



COORDINATION STRUCTURES IN DEPENDENCY TREEBANKS

MARTIN POPEL, DAVID MAREČEK, JAN ŠTĚPÁNEK, DANIEL ZEMAN, ZDENĚK ŽABOKRTSKÝ

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



INTRODUCTION

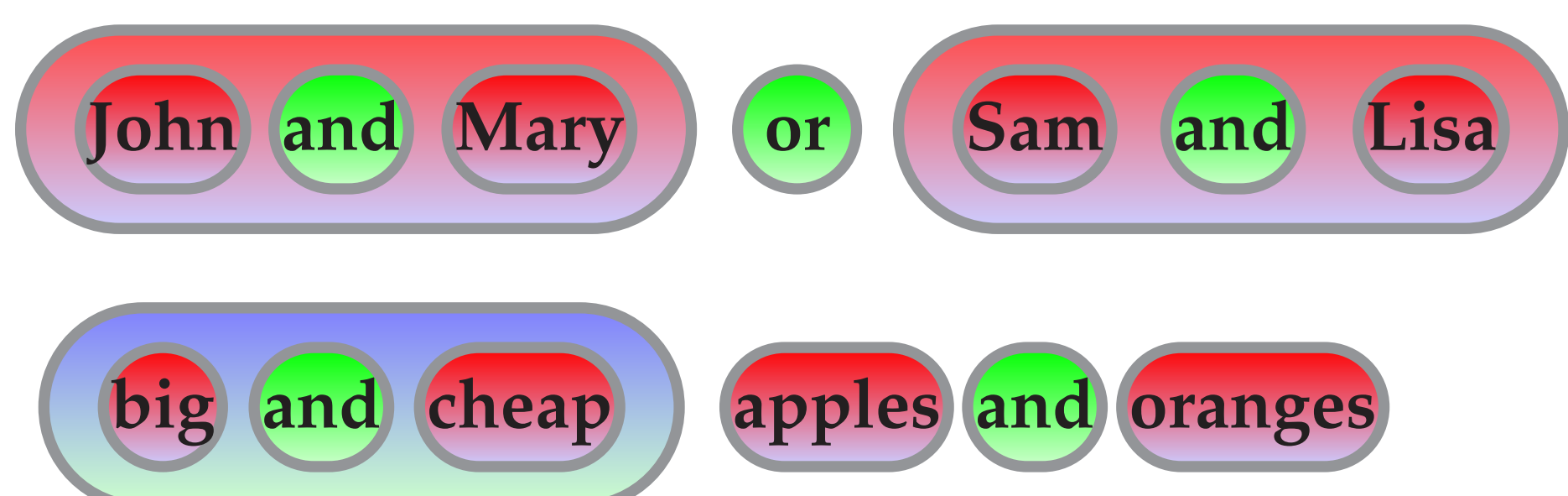
Coordination structures (CS)

are difficult to represent in dependency treebanks:

- coordination vs. dependency are fundamentally different relations
- nested coordination
- shared vs. private modifiers
- multiconjunct CS, punctuation, etc.

Examples:

