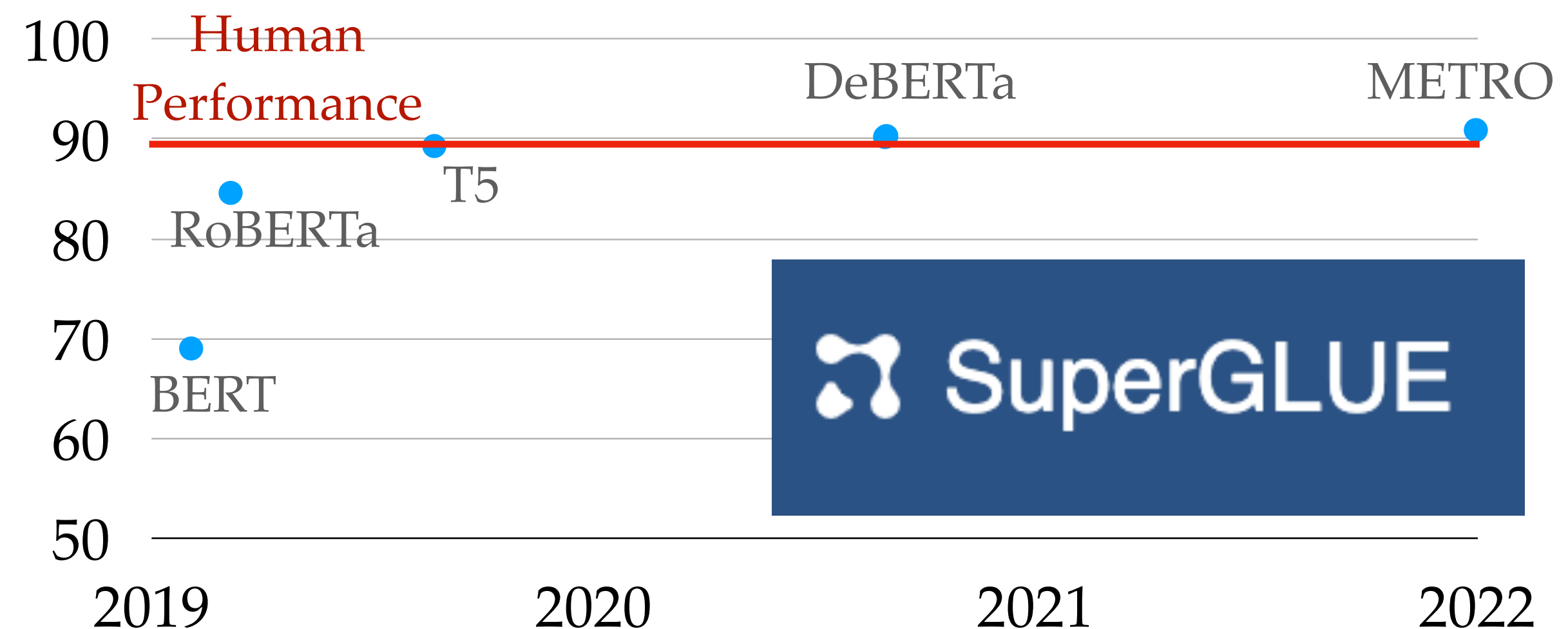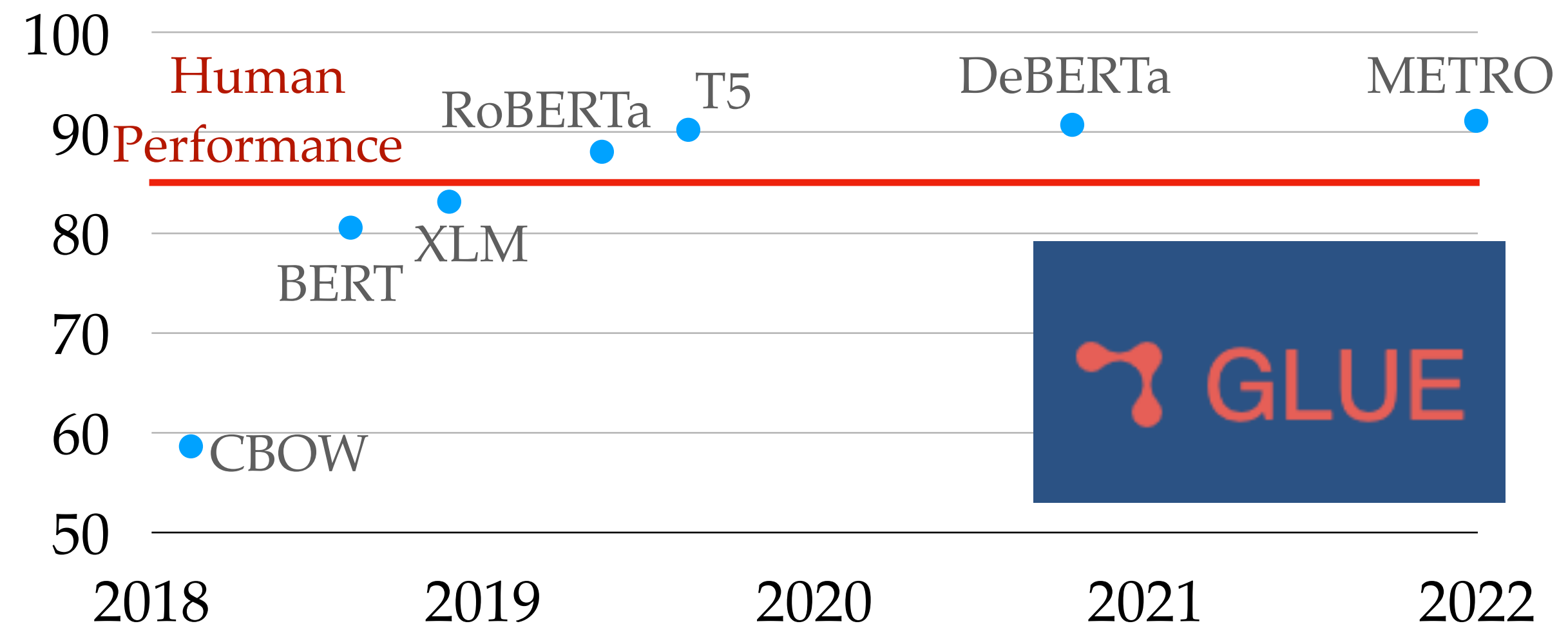# Crosslingual Sharing for Low-Resource Natural Language Processing

**Phoebe Mulcaire**

(in collaboration with Swabha Swayamdipta, Jungo Kasai, Nikolaos Pappas and Noah A. Smith)

# large-scale NLP is wildly successful

- BERT/GPT-3/OPT/etc.: billions of params, trained on billions (even hundreds of billions!) of tokens
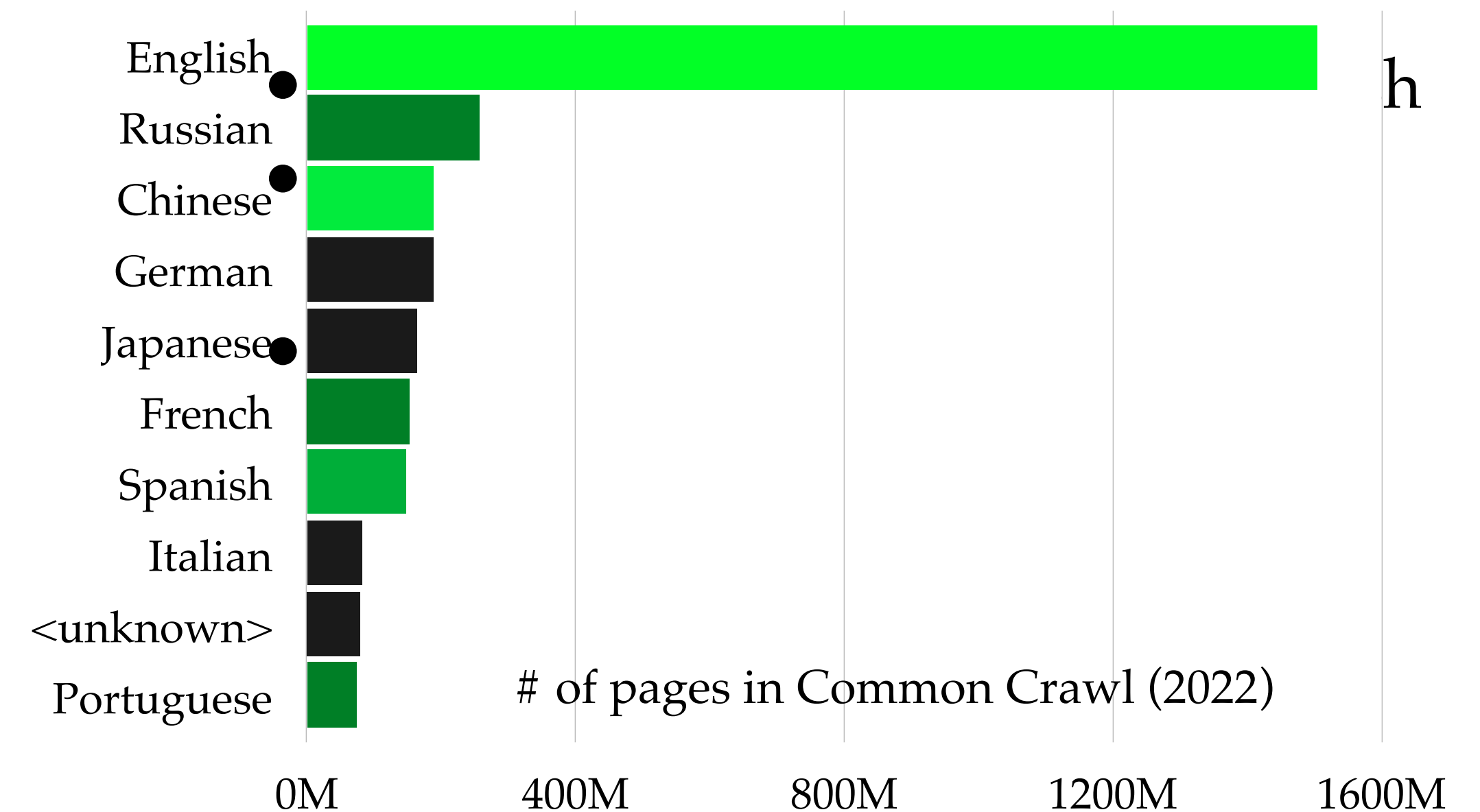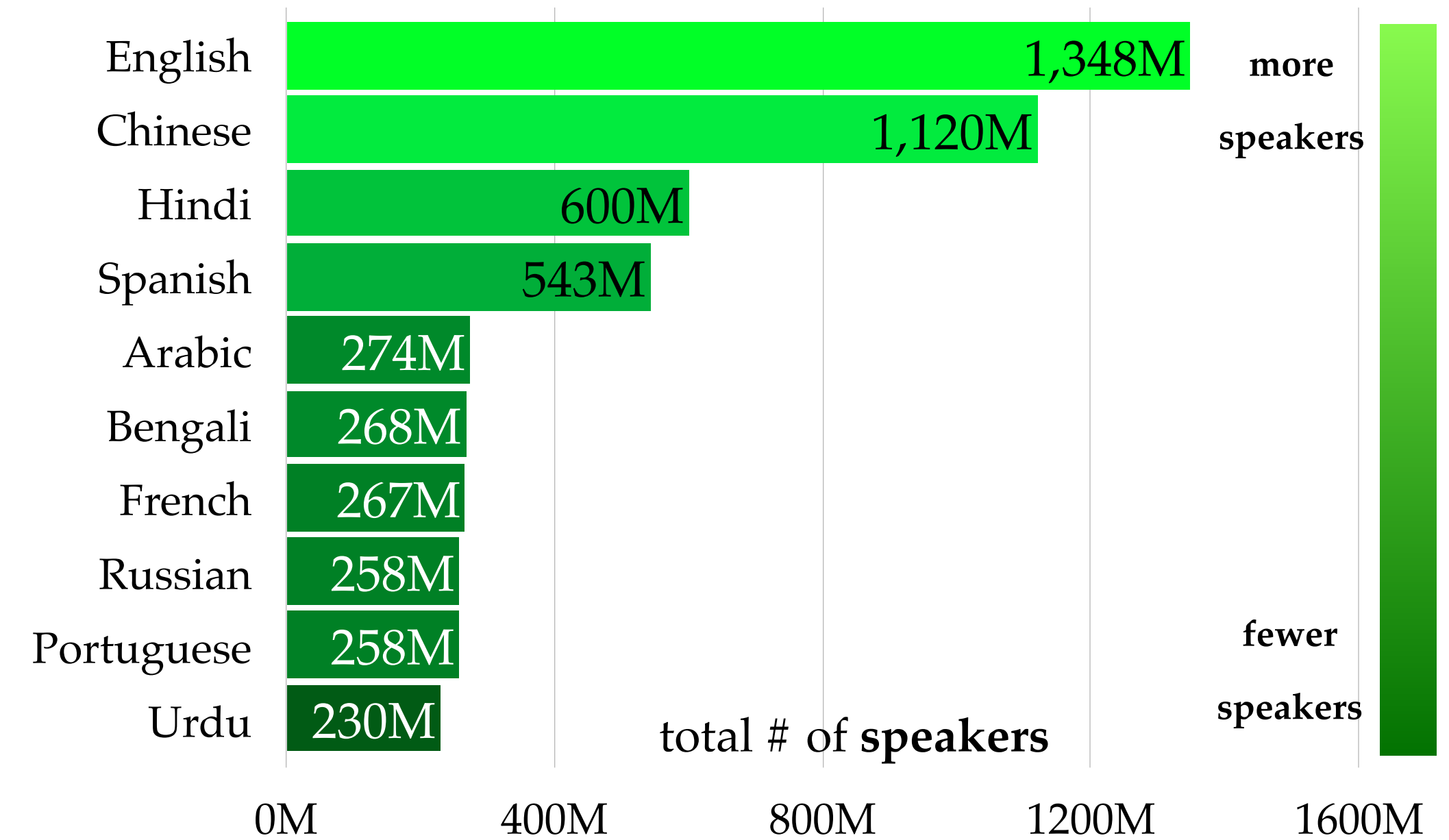
- result: English NLP is "solved"

