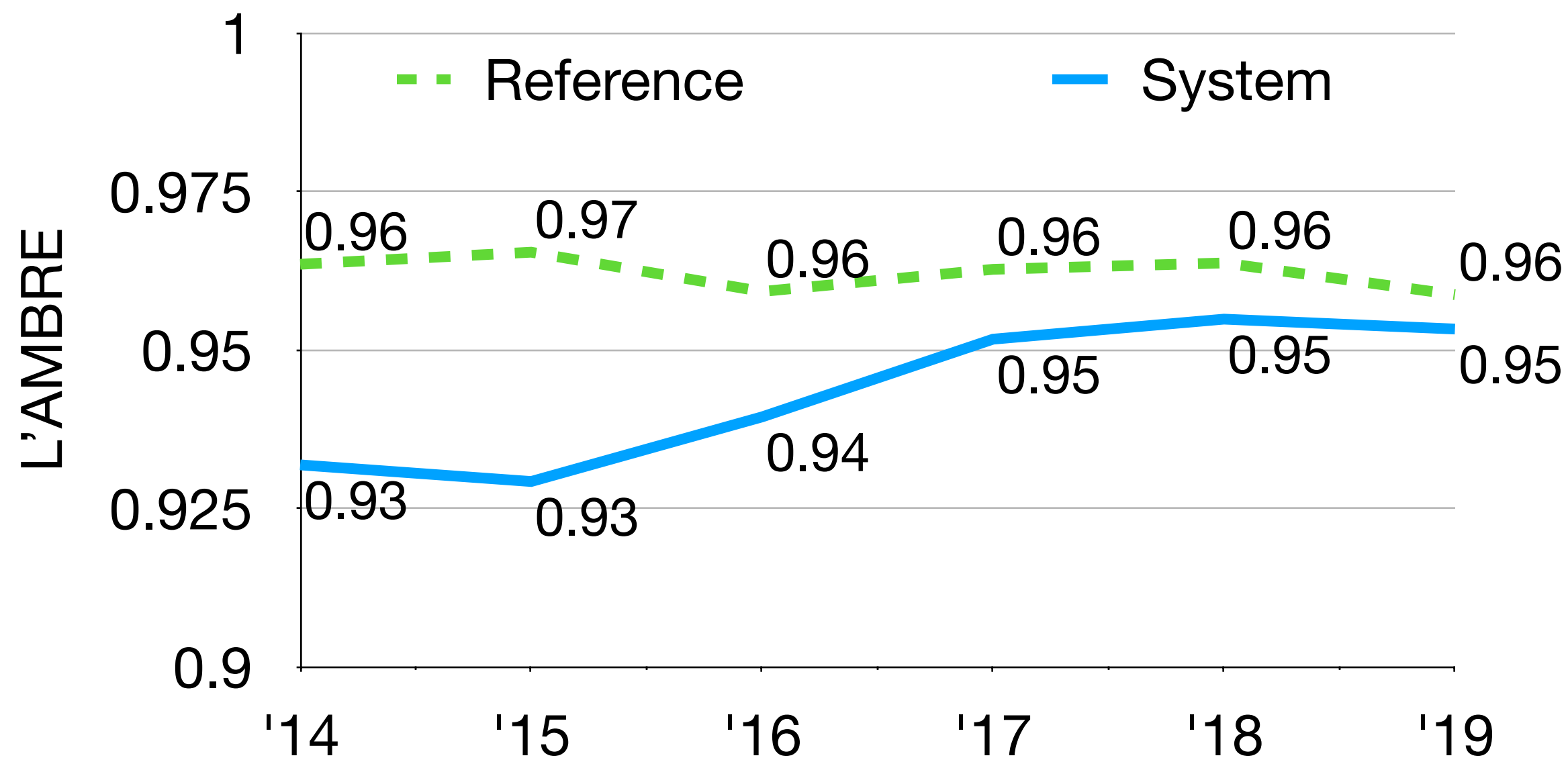


Fine-grained Diachronic Analysis

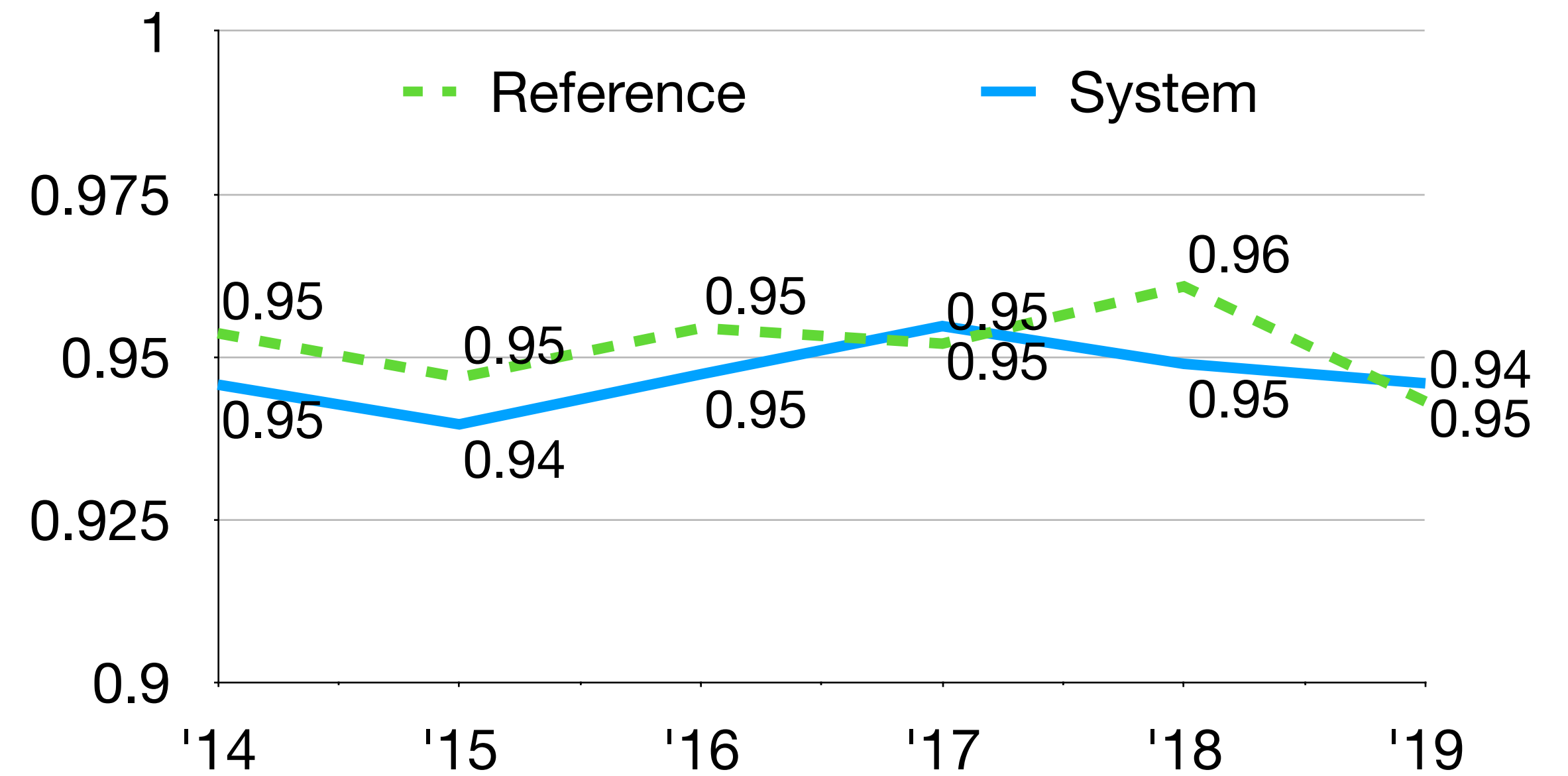


Diachronic Analysis

Russian WMT



German WMT



Summary

- Our contributions,
 - **L'AMBRE**: a multilingual metric to measure wellformedness
 - **Robust parsers**: make parsers and taggers robust by adding morphological noise
 - **Case Study on NMT**: applications in evaluating WMT system outputs

github.com/adithya7/lambre