notation: conjunct conjunction shared modifier



Problem

- large inter-treebank variation
- obstacle for multi-lingual parsing

Our Goal

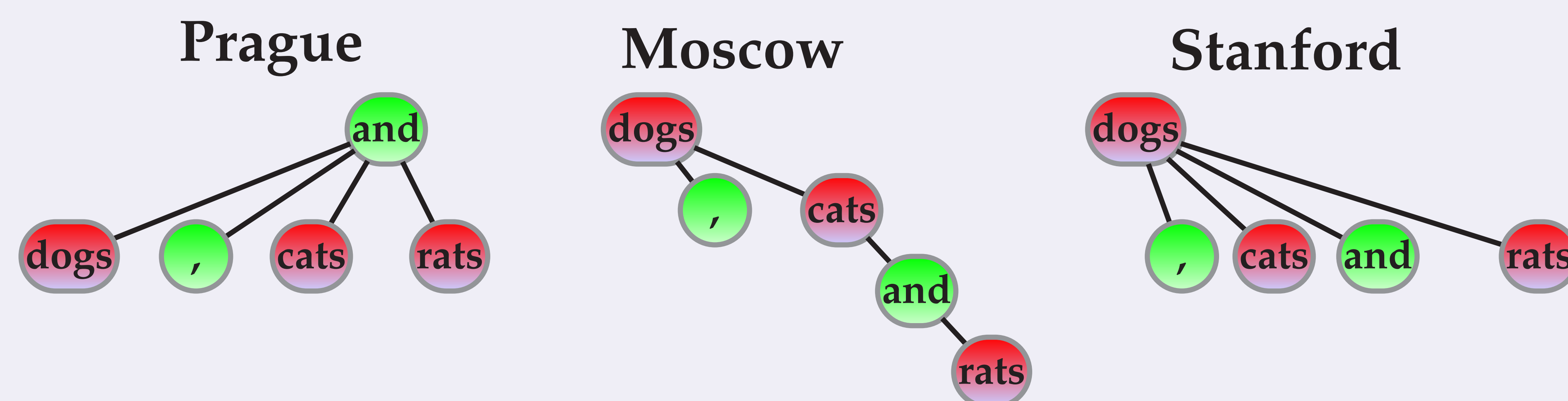
- explore the CS variations in a systematic way
- convert the treebanks into a common CS style

NOVEL TAXONOMY OF CS STYLES

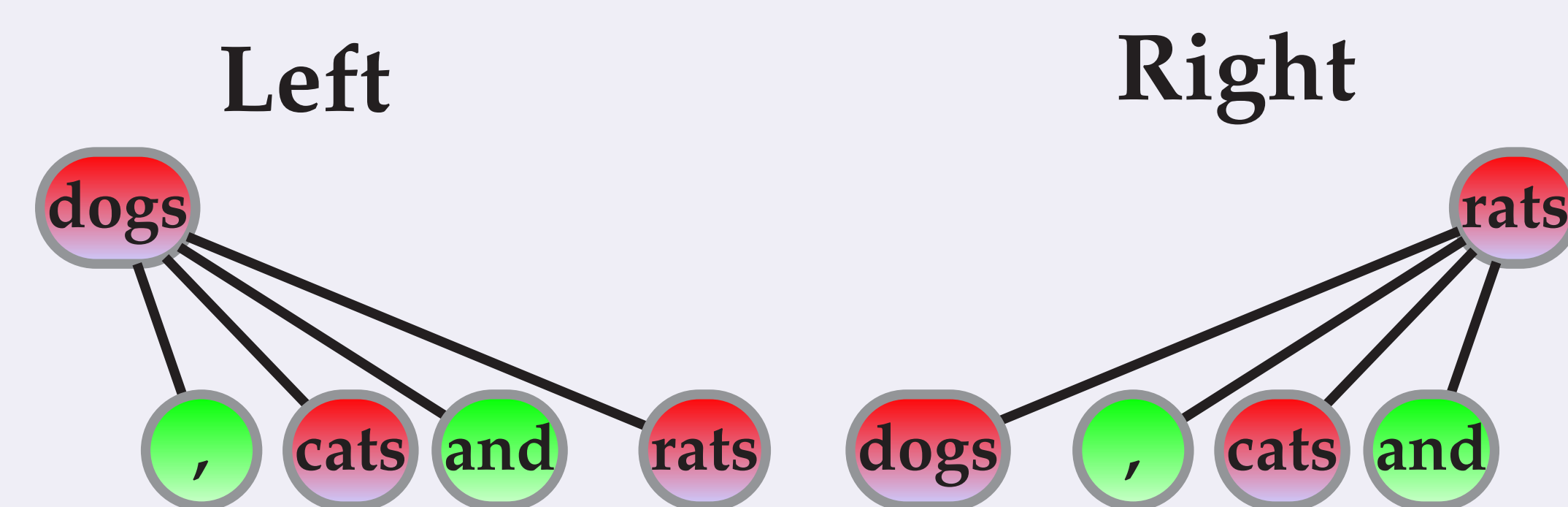
We identified

- 5 dimensions in CS tree shape variations
- 3 dimensions in CS labeling
- a few additional subtle variations
- in theory over one thousand possible styles
- 16 styles found in the real treebanks