# but… there's a resource gap

- Ethnologue records >7000 living languages[1]

- English is the most widely spoken language… but there's a fat tail

- the most widely *used* languages ≠ the ones with the most *resources*[2] (or research)

total # of **speakers**

| Language | Speakers |
|----------|----------|
| English | 1,348M |
| Chinese | 1,120M |
| Hindi | 600M |
| Spanish | 543M |
| Arabic | 274M |
| Bengali | 268M |
| French | 267M |
| Russian | 258M |
| Portuguese | 258M |
| Urdu | 230M |

# of pages in Common Crawl (2022)

- English
- Russian
- Chinese
- German
- Japanese
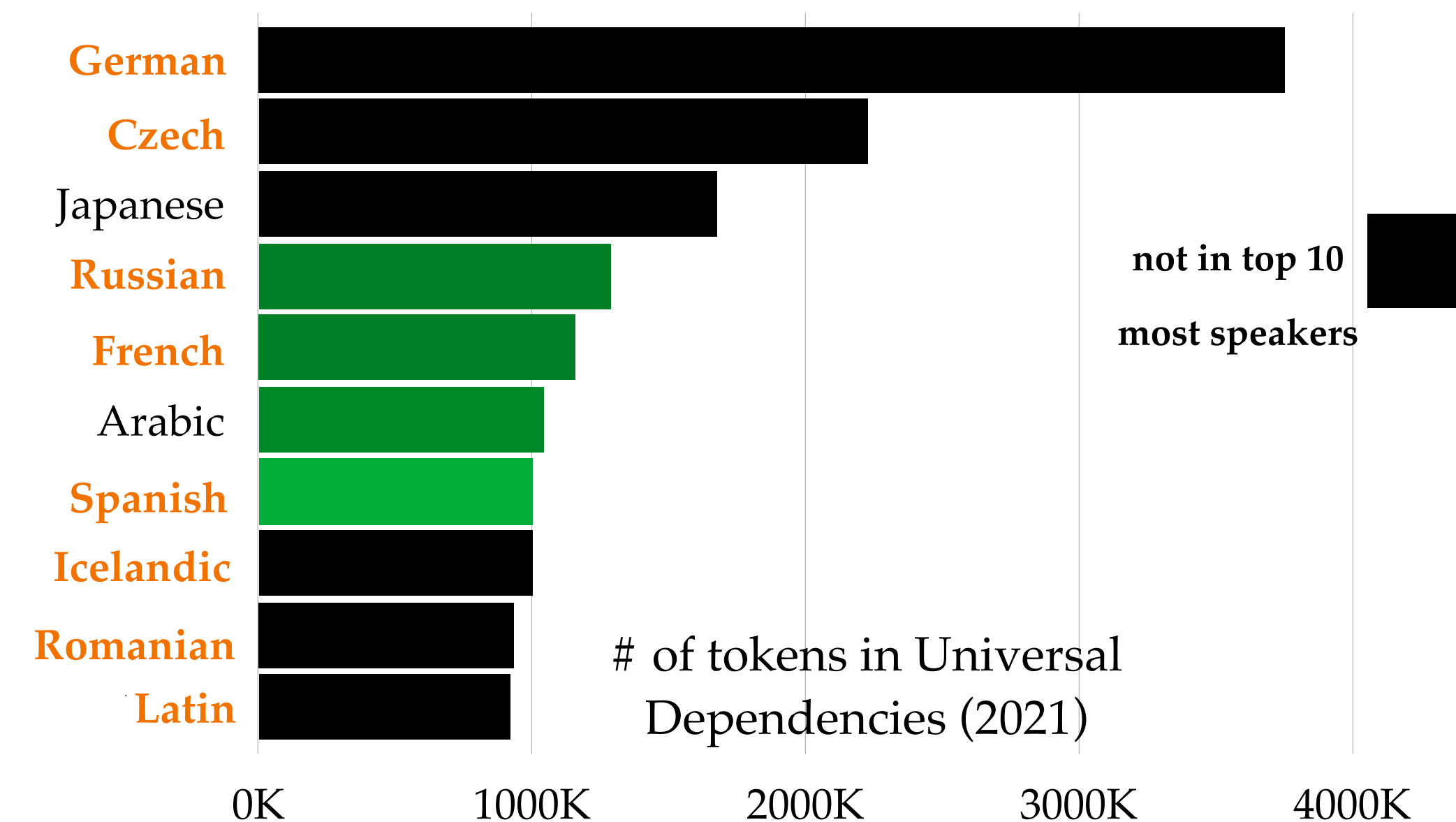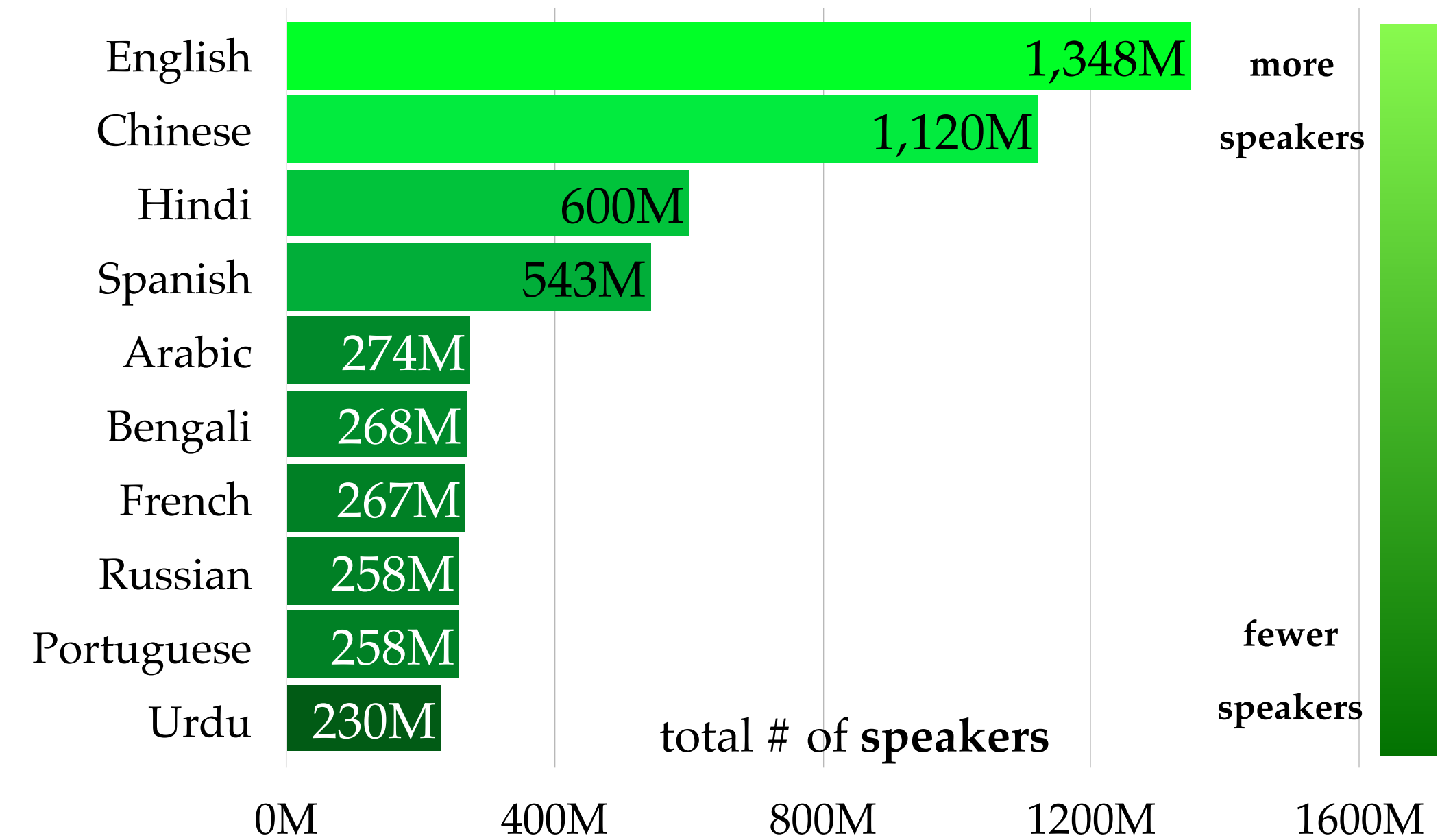- French
- Spanish
- Italian
- <unknown>
- Portuguese

# but… there's a resource gap

- Ethnologue records >7000 living languages[1]

- English is the most widely spoken language… but there's a fat tail

- the most widely *used* languages ≠ the ones with the most *resources*[2] (or research)

- most-resourced languages don't reflect world's linguistic diversity

- we can't replicate English NLP for every language

[1] Ethnologue: www.ethnologue.com

[2] Common Crawl

4

English 1,348M
Chinese 1,120M
Hindi 600M
Spanish 543M
Arabic 274M
Bengali 268M
French 267M
Russian 258M
Portuguese 258M
Urdu 230M

more speakers

fewer speakers

total # of **speakers**

0M   400M   800M   1200M   1600M

German
Czech
Japanese
Russian
French
Arabic
Spanish
Icelandic
Romanian
Latin

not in top 10 most speakers

# of tokens in Universal Dependencies (2021)

0K   1000K   2000K   3000K   4000K

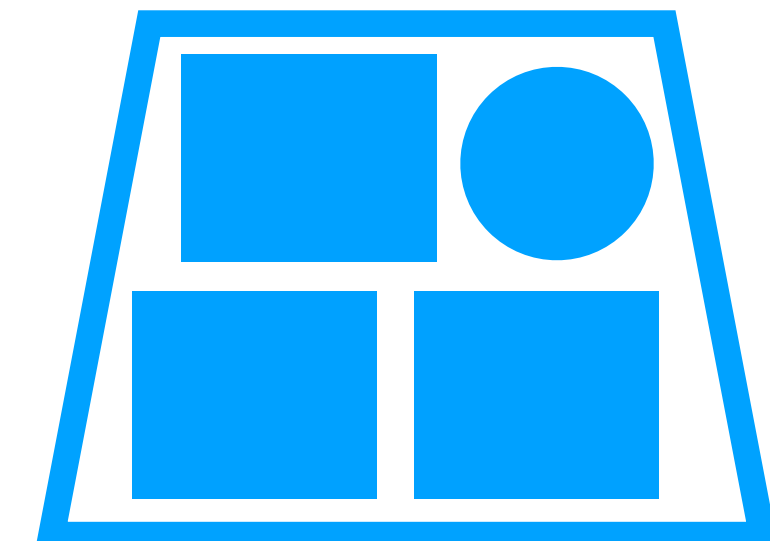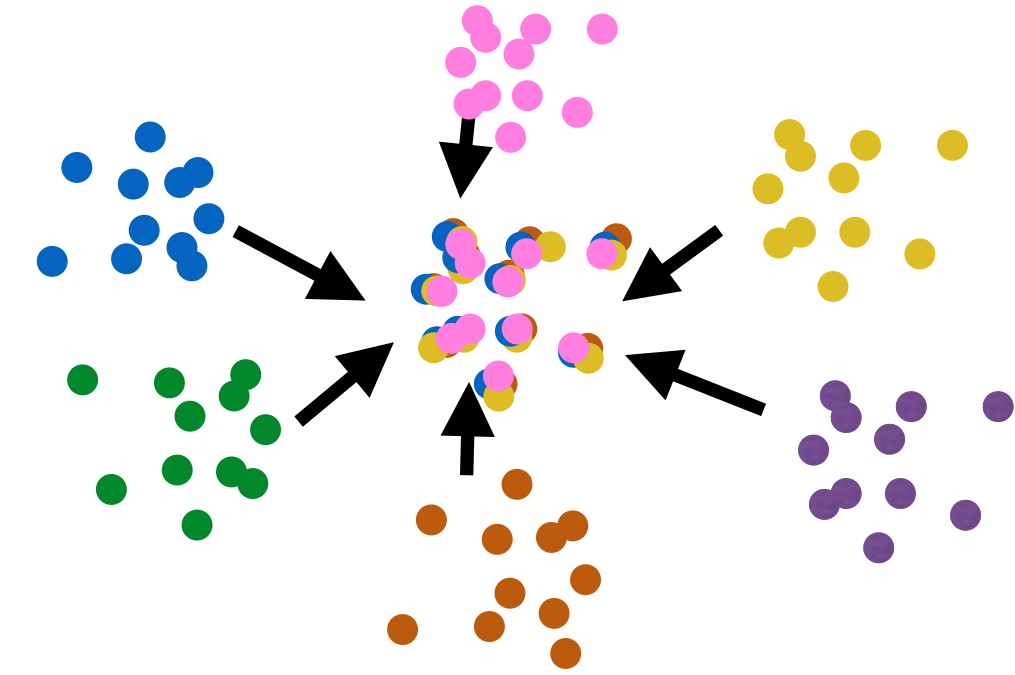# language-universal NLP

we want systems that:

