

Morphological Processing of Low-Resource Languages:

Where We Are and What's Next



Adam Wiemerslage*, Miikka Silfverberg*, Changbing Yang,
Arya D McCarthy, Garrett Nicolai, Eliana Colunga, and Katharina Kann



For a survey of recent work in computational morphology for low-resource languages, please refer to our paper!

Task: *tUMPC*

Questions

Predict the full paradigm for a given word in context

“My best friend broke our lamp” “Some geese are flying over my head”

break	<u>broke</u>	breaking
breaks	broken	

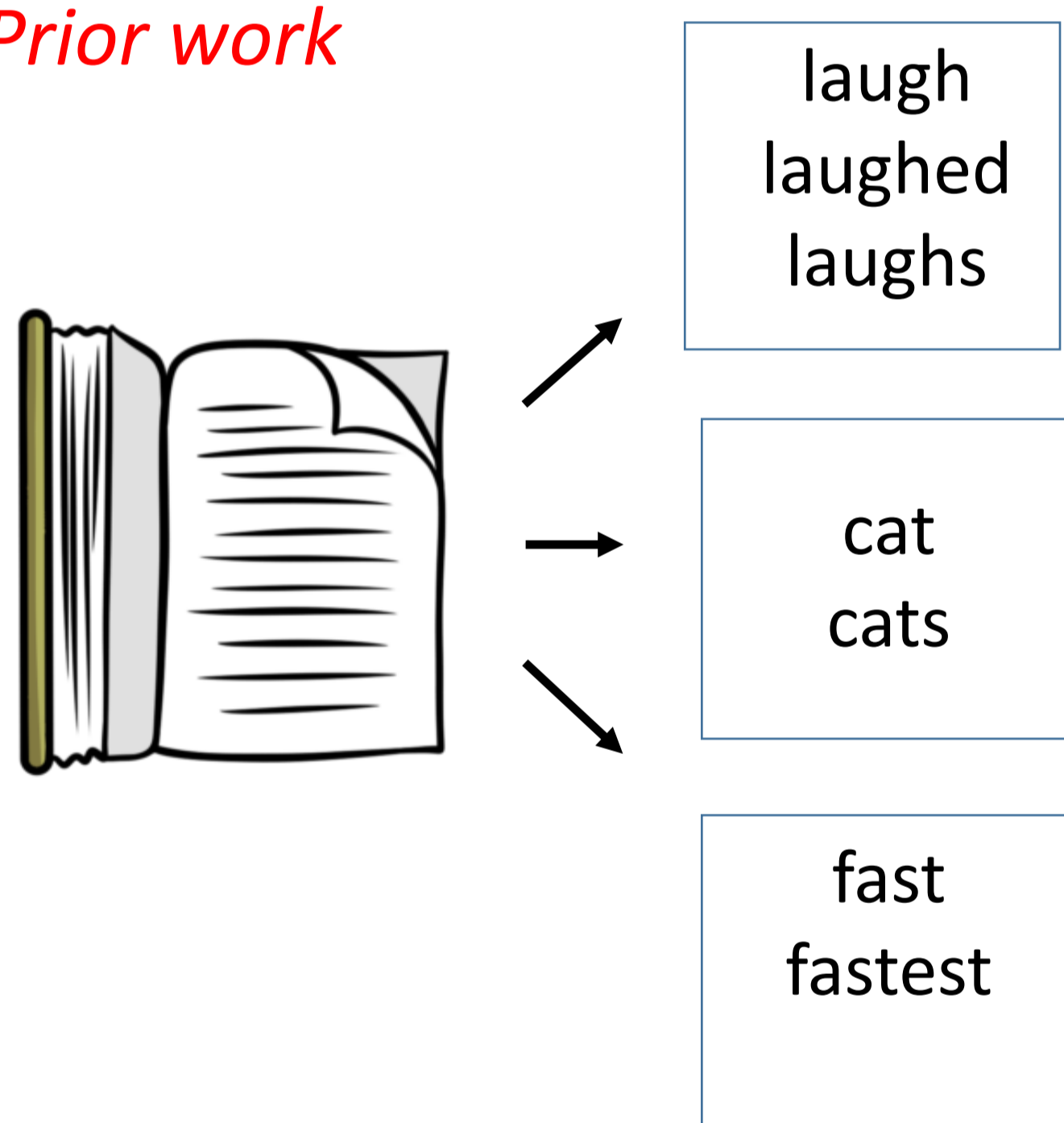
goose	<u>geese</u>
-------	--------------

1. How to handle *all* words (and POS) in the corpus when learning?
2. Which existing systems work best in a pipeline?
3. What kind of corpus is easiest to learn from?