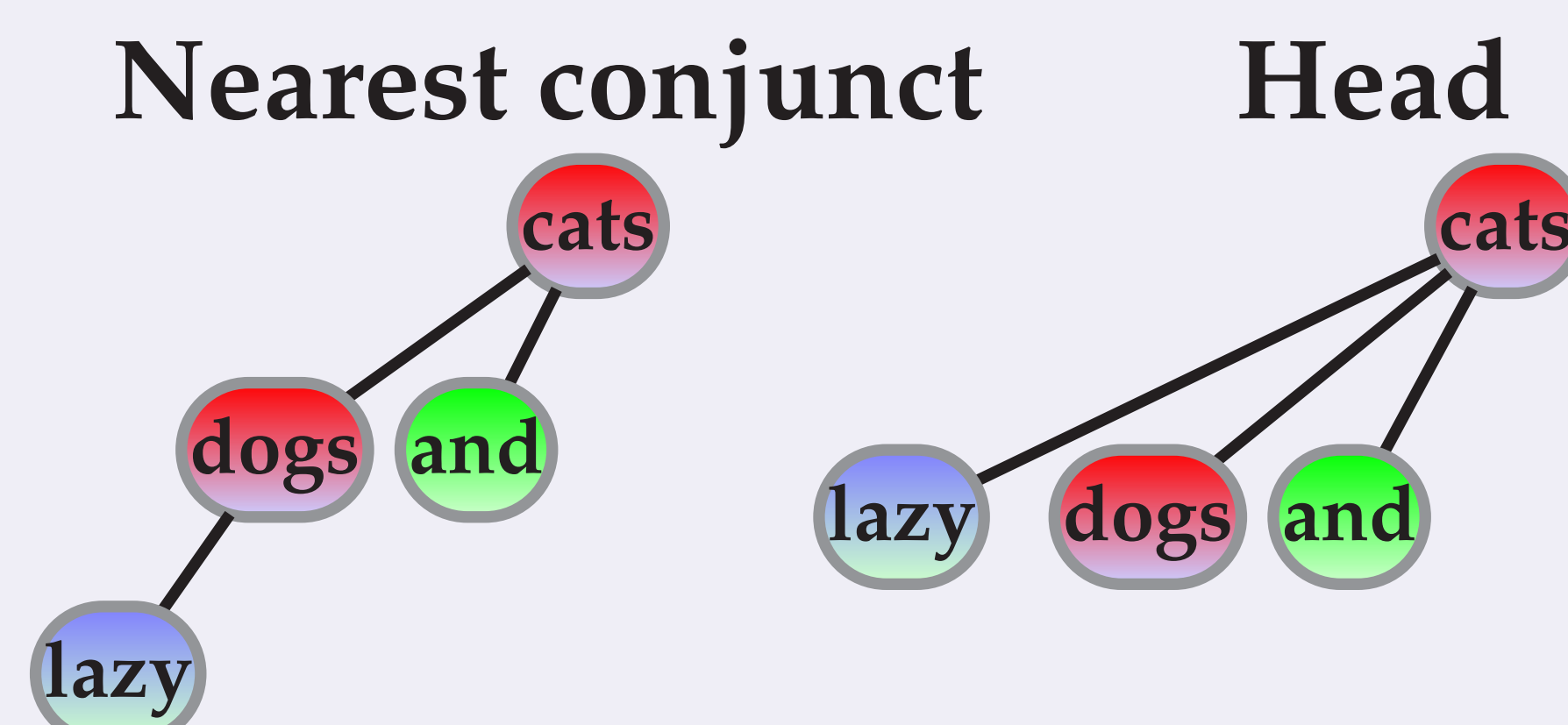
Family – configuration of conjuncts



Choice of head



Shared modifiers below



CONVERTIBILITY

Different CS styles do not have equivalent expressive power ⇒ no chance for a lossless conversion.