- don't rely on large amounts of language-specific resources

- don't rely on large amounts of language-specific researcher effort (e.g. custom architecture choices)
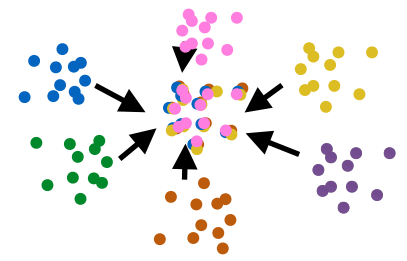
our focus:

- *crosslingual sharing*

- …via *polyglot* models

- …for *low-resource* settings

# outline

- *Polyglot Semantic Role Labeling (Ch. 2)*

  - supervised, linguistic structure prediction

- *Polyglot Language Modeling (Ch. 3)*

  - language models for word representations

- *Grounded Compositional Output Embeddings (Ch. 4-5)*
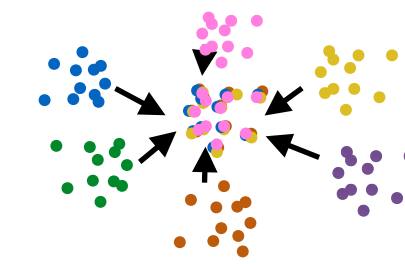
  - low-resource language models

# semantic role labeling

I think Peter even made some deals with the gorillas .
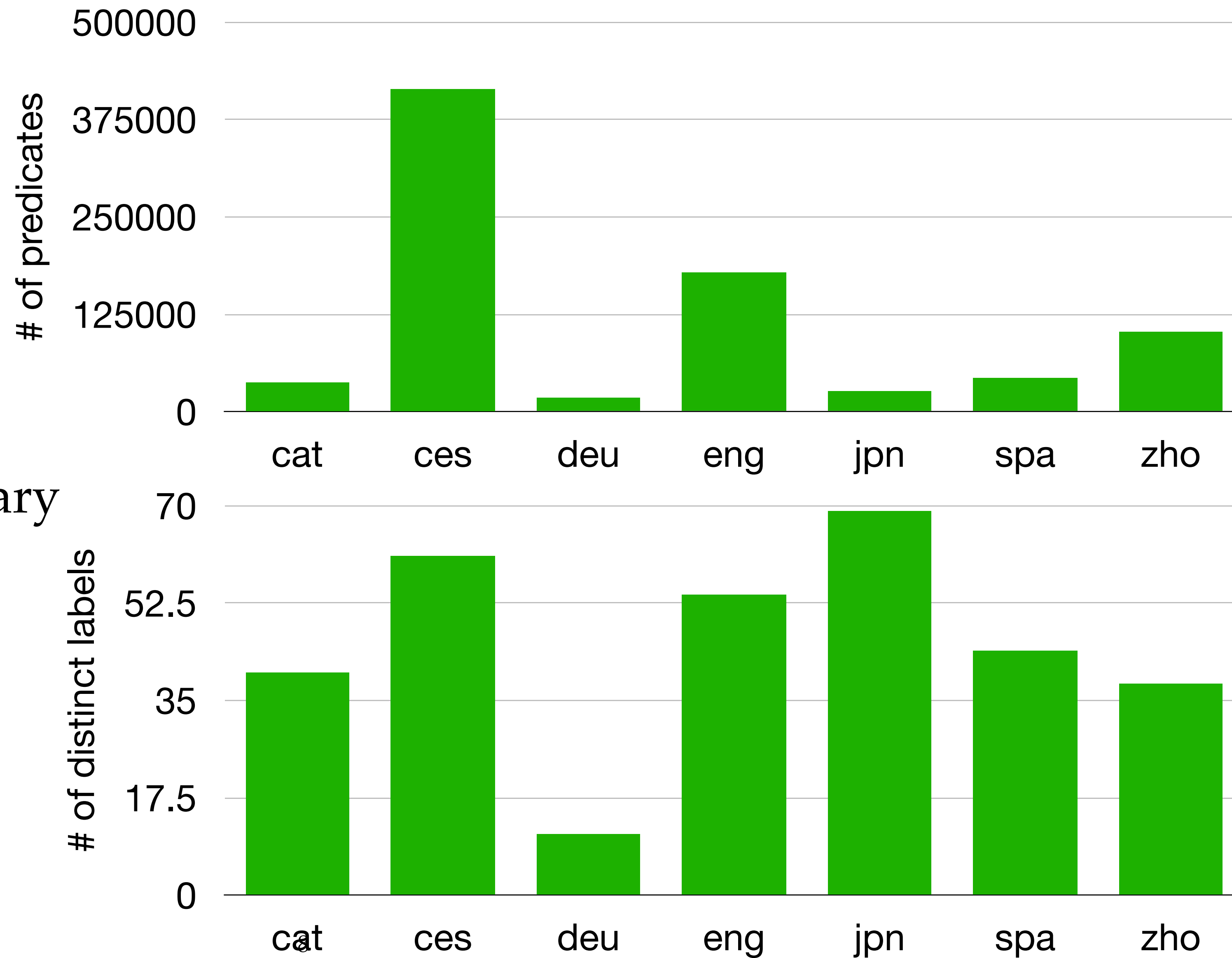O   O    **A0**  **AM-ADV**  O   O    **A1**  **AM-ADV**  O   O

Pero el suizo difícilmente atacará a Rominger en la montaña .
O  O **arg0-agt**  **argM-adv**   O   O **arg1-pat** **argM-loc**  O   O

Četrans oslovil sedm velkých evropských výrobců nákladních automobilů.
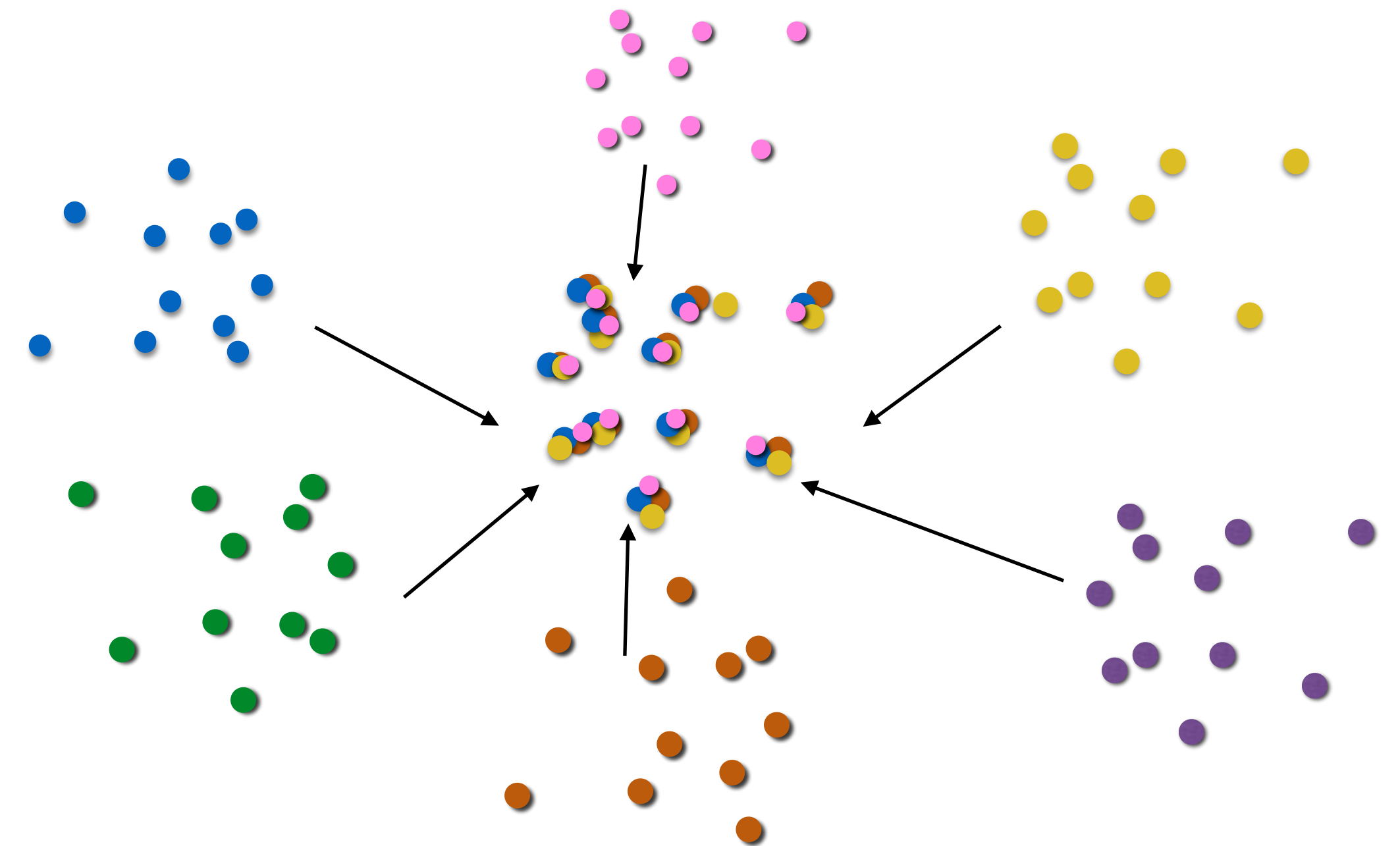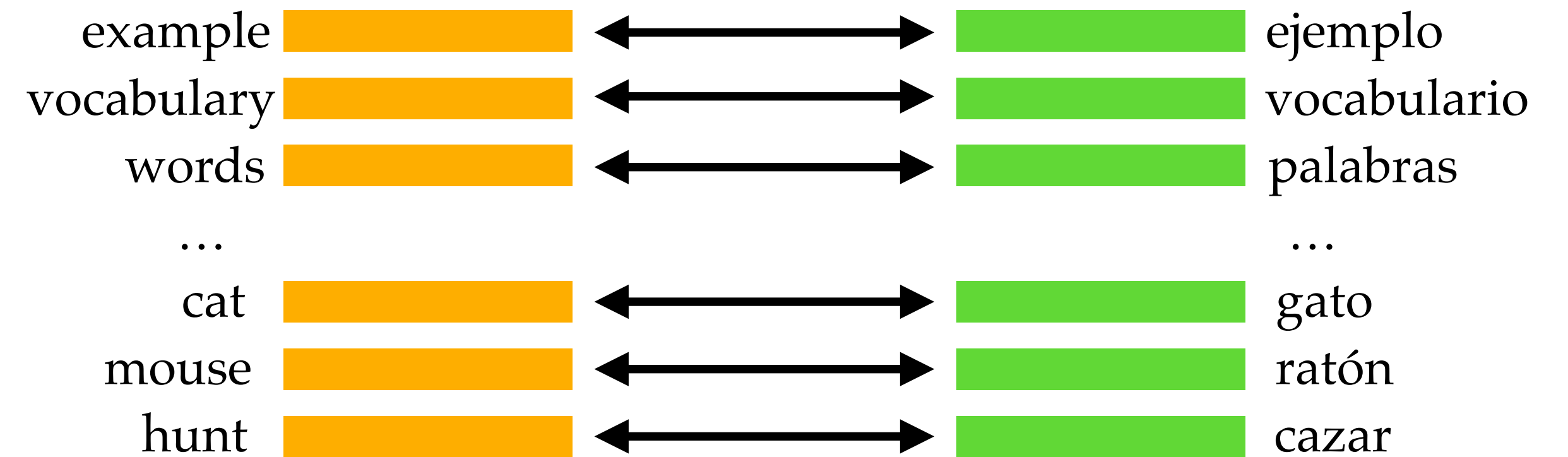O    O   **RSTR**  **RSTR**    **RSTR**    O     O    **PAT**

# CoNLL 2009

- format is the same, but:

- data is wildly imbalanced between languages—independent models (e.g. Zhao et al., 2009) would vary

- output labels vary—annotation projection (e.g. Padó and Lapata, 2005) is ruled out
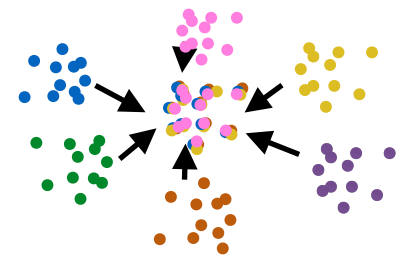
# our approach: polyglot semantics