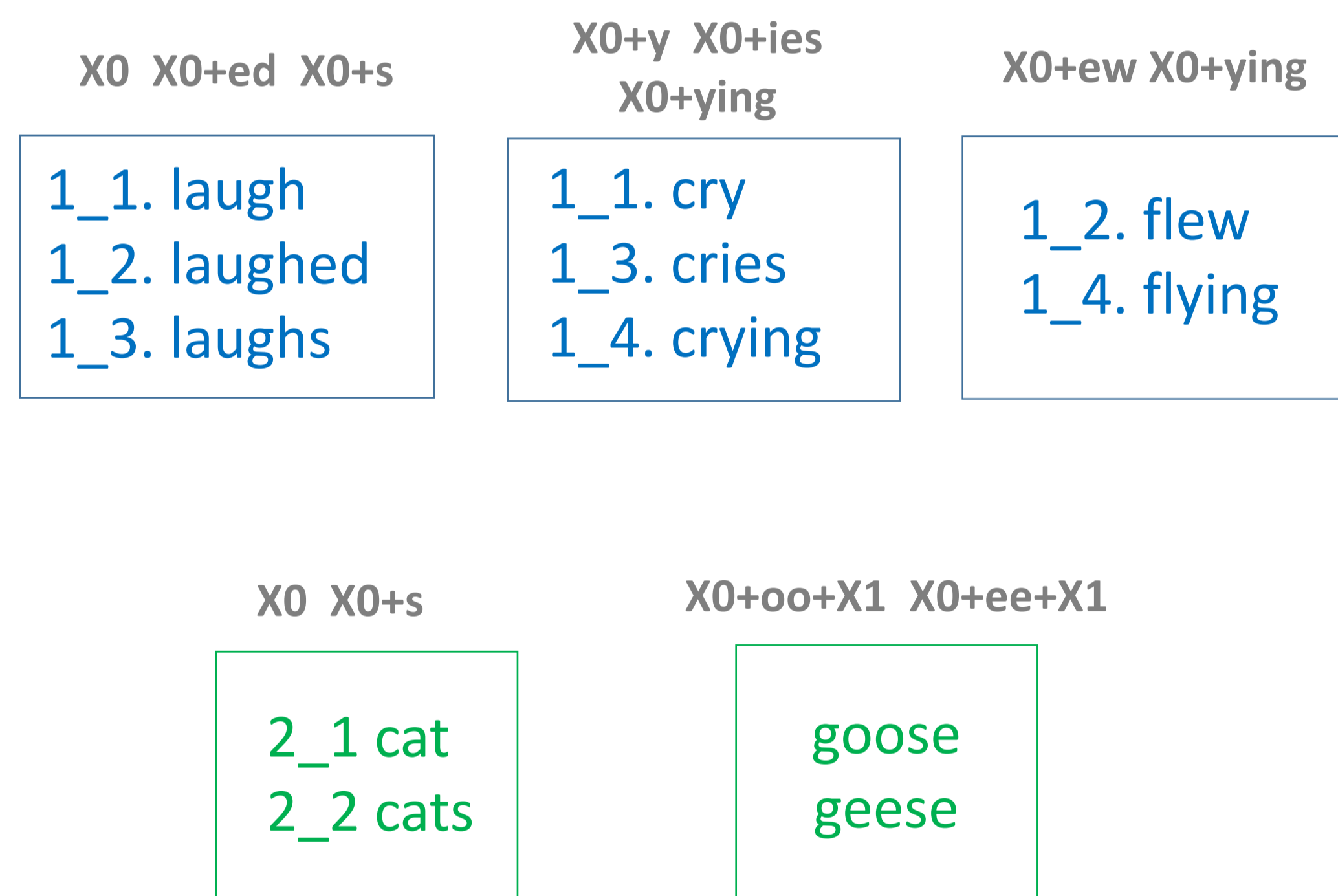
Learning

1. Cluster (partial) paradigms

Prior work

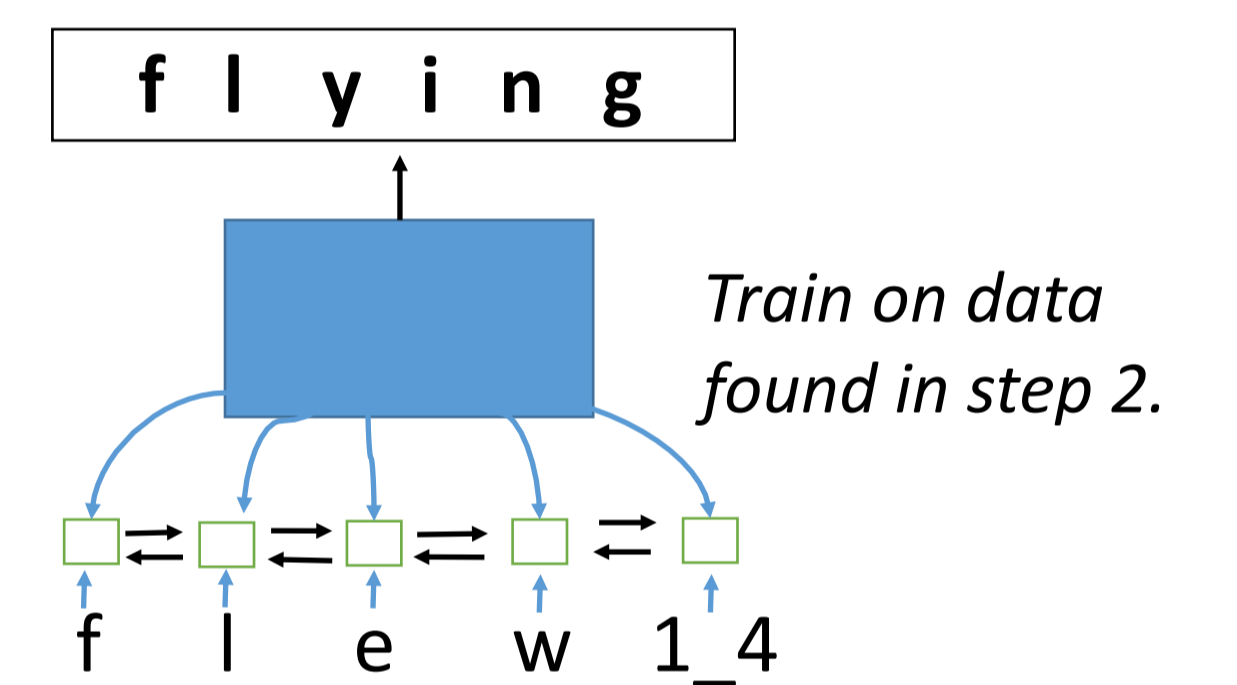