- We developed an algorithm that decomposes a CS in one style and assembles it in another style.
- Empirical roundtrip accuracy: usually > 99%

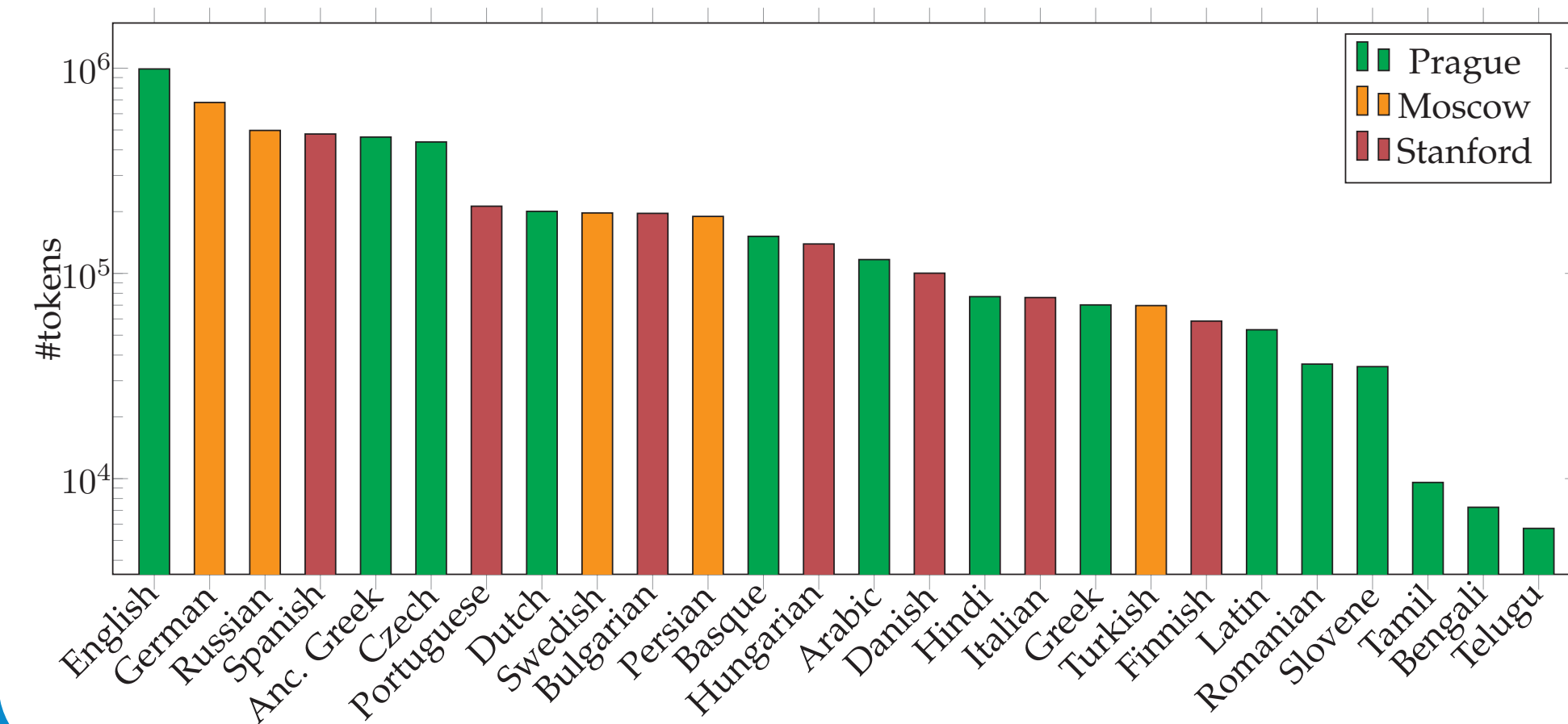
Roundtrip means e.g. Prague → Moscow → Prague evaluated by unlabeled attachment score.