- multilingual word vectors

- produce word vectors for each language based on co-occurrence statistics

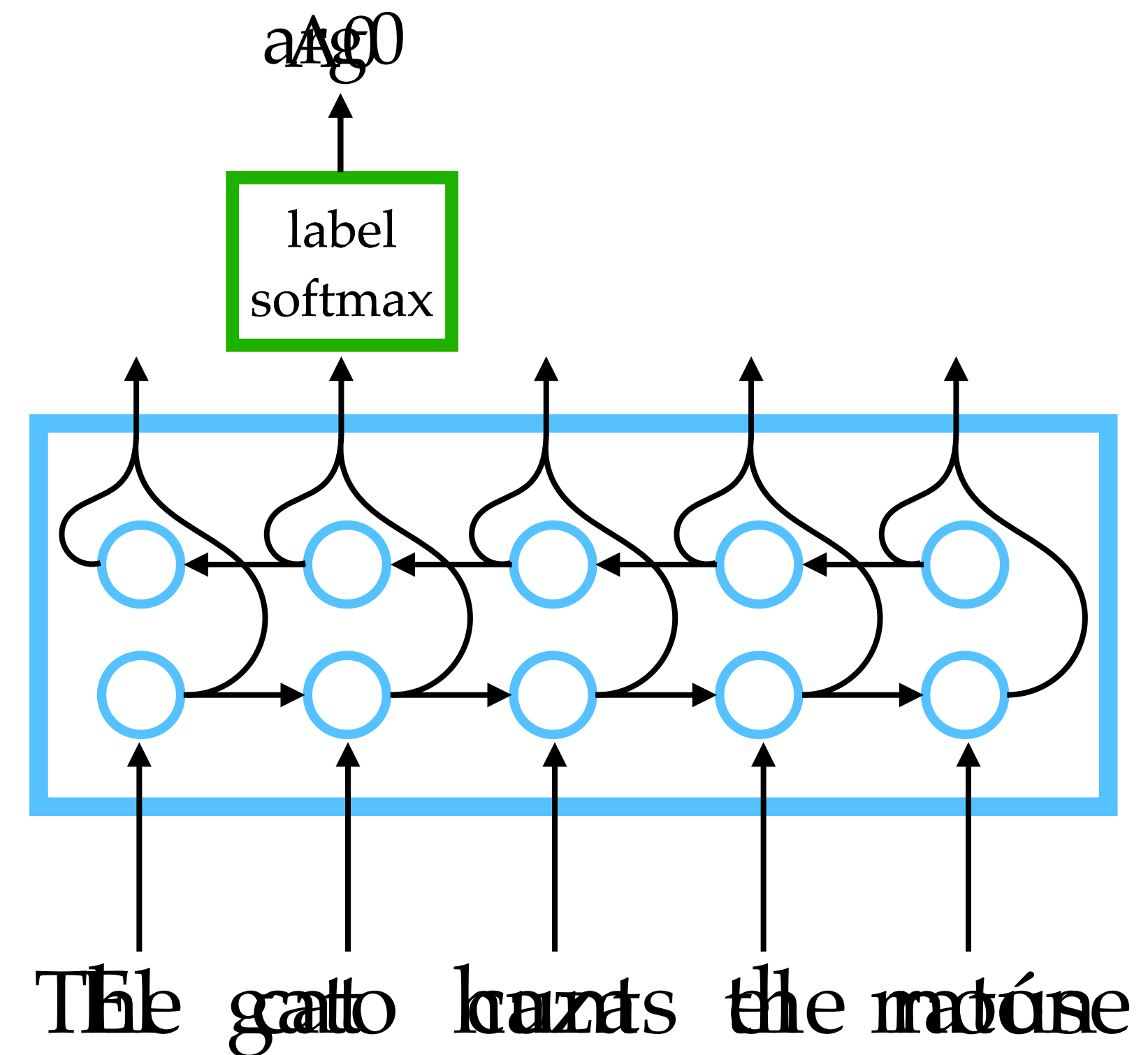- align to match English using a bilingual dictionary[1]

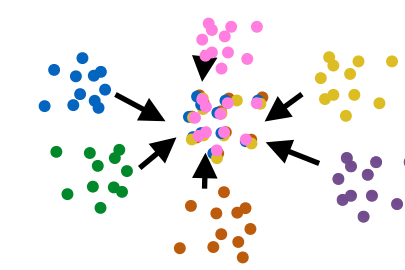| example | | ejemplo |
| vocabulary | | vocabulario |
| words | | palabras |
| ... | | ... |
| cat | | gato |
| mouse | | ratón |
| hunt | | cazar |

[1] Faruqui et al. (2014); Ammar et al. (2016)

# our approach: polyglot semantics

- task model based on a (then) SOTA monolingual model[1]

- multilingual word vector inputs

- sharing in parameters: deep bi-LSTM
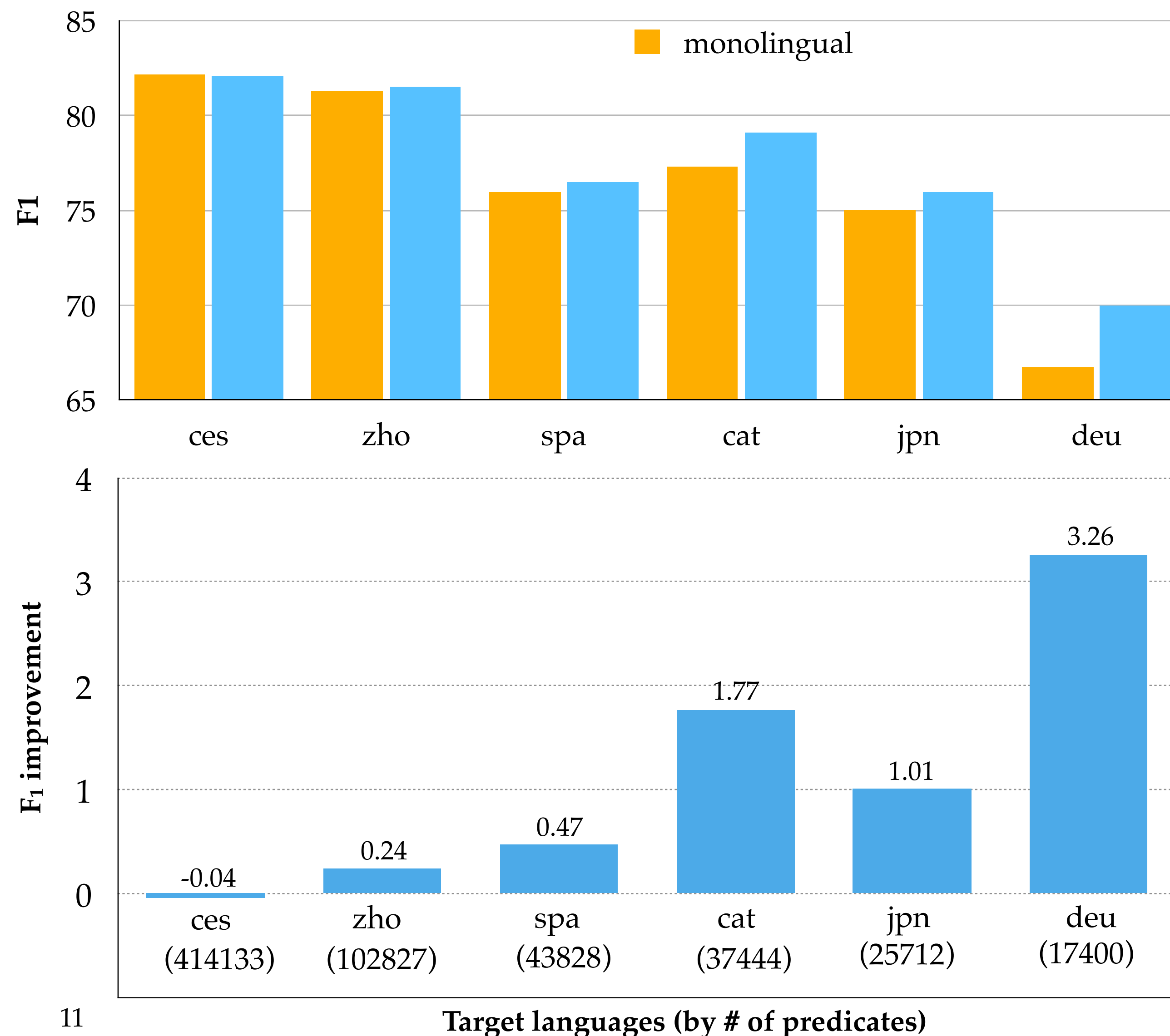
- independent label embeddings



[1] He et al. (2017)

# polyglot SRL: experiment and results

- for each non-English language, train

  - a monolingual model

  - a polyglot model with English

- most languages improve from polyglot training

- lower-resource languages benefit more

# polyglot SRL: takeaways

- can represent data from multiple languages in a shared representation space

- by sharing data across languages, you can improve performance

- lower-resource languages benefit more

- different annotation schemes are not a strict barrier

# problems with word representations

- word vectors are great, but limited:

  - poorly handle polysemy

  - similar words trained independently

# solution: contextualized word representations

- train a language model first

Here's an e x a m p l e

char-level CNN

language model → next word predictions

sentence from a corpus

# solution: contextualized word representations

- train a language model first

- feed hidden states to the task model as input



Here's an  e x a m p l e
sentence from a supervised dataset

char-level CNN

language model → task model → task predictions

ELMo (Peters et al. 2018); c.f. BERT (Devlin et al. 2018)

# an intuitive approach: alignment of averages

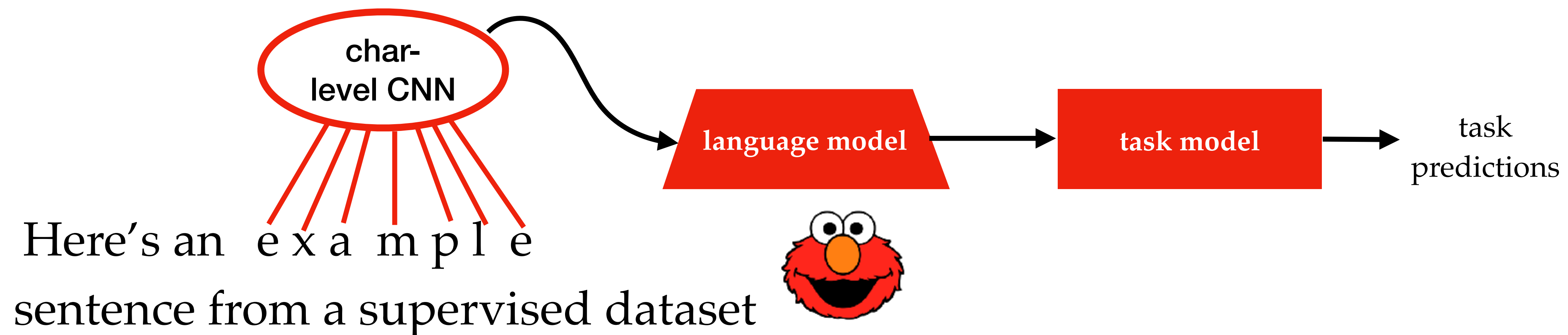- train separate language models for each language

- align "average" embedding across contexts with a bilingual dictionary



English sentence → **English language model** → English next word predictions

Spanish sentence → **Spanish language model** → Spanish next word predictions

English-Spanish dictionary

English-Spanish alignment transform → **polyglot** task model → task predictions

16

# polyglot contextualization: Rosita

- train a language model first

- feed hidden states to the task model as input

- for a multilingual model, we need a multilingual language model!

English next word
predictions

**char-
level CNN**

**polyglot**
language model

**polyglot**
task model

English task
predictions

Here's an e x a m p l e

multilingual BERT (Devlin et al. 2018) and XLM
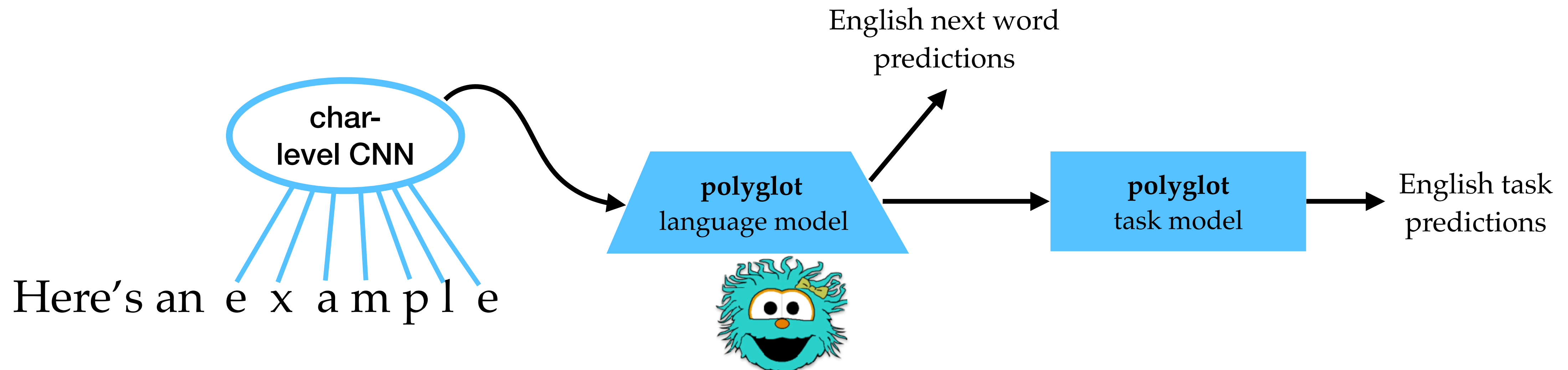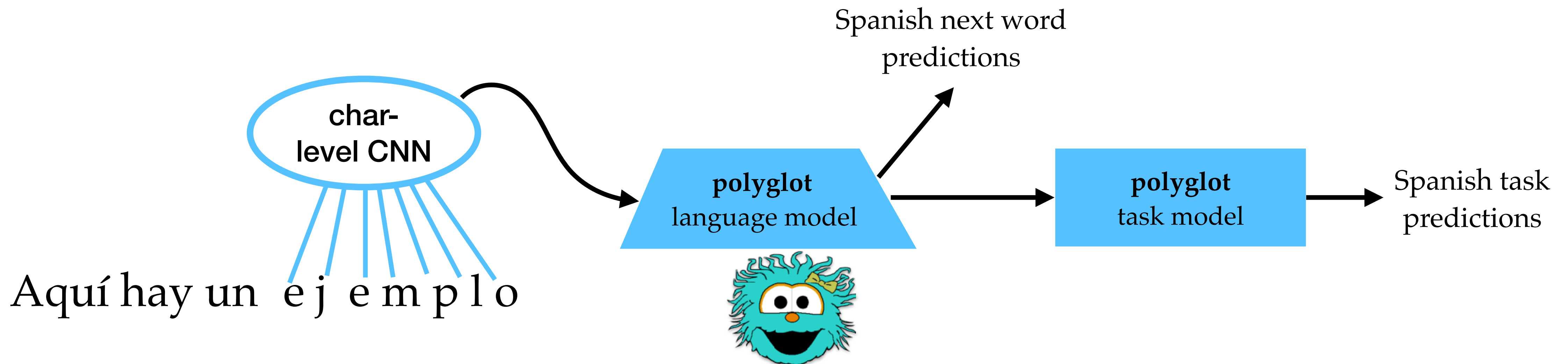(Lample and Conneau, 2019) also use this approach

# polyglot contextualization: Rosita

- train a language model first

- feed hidden states to the task model as input

- for a multilingual model, we need a multilingual language model!