2. Align slots (and cluster POS)



3. Train inflector

Prior work



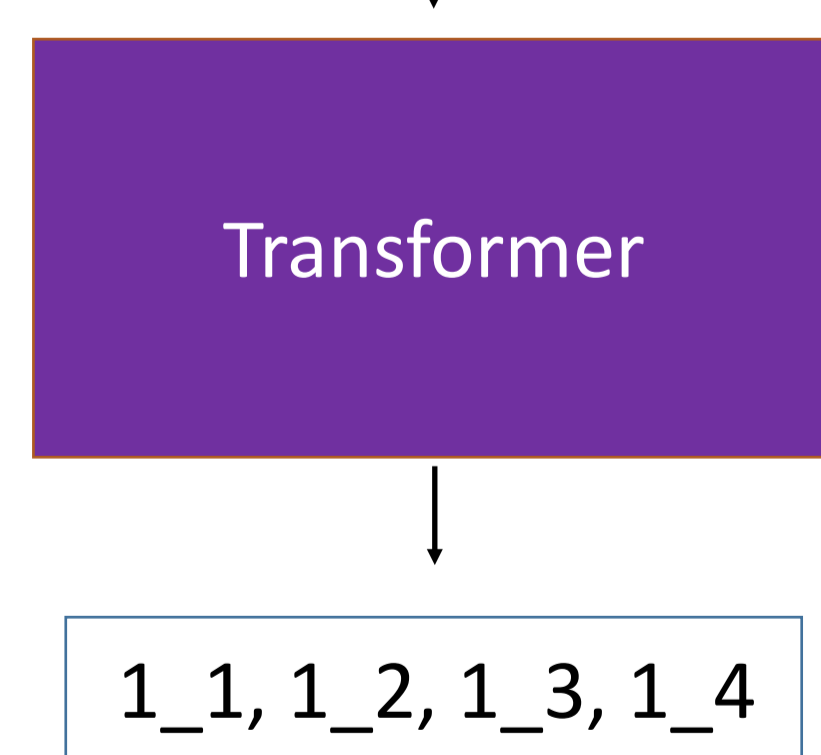
4. Train tagger



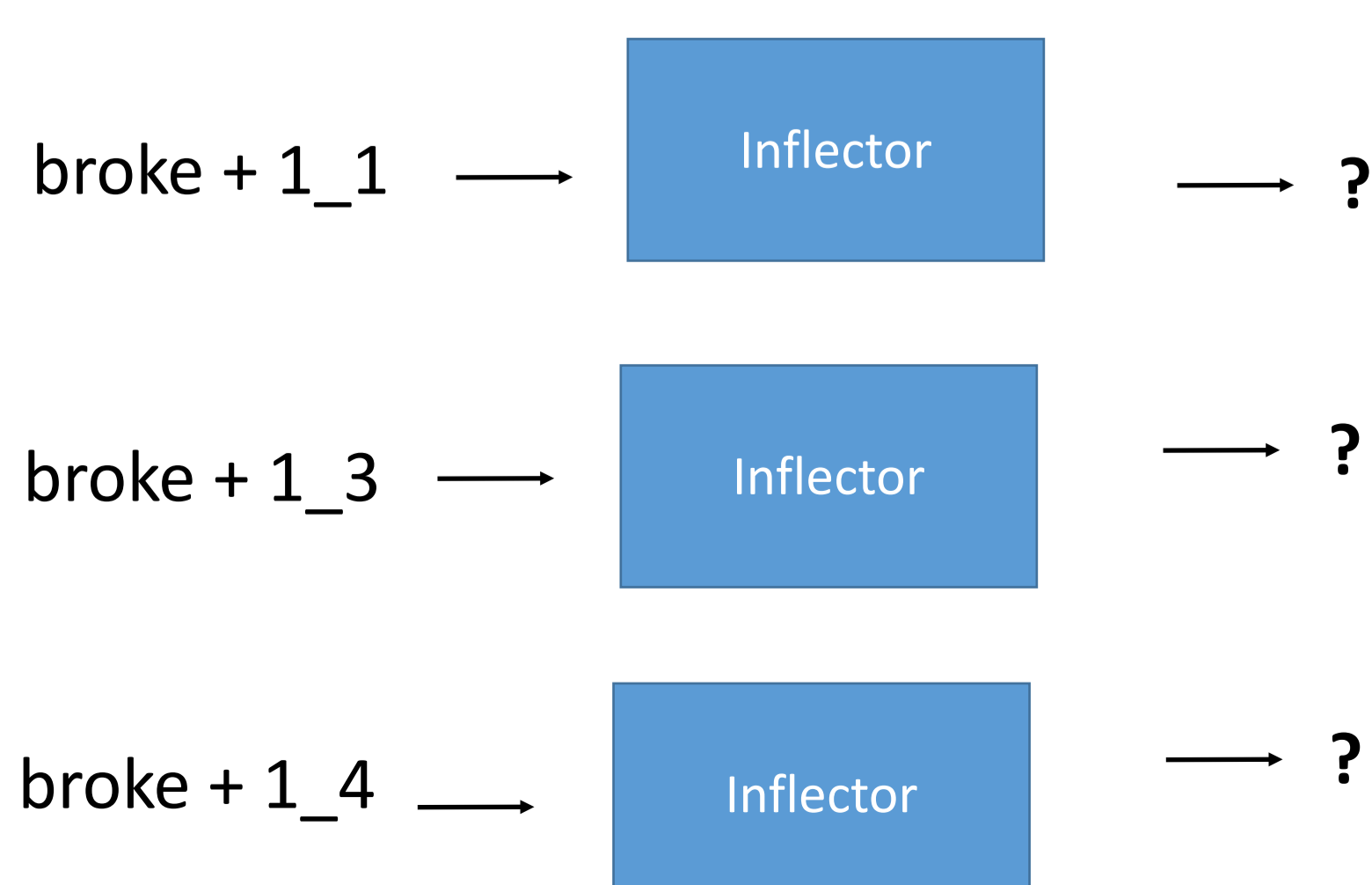