CONCLUSIONS

- a survey of coordination styles in 26 treebanks
- a general taxonomy which covers most of the variations
- 26 treebanks converted into a common style available at <http://ufal.ms.mff.cuni.cz/hamledt/>
- relatively high convertibility accuracy should allow future experiments with learnability of CS by parsers

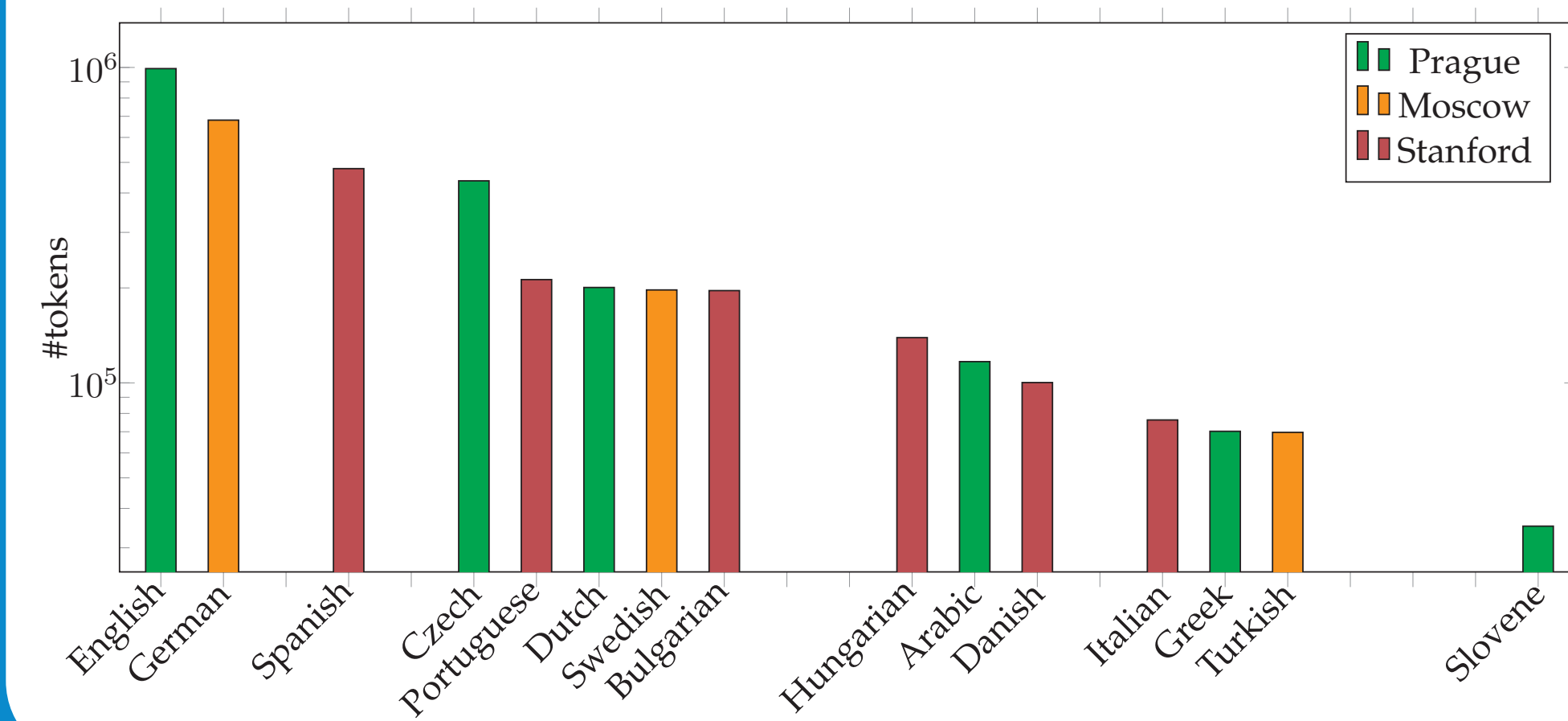
ANALYZED TREEBANKS

26 treebanks from HAMLEDT



SIMILAR COLLECTIONS OF TREEBANKS

CoNLL 2006–2009



GOOGLE Universal Treebanks v1.0