Spanish next word
predictions

**char-level CNN**

**polyglot** language model

**polyglot** task model

Spanish task predictions
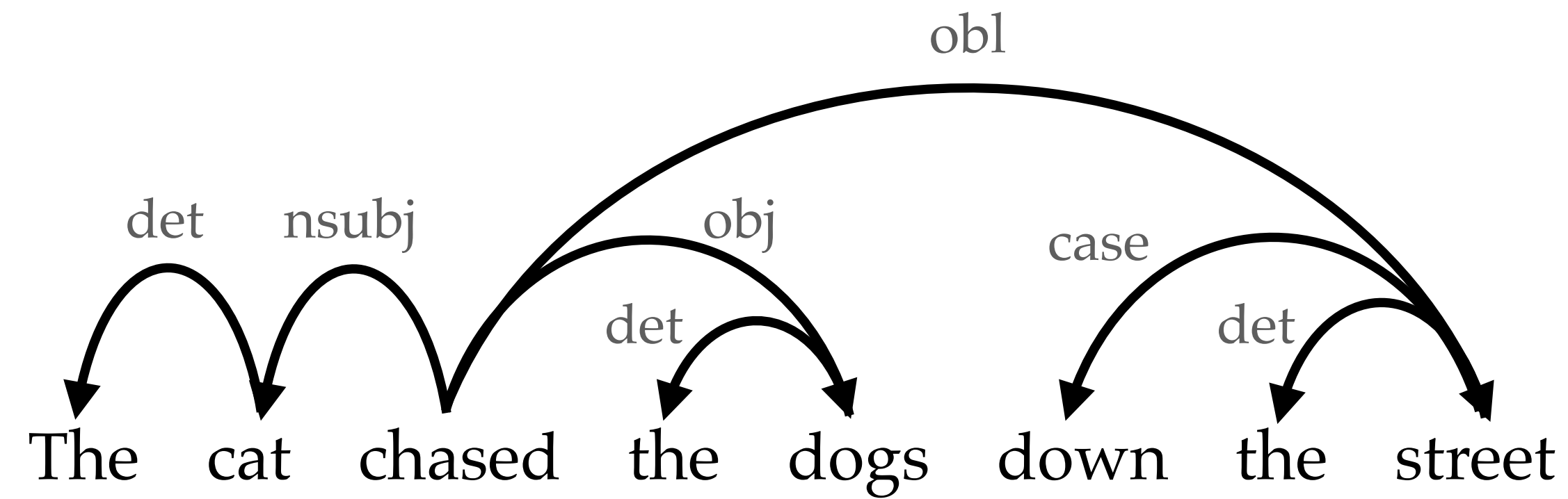
Aquí hay un e j e m p l o

multilingual BERT (Devlin et al. 2018) and XLM
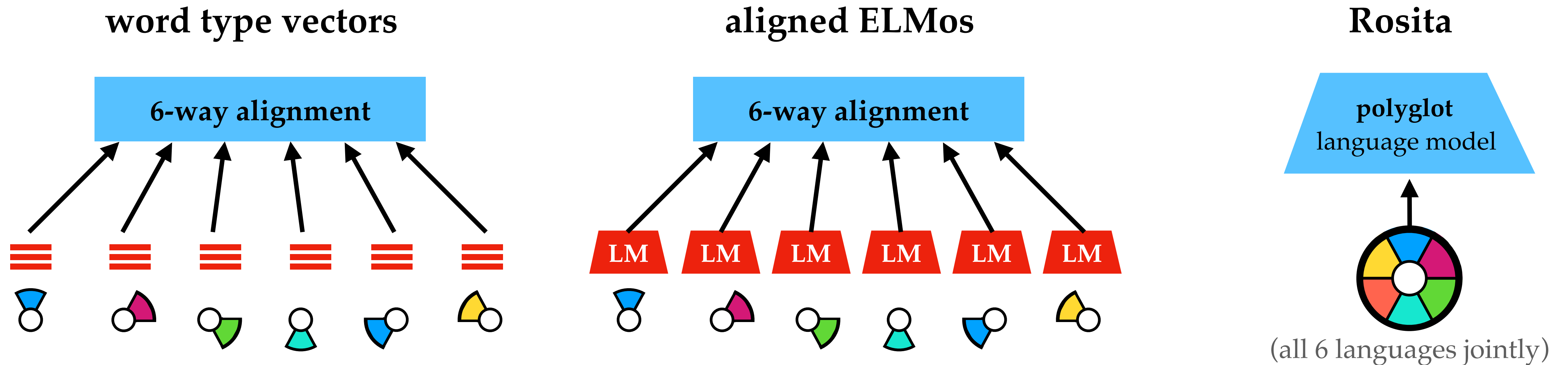(Lample and Conneau, 2019) also use this approach

# polyglot LMs: experiments

- Universal Dependencies syntax parsing (which *does* match across languages)

The cat chased the dogs down the street

*det, nsubj, obl, obj, det, case, det*

# polyglot LMs: experiments

- Universal Dependencies syntax parsing (which *does* match across languages)

- "zero-target" evaluation:

  - language models (or word vectors) combining six languages



**word type vectors**

**aligned ELMos**

**Rosita**

(all 6 languages jointly)
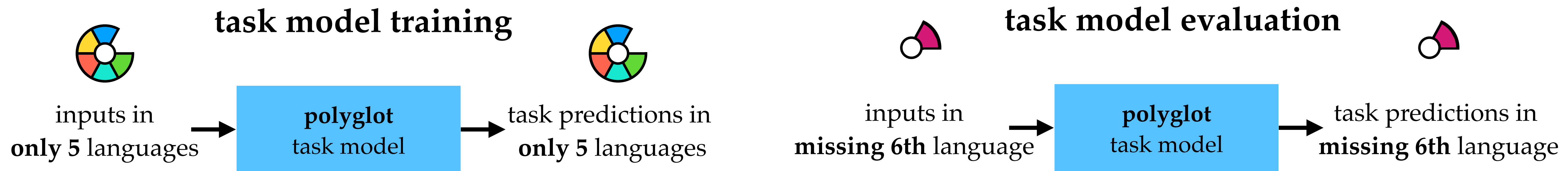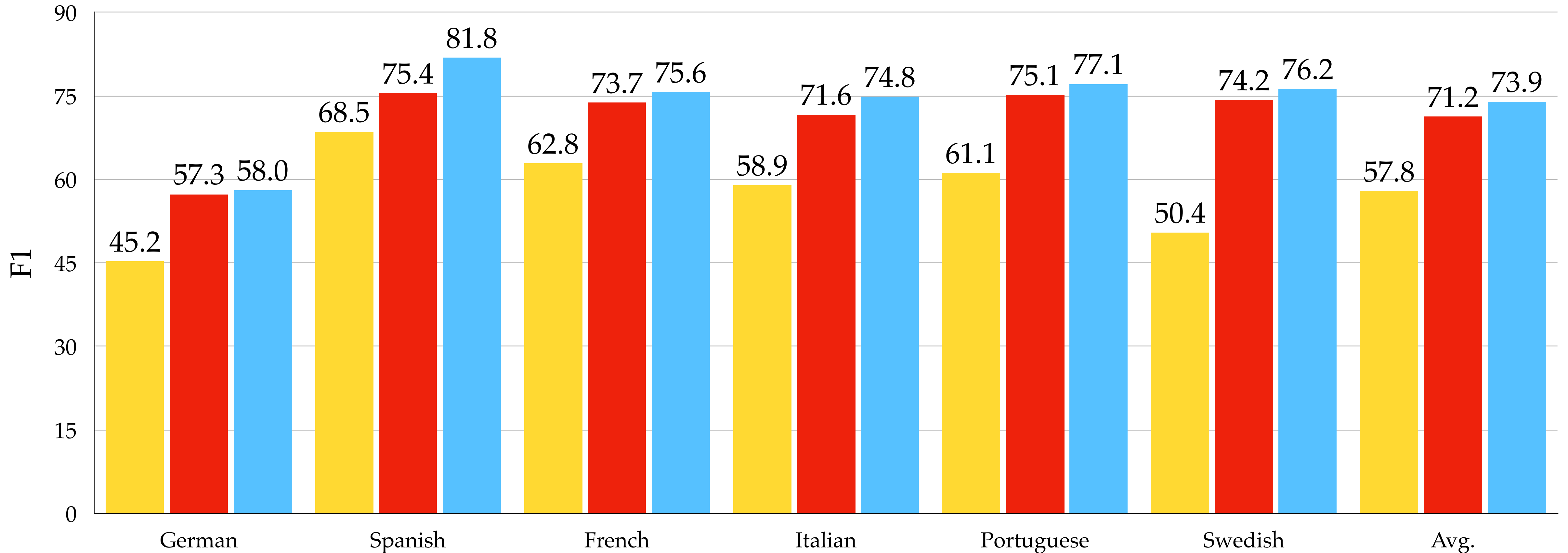
# polyglot LMs: experiments

- Universal Dependencies syntax parsing (which *does* match across languages)

- "zero-target" evaluation:

    - language models (or word vectors) combining six languages

    - six parsers, each trained on only five—evaluate on the missing language

**task model training**

inputs in
**only 5** languages → | **polyglot** task model | → task predictions in
**only 5** languages

**task model evaluation**

inputs in
**missing 6th** language → | **polyglot** task model | → task predictions in
**missing 6th** language
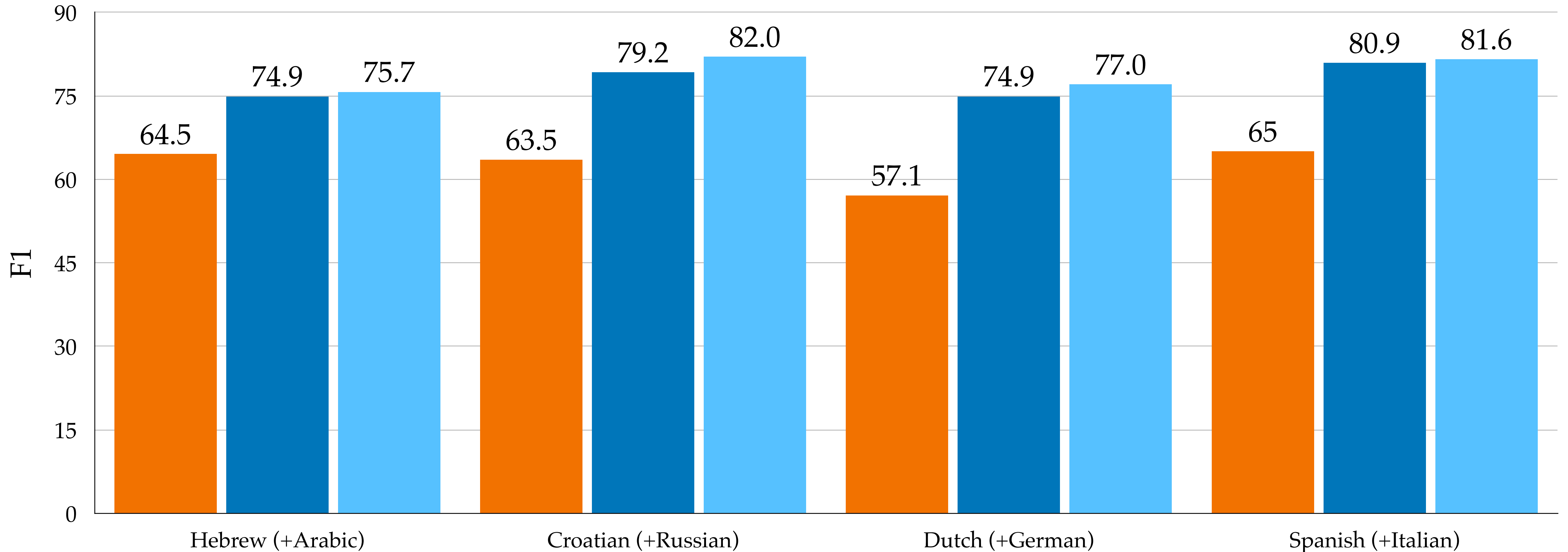
# polyglot LMs: zero-target results



type vectors vs aligned LMs vs polyglot LMs: Universal Dependencies parsing F1

# polyglot LMs: diverse languages



Legend: ■ mono ELMo  ■ Rosita (tgt+English)  ■ Rosita (tgt+similar)