Inference

1. Predict slots

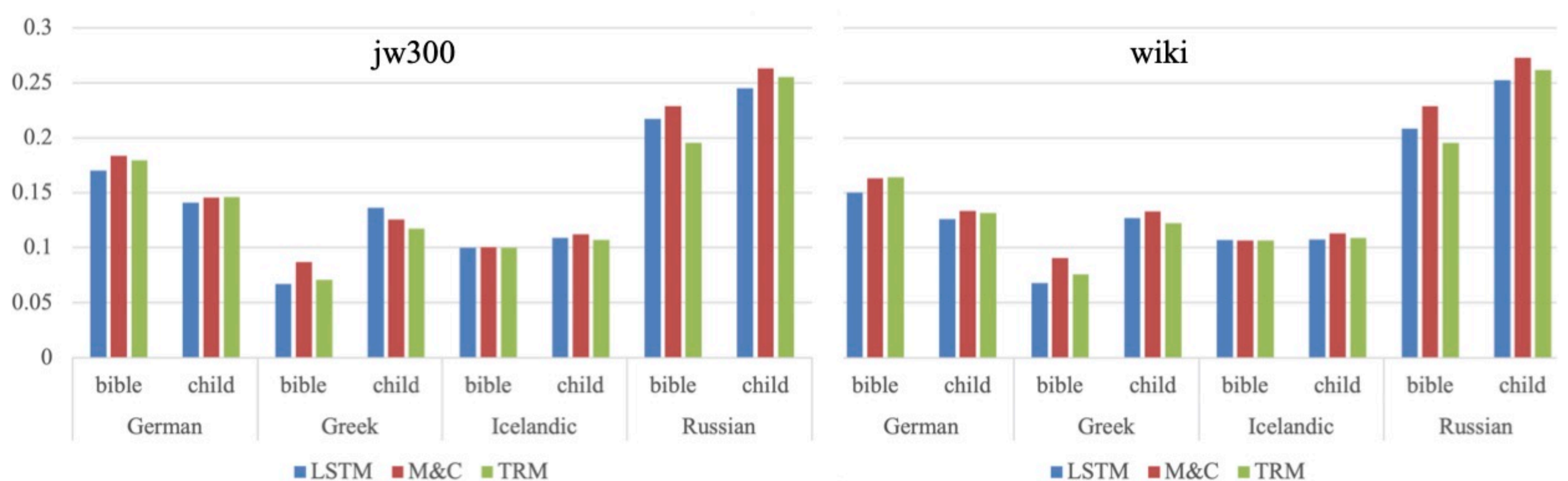
“My best friend broke our lamp”



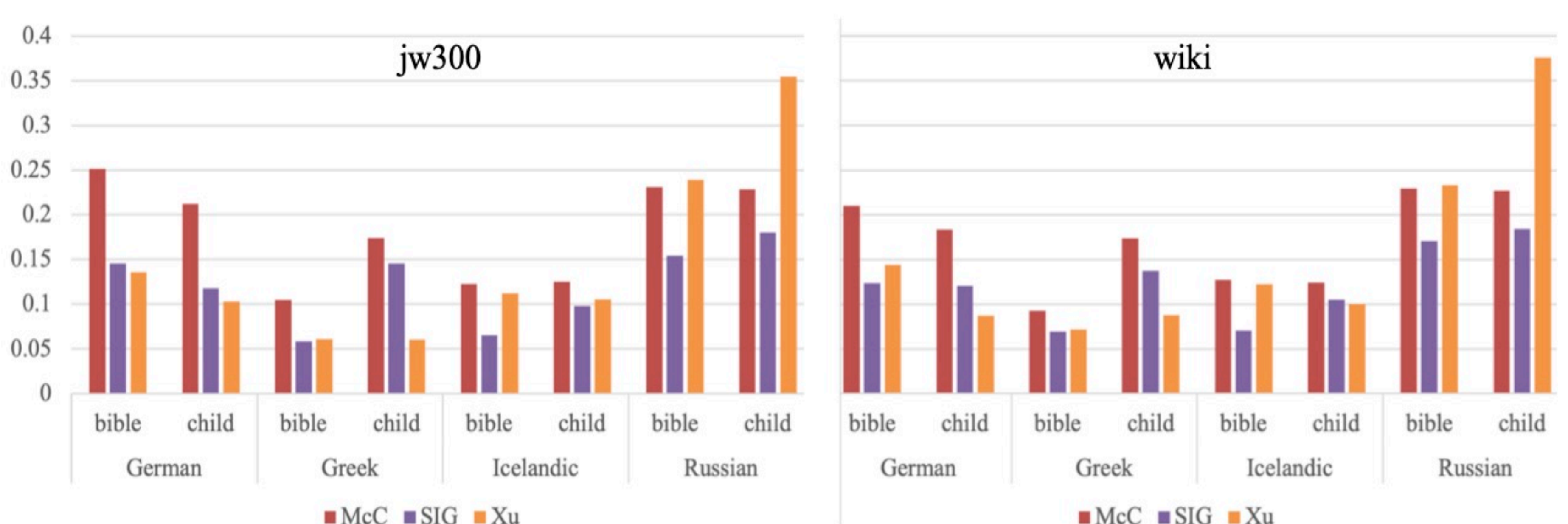
2. Inflect



Results



Choice of inflection system leads to low variation



Choice of paradigm clustering algorithm leads to high variation