- Hebrew (+Arabic): 64.5, 74.9, 75.7
- Croatian (+Russian): 63.5, 79.2, 82.0
- Dutch (+German): 57.1, 74.9, 77.0
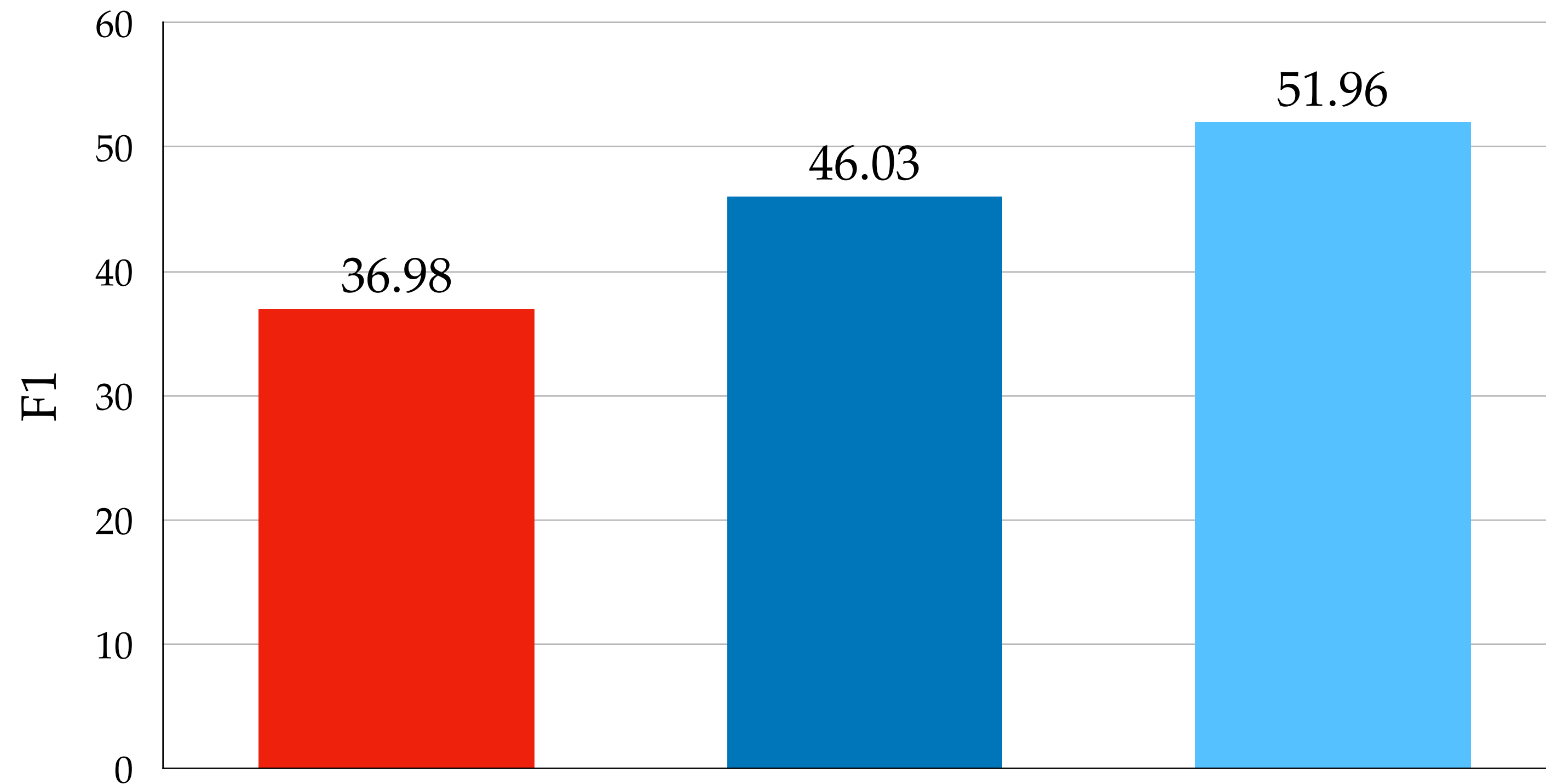- Spanish (+Italian): 65, 80.9, 81.6

F1 (y-axis)

sharing with English vs a similar language: Universal Dependencies parsing F1

23

# polyglot LMs: true low-resource



aligned ELMos (Kazakh+Turkish)    Rosita (Kazakh+English)    Rosita (Kazakh+Turkish)

Kazakh, a real low-resource language: Universal Dependencies parsing F1
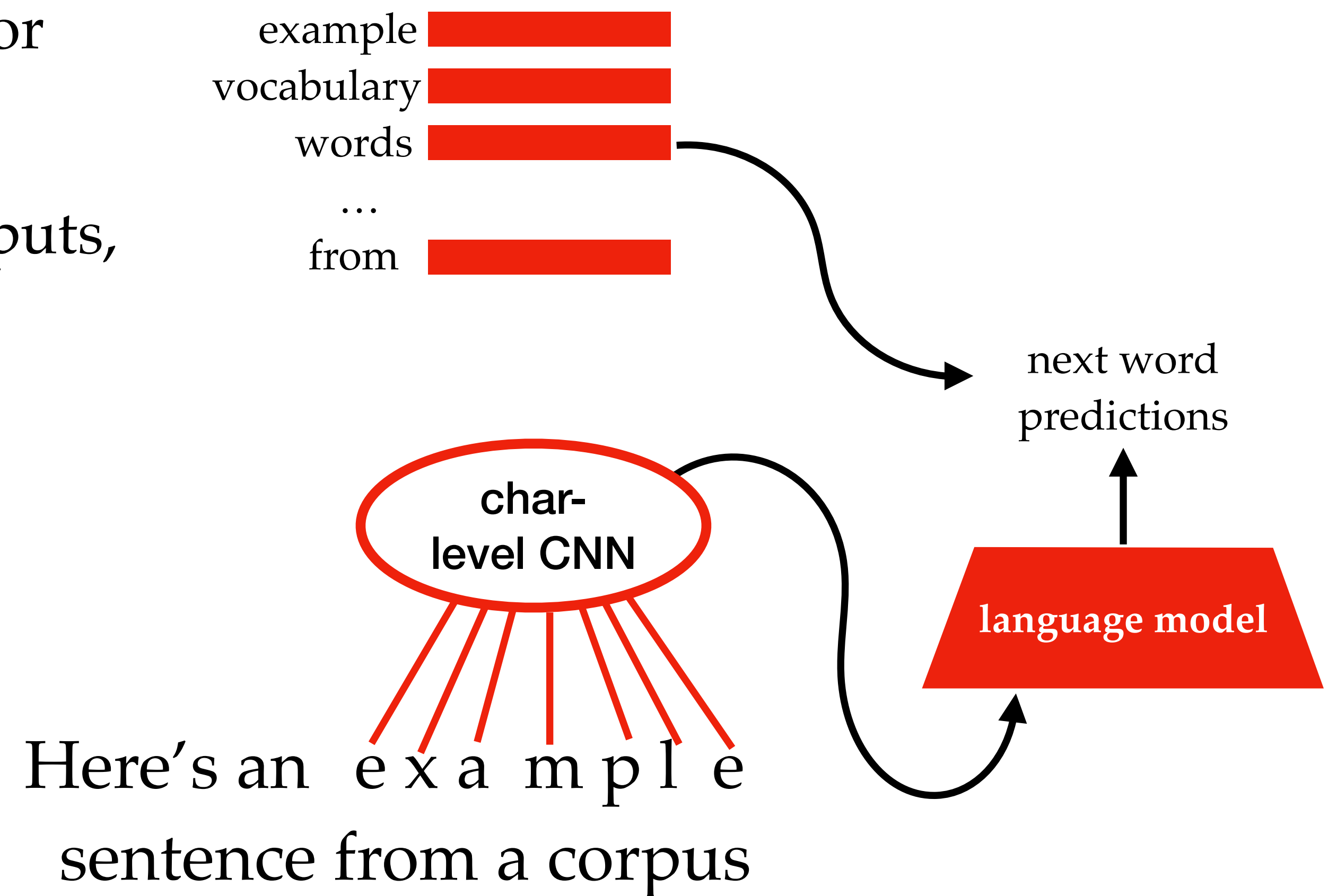
# polyglot LMs: takeaways

- polyglot modeling with contextualized representations works!

- don't need *any* explicit crosslingual supervision for multilinguality!

- polyglot training captures something alignment doesn't

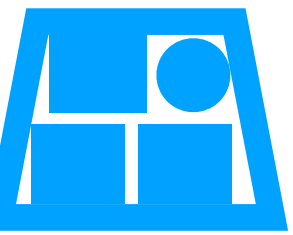- lots more experiments in the papers/Chapter 3

Mulcaire et al. 2019a: Polyglot Contextualized Representations Improve Crosslingual Transfer

Mulcaire et al. 2019b: Low-Resource Parsing With Crosslingual Contextualized Representations

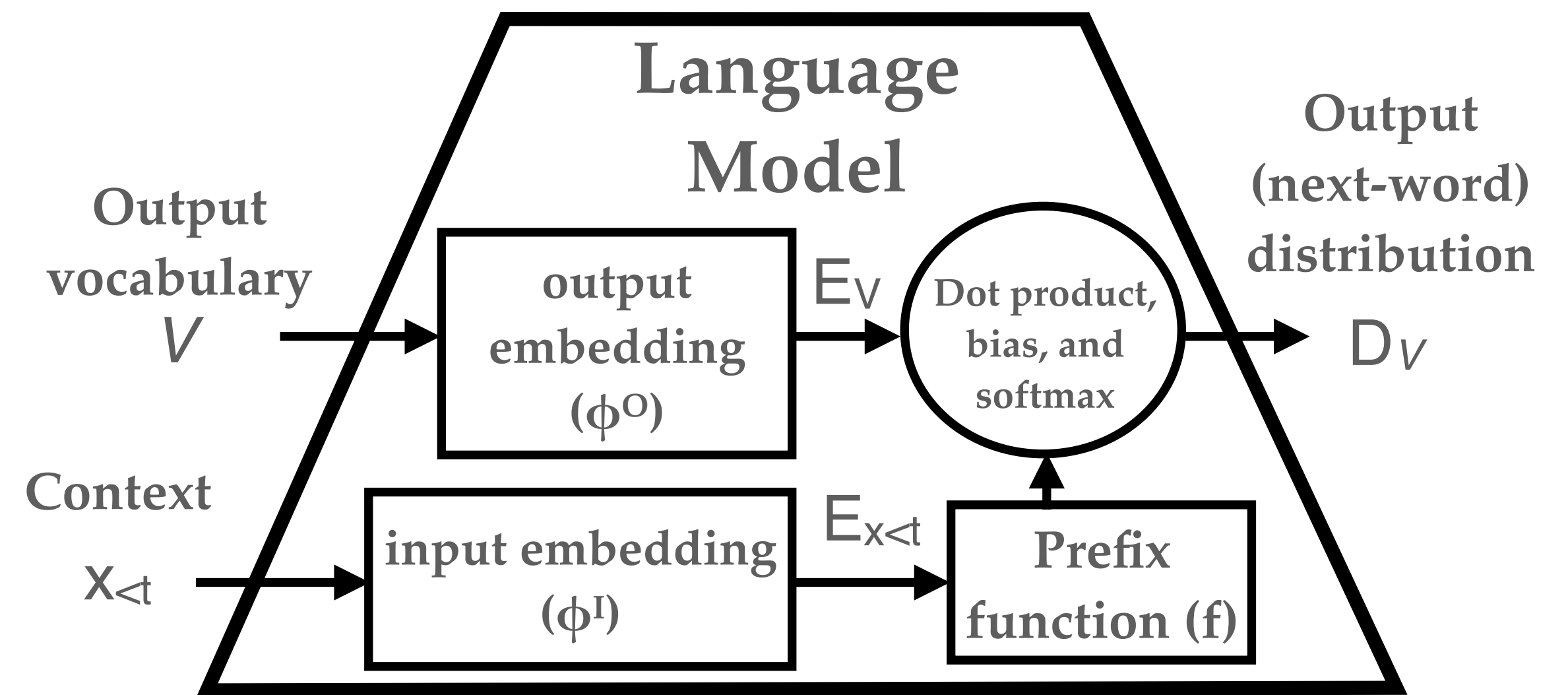# what if our language model training data is small?

- rare/out-of-domain words might get poor representations

- ELMo and Rosita have compositional inputs, but outputs are just type embeddings

- improve language models:

  - handle unknown words in test

  - improve rare word representations
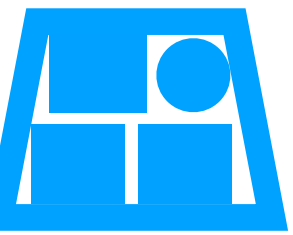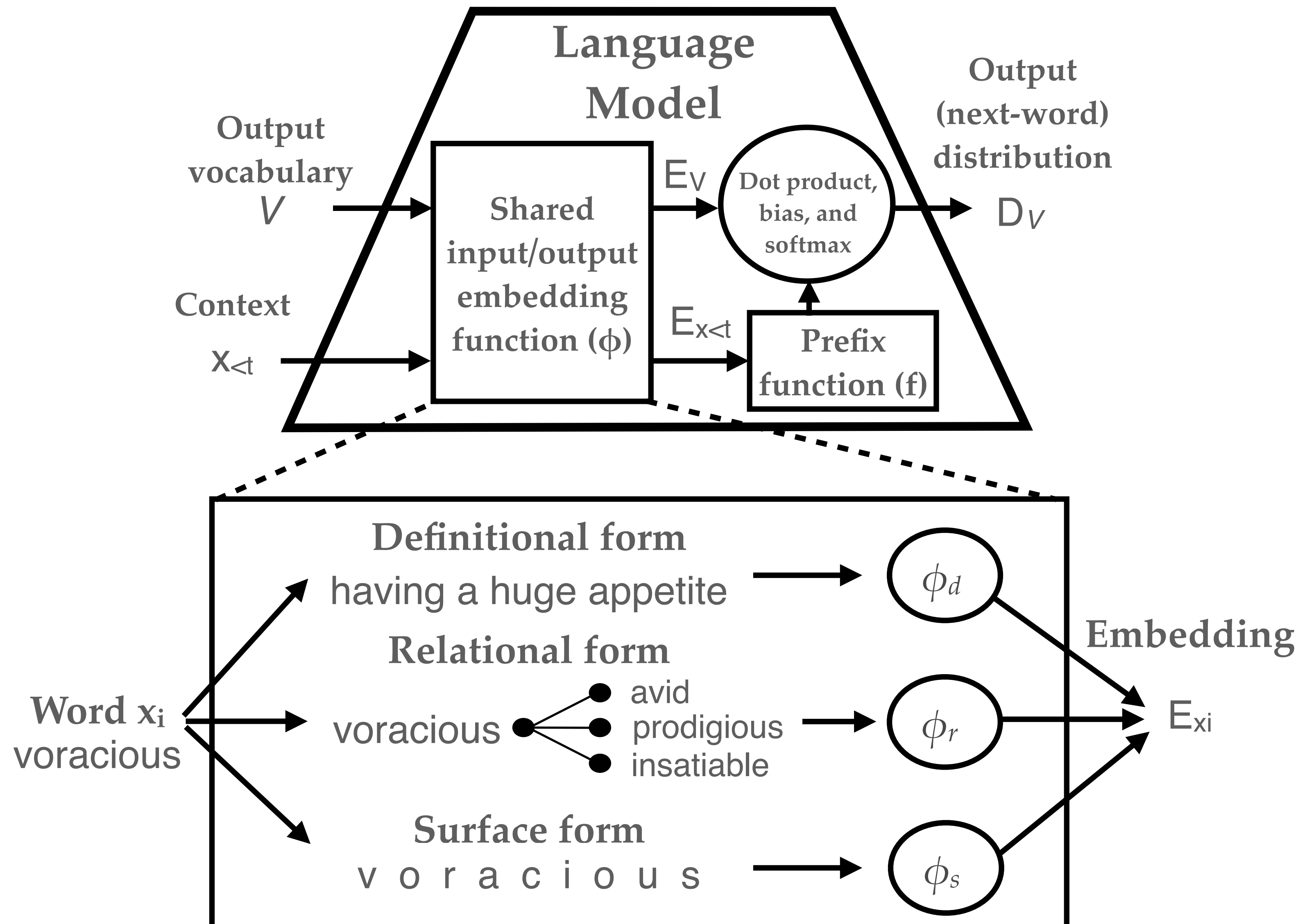
  - sample-efficient learning

example
vocabulary
words
...
from

next word
predictions

char-
level CNN

language model

Here's an  e x a m p l e
sentence from a corpus

# pieces of a language model

- input embedding, output embedding, prefix function

- traditional/lookup: input and output are lookup tables

- ELMo: input is a CNN, output is lookup

- many other possibilities: tied, bilinear, adaptive



**Language Model**

Output vocabulary $V$ → output embedding $(\phi^O)$ → $E_V$ → Dot product, bias, and softmax → $D_V$ Output (next-word) distribution

Context $x_{<t}$ → input embedding $(\phi^I)$ → $E_{x_{<t}}$ → Prefix function (f)
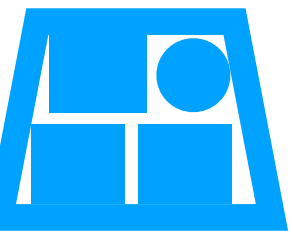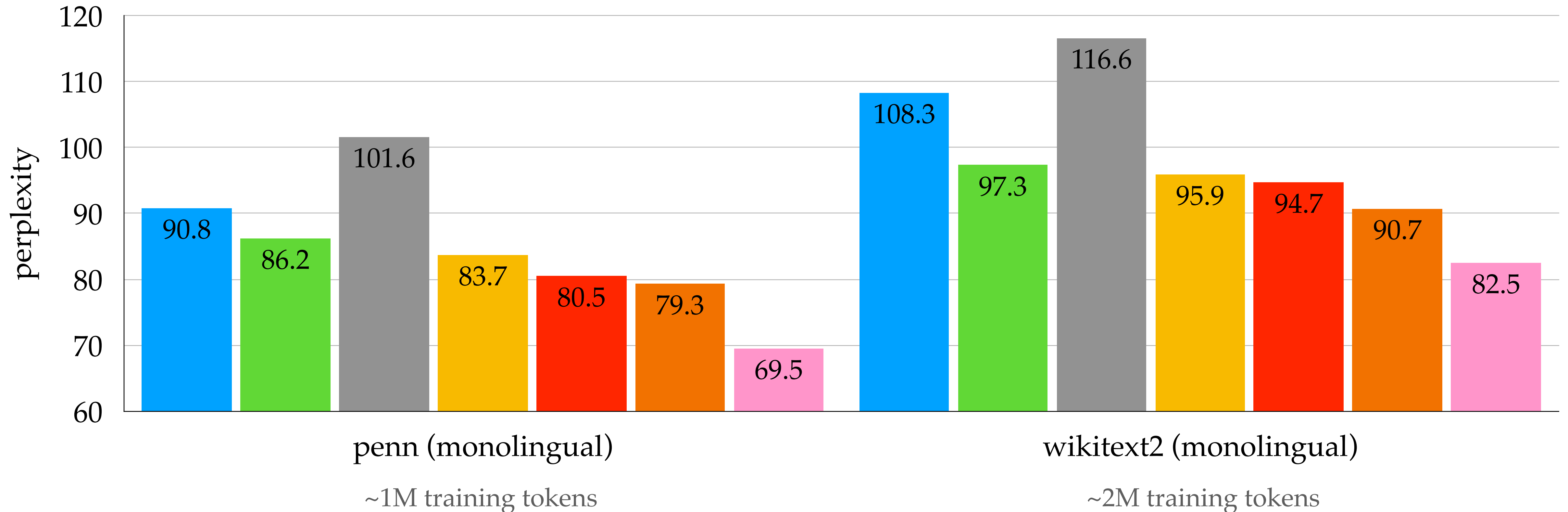
# grounded compositional outputs (GroC)



- use the same composition function for input and output

- combine surface form with relational and definitional features (from WordNet)

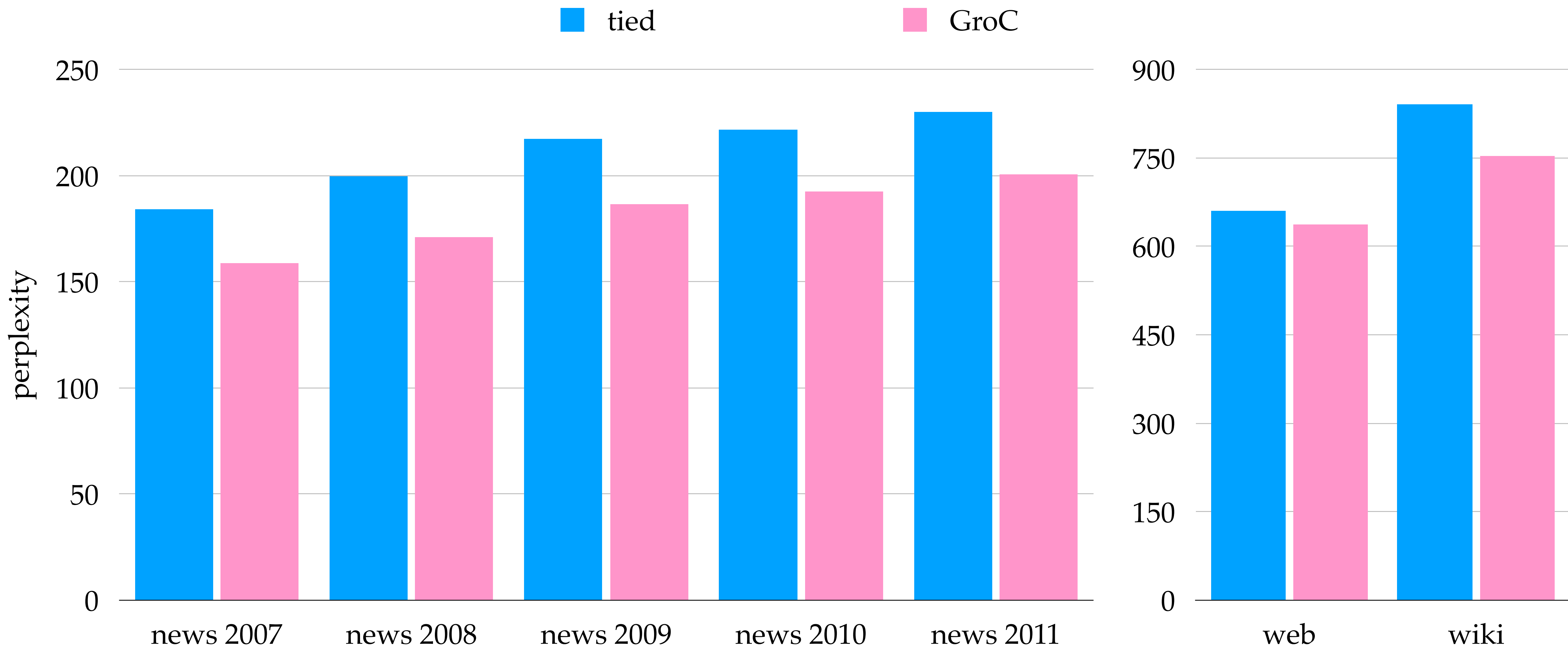- (also have a residual network applied to output in some cases)

# conventional language modeling

- perplexity: lower is better!



Legend: Lookup, Tied, Conv., Bilinear, Deep residual, Adaptive, GroC (ours)

penn (monolingual): 90.8, 86.2, 101.6, 83.7, 80.5, 79.3, 69.5 — ~1M training tokens

wikitext2 (monolingual): 108.3, 97.3, 116.6, 95.9, 94.7, 90.7, 82.5 — ~2M training tokens
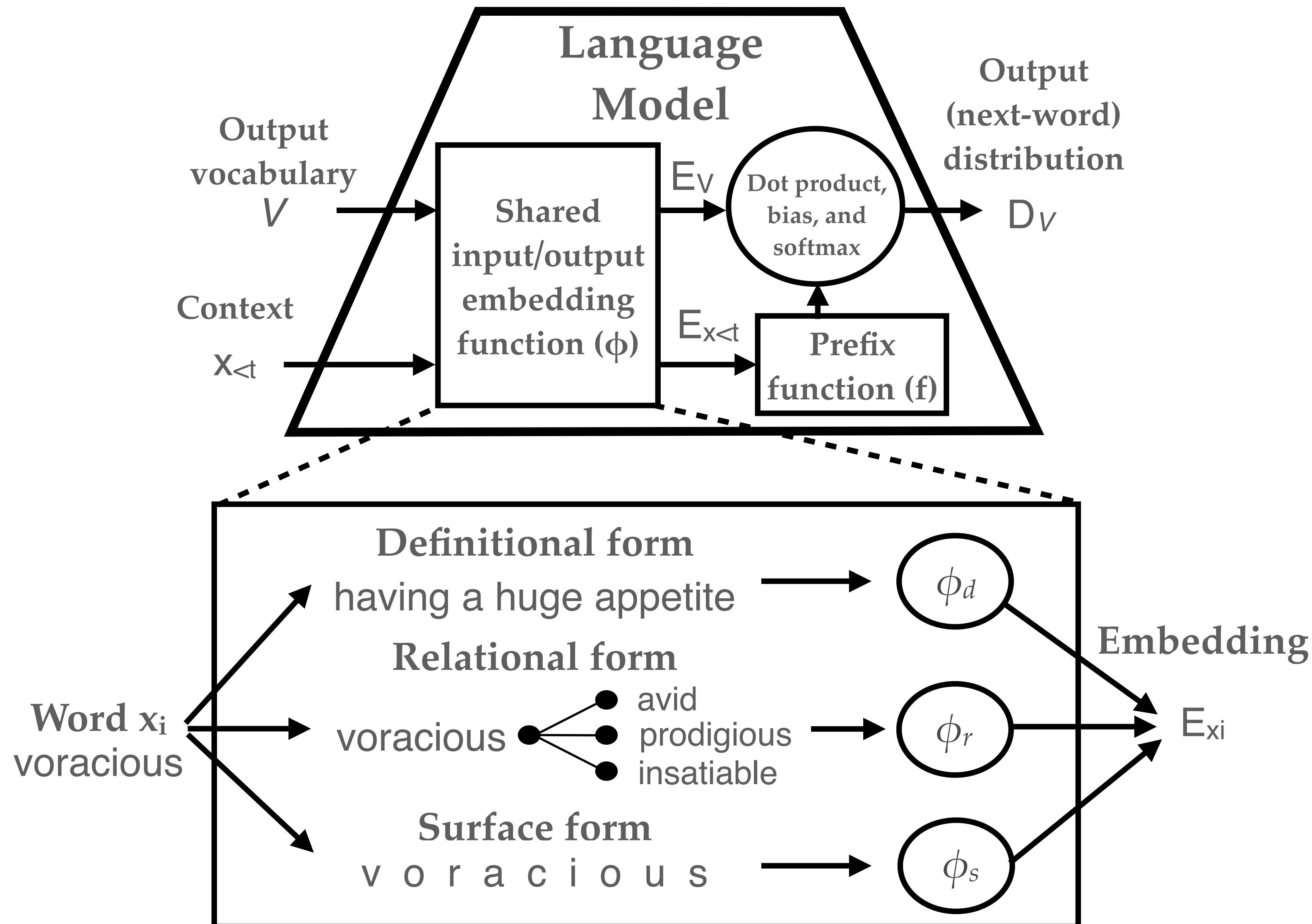
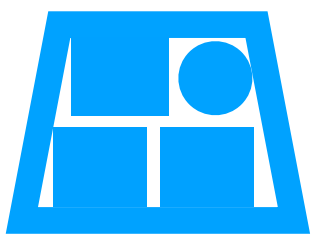# zero-resource cross-domain adaptation

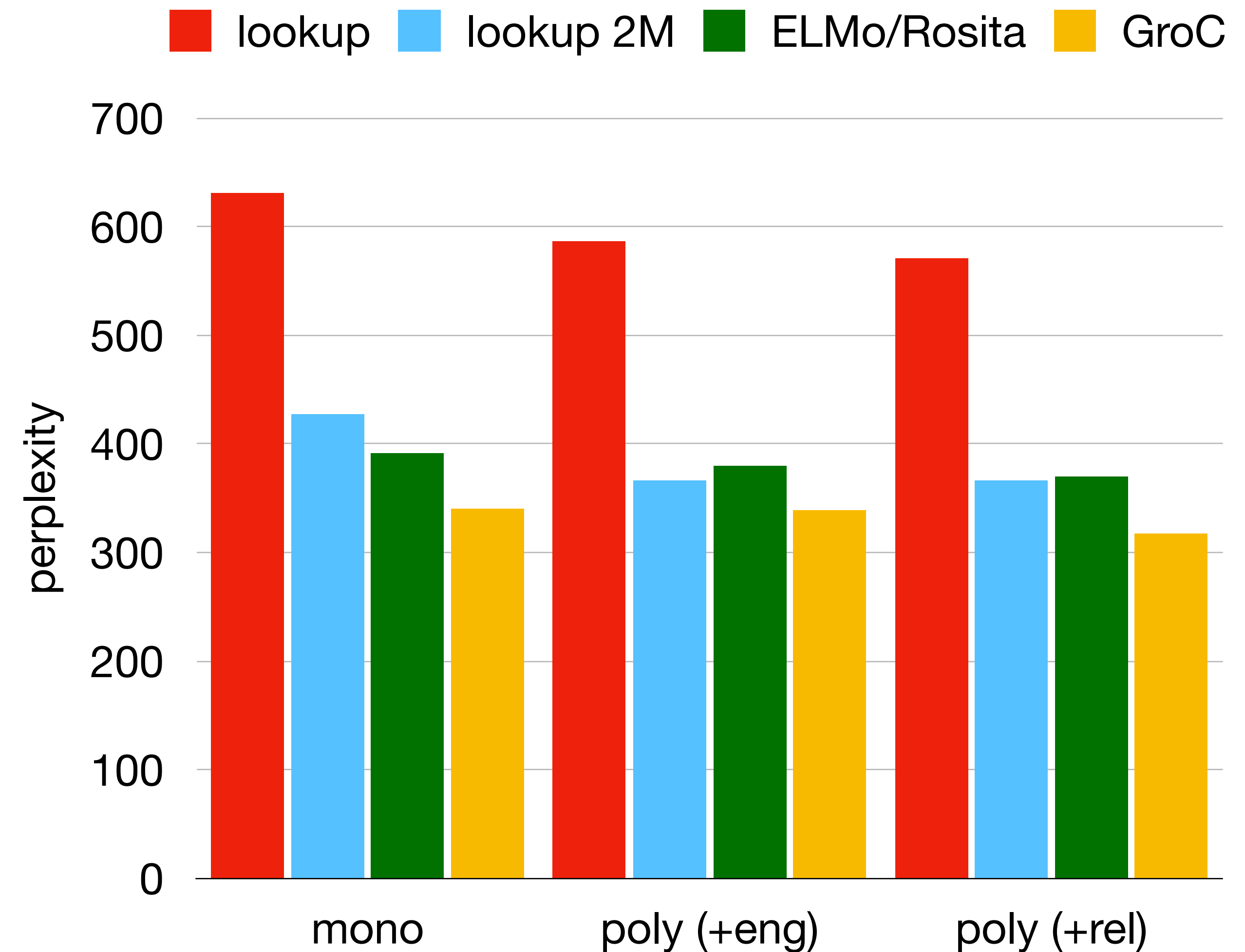# polyglot vocab-independence: multilingual GroC



- share all parameters

- use a multilingual lexicon for relational and definitional features (Open Multilingual WordNet)— this only covered some languages

# multilingual GroC: results

- lookup vs ELMo/Rosita vs GroC

- monolingual / +English / +related

- multilingual GroC is reliably the best method across 9 target languages

- related languages help more

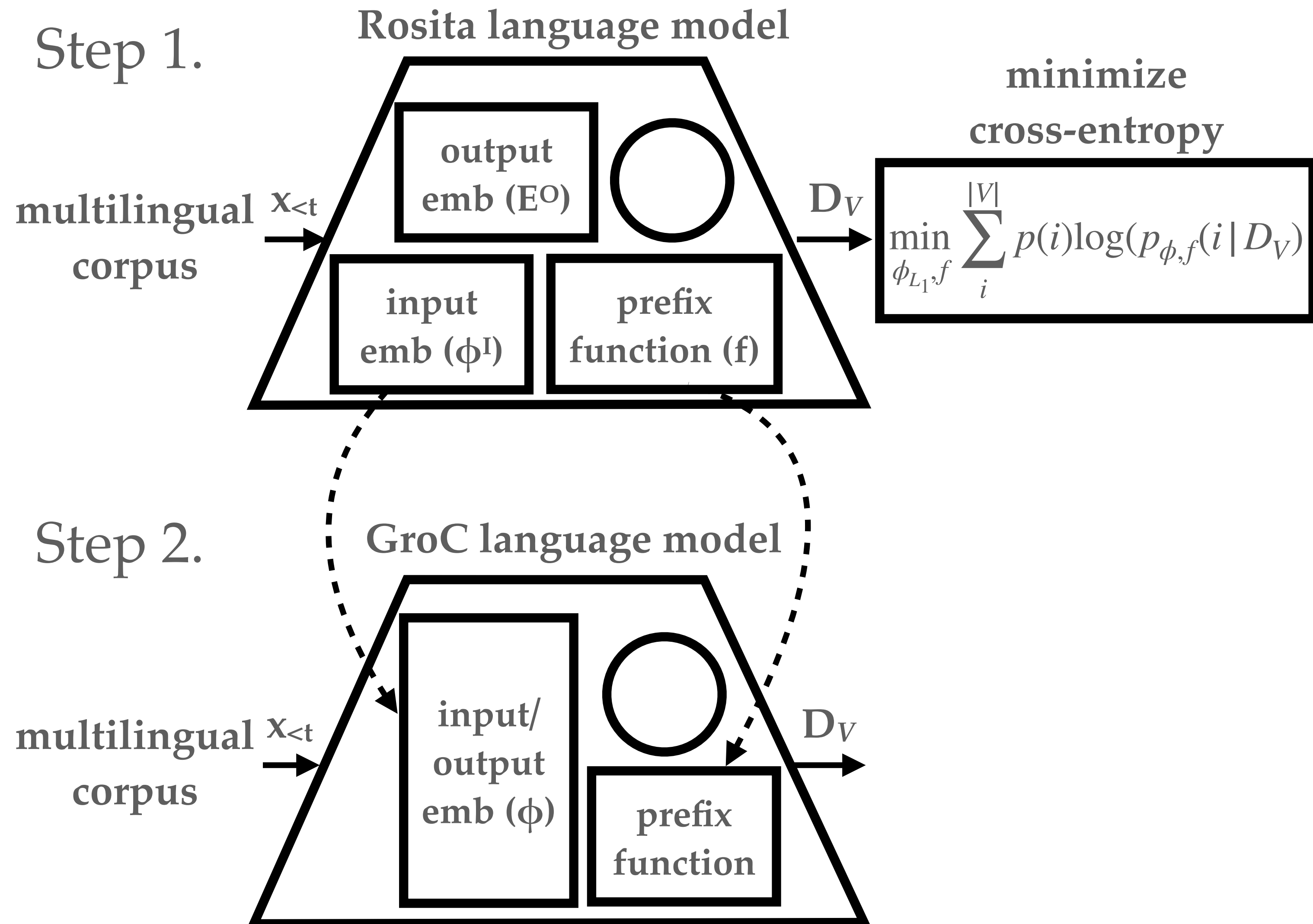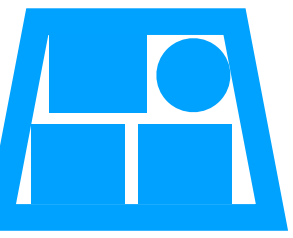- still outperforms lookup with 0.5x the data!
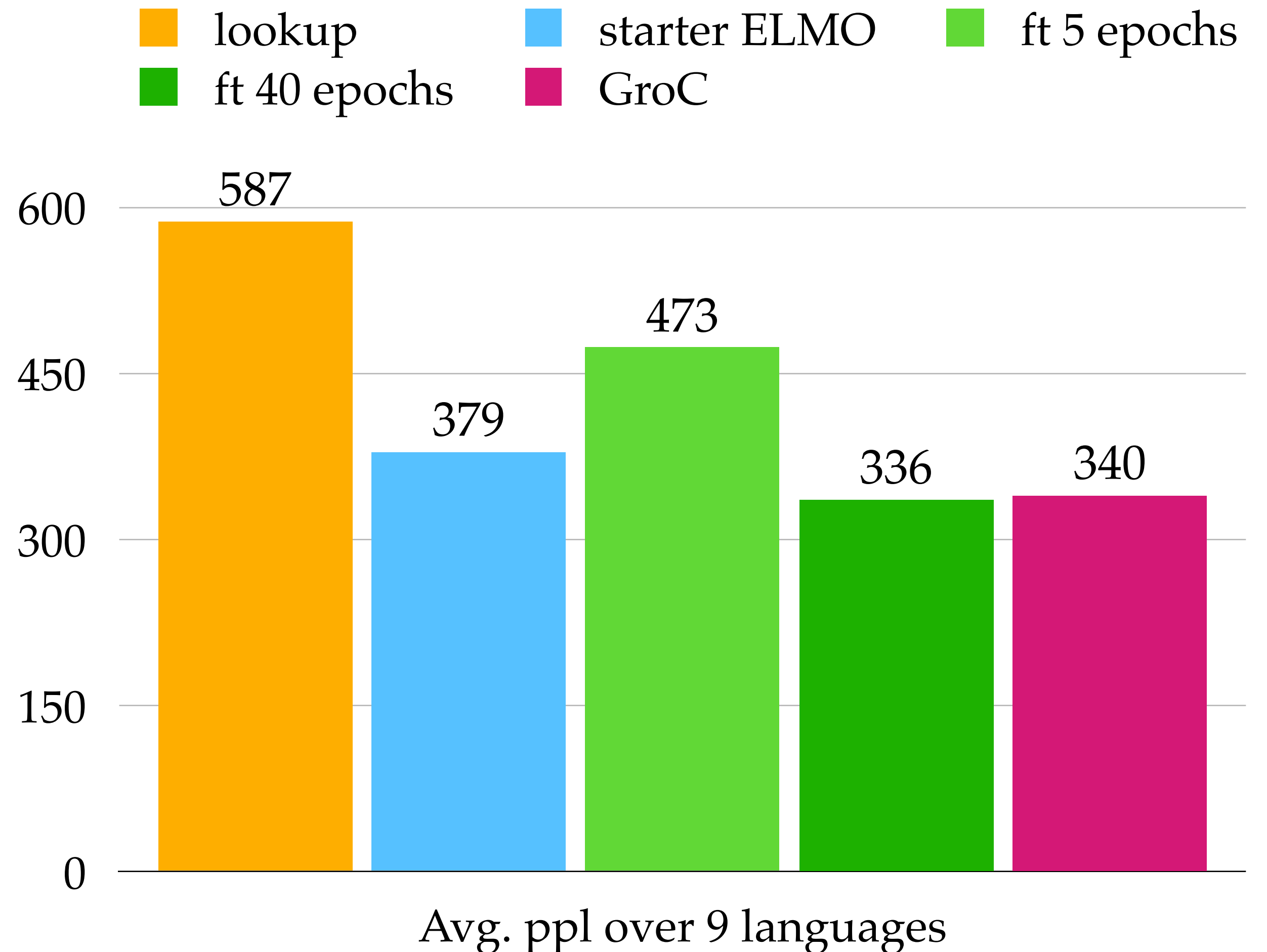
# initializing compositional outputs

- ELMo/Rosita trains faster than GroC

- train an Rosita-like LM

- turn the compositional *input* embedding into a *shared* input-output embedding

- produce a GroC model cheaply—but need to finetune

**Step 1.**

**Rosita language model**

**minimize**

**cross-entropy**

**multilingual corpus** $\mathbf{x}_{<t}$

output emb (E$^O$)

input emb ($\phi^I$)

prefix function (f)

$\mathbf{D}_V$

$$\min_{\phi_{L_1}, f} \sum_i^{|V|} p(i) \log(p_{\phi, f}(i \mid D_V))$$

**Step 2.**

**GroC language model**

**multilingual corpus** $\mathbf{x}_{<t}$

input/ output emb ($\phi$)

prefix function

$\mathbf{D}_V$

# initializing compositional outputs: results

- needs finetuning, but not much

- can beat GroC-from-scratch with less total training time!

- holds promise for application of GroC-like representations to large-scale language models!

**Legend:** lookup, starter ELMO, ft 5 epochs, ft 40 epochs, GroC

587 · 379 · 473 · 336 · 340

Avg. ppl over 9 languages

# conclusion

• crosslingual sharing works

• low-resource NLP is hard, but tractable—*if* we use sharing

• related languages and vocab-independence are useful

# thank you!

Collaborators: