

Grounded Compositional Outputs for Adaptive Language Modeling

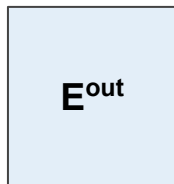
Nikolaos Pappas, Phoebe Mulcaire, Noah A. Smith

EMNLP 2020



Neural language models

Output embedding

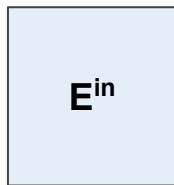


Classifies tokens

Prefix encoder



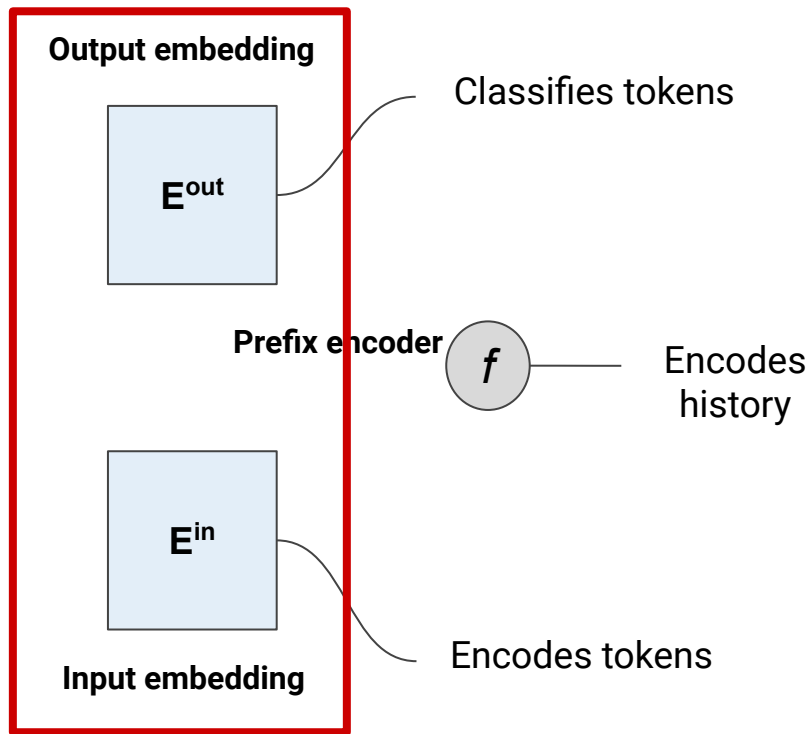
Encodes history



Input embedding

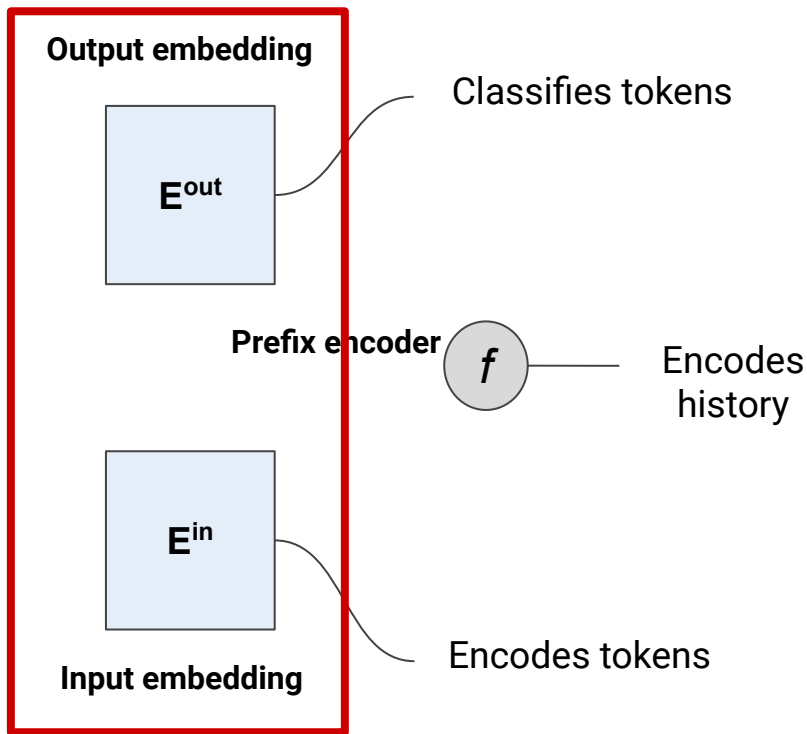
Encodes tokens

Neural language models: Limitations



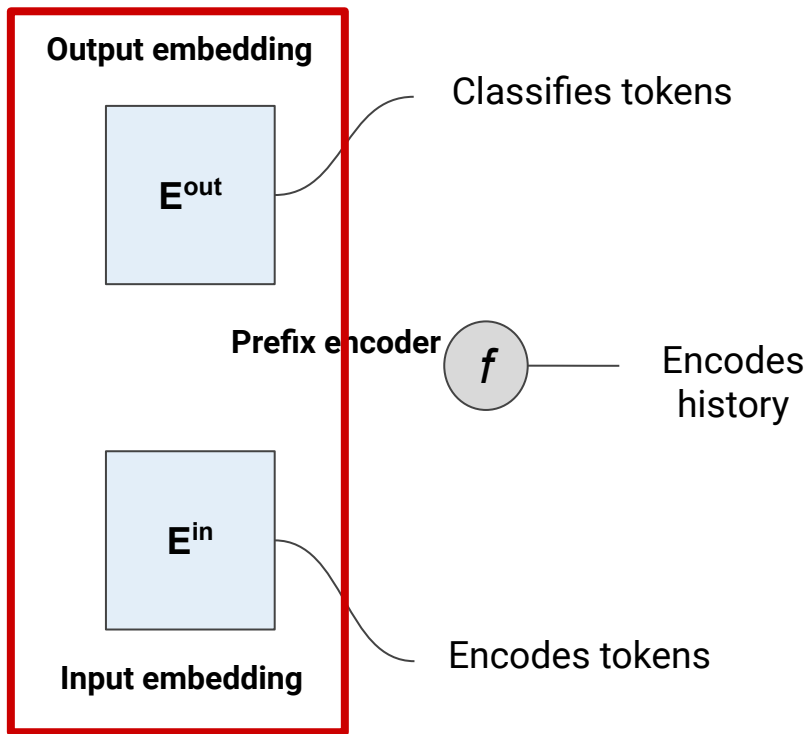
- Parameterization depends on the vocabulary

Neural language models: Limitations



- Parameterization depends on the vocabulary
- Handle rare or new words poorly

Neural language models: Limitations



- Parameterization depends on the vocabulary
- Handle rare or new words poorly
- Cannot be gracefully modified once trained

Motivation

- Can we use lexicons to better generalize to rare words?
- How to decouple training and testing vocabularies?
- What output form is most suitable for domain adaptation?



Previous work on rare words

Subword tokenization

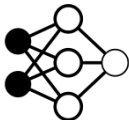
- Character-level models ([Cherry et al., 2018](#); [Al-Rfou et al., 2019](#))
 - ❌ Costly prefix encoders and training
- Data-driven vocabulary selection ([Sennrich et al., 2016](#); [Radford et al., 2018](#))
 - ❌ Linguistically simplistic
 - ❌ Rely on lookup tables

Interpolation with a neural cache

- Local or unbounded neural cache ([Graves et al., 2017a,b](#))
 - ✓ Low-cost adaptation to rare/new words

$$p(x_t|h_{1:t}, x_{1:t}) = (1 - \lambda)p_{vocab}(x_t|h_t) + \lambda p_{cache}(x_t|h_{1:t}, x_{1:t})$$

Neural model



Local neural cache

$$\propto \sum_{i=1}^{t-1} \mathbb{1}_{\{w=x_{i+1}\}} \exp(\theta h_t^\top h_i)$$

Sharing across words

Sharing across words

Lookup tables
(Zaremba et al., 2014;
Press & Wolf, 2017)

Input



$$E^{out} \neq E^{in}$$

Output



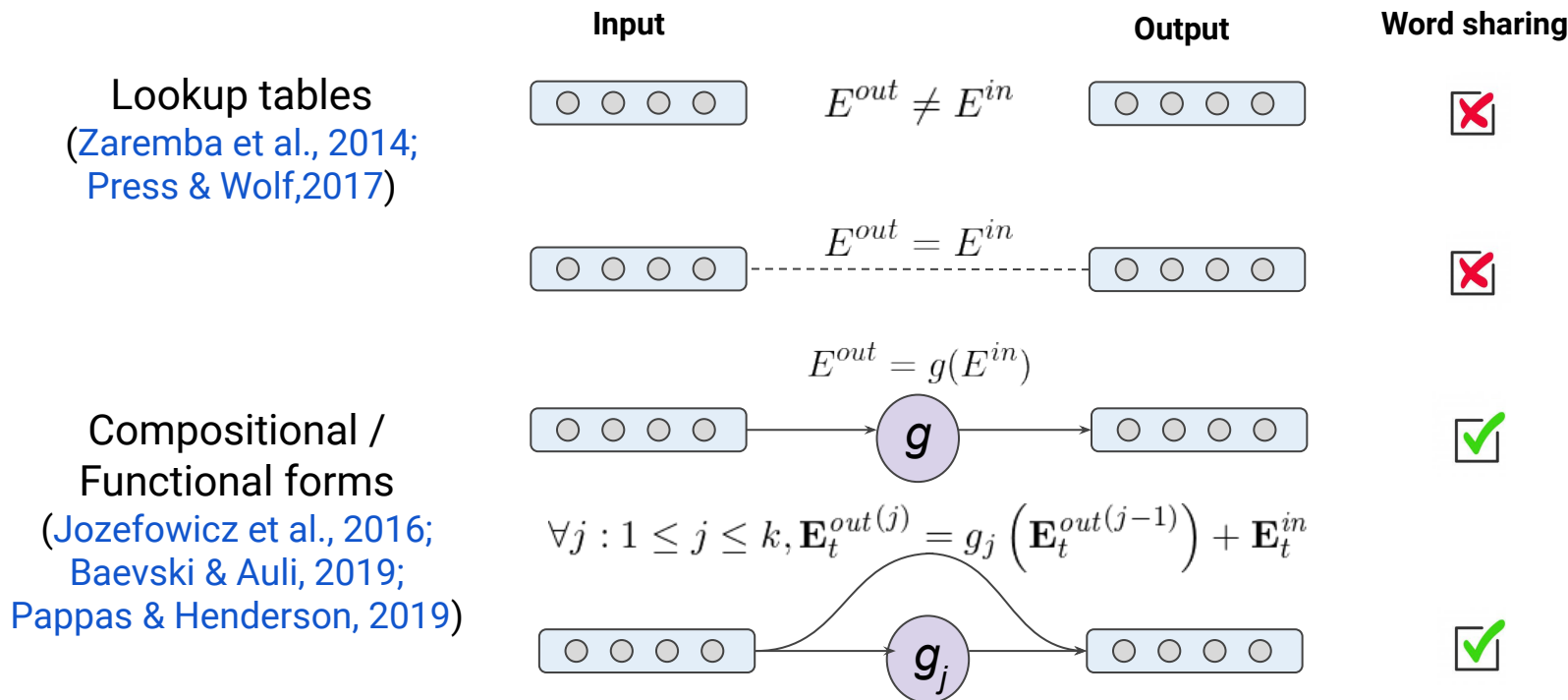
Word sharing



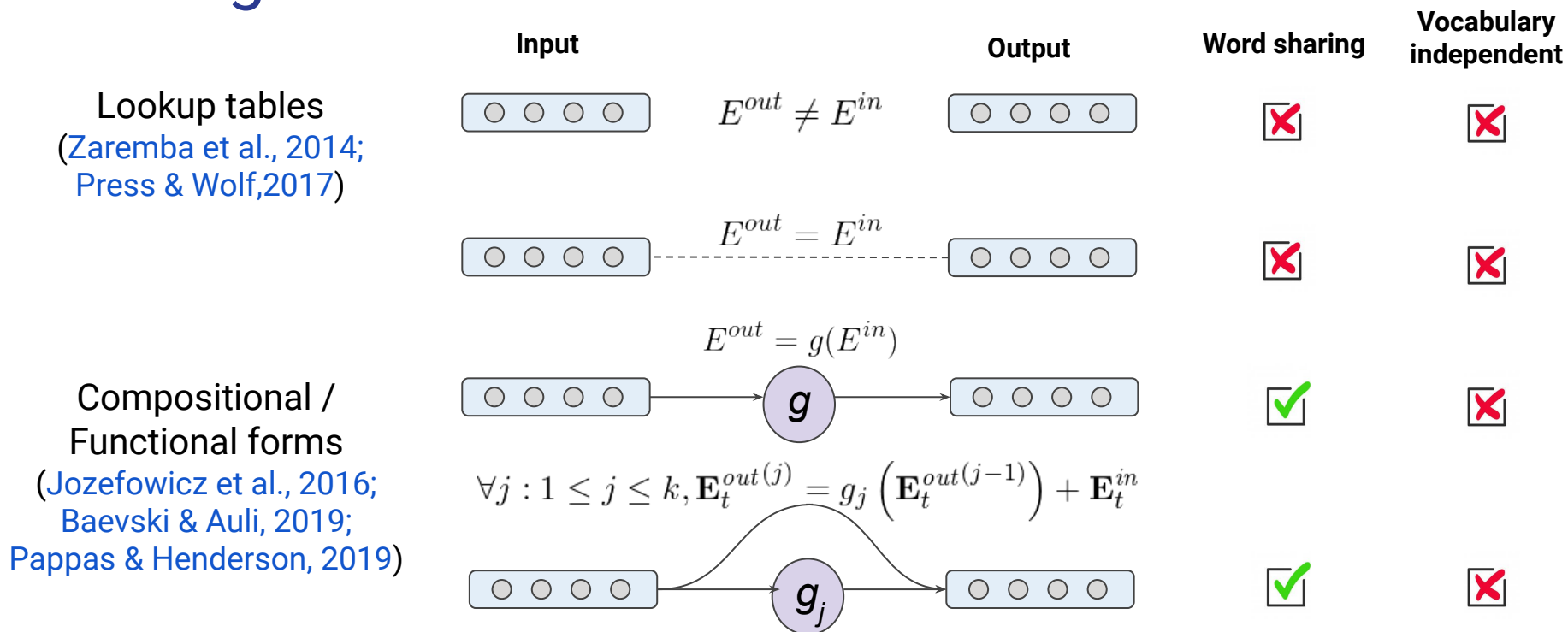
$$E^{out} = E^{in}$$



Sharing across words

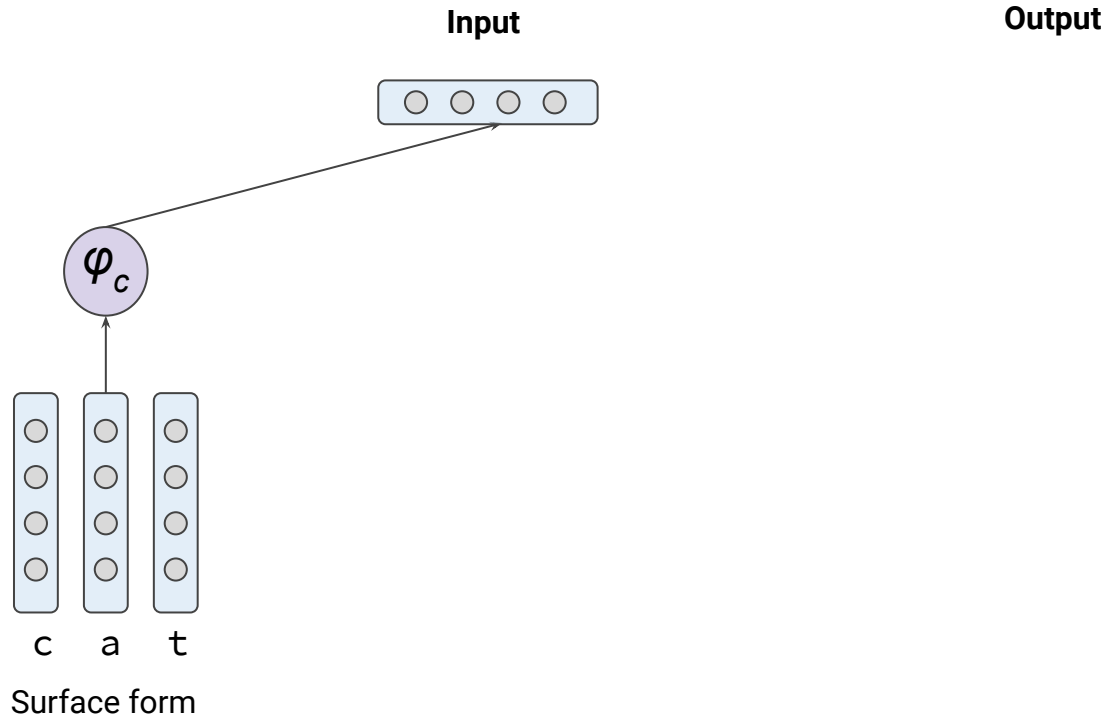


Sharing across words

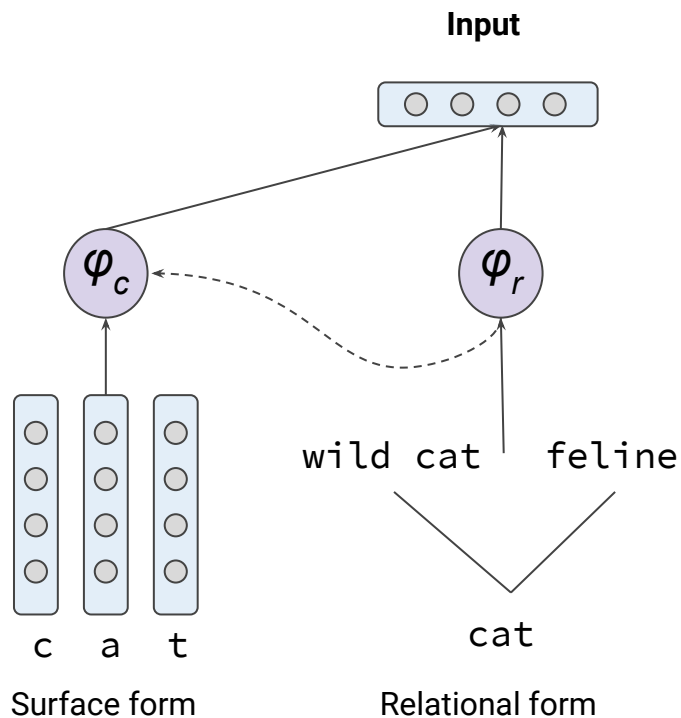


Our method – GroC

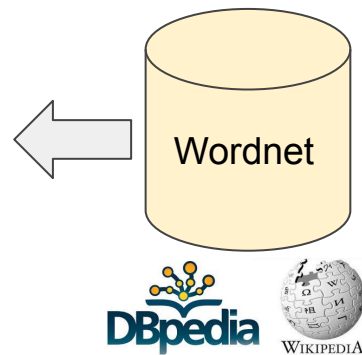
Grounded compositional forms



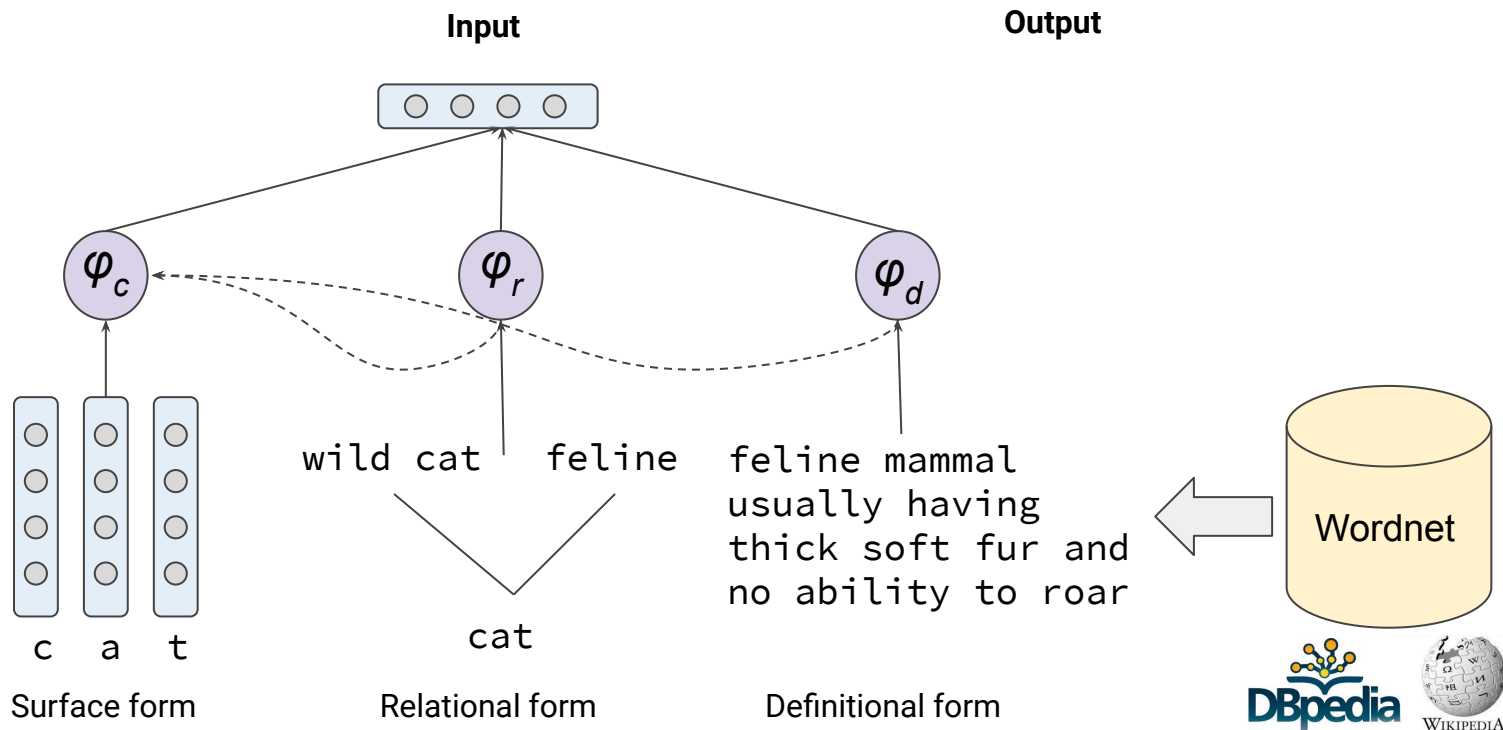
Grounded compositional forms



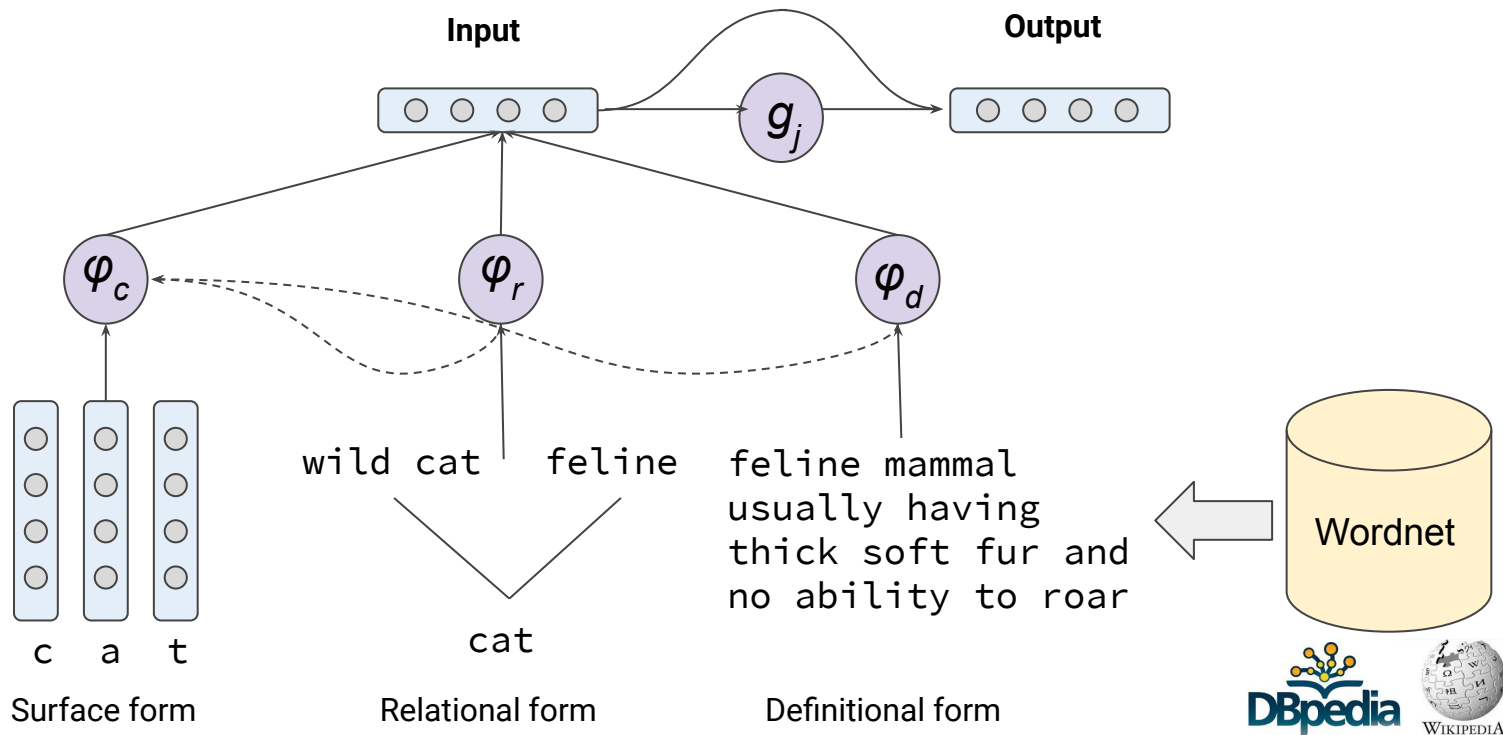
Output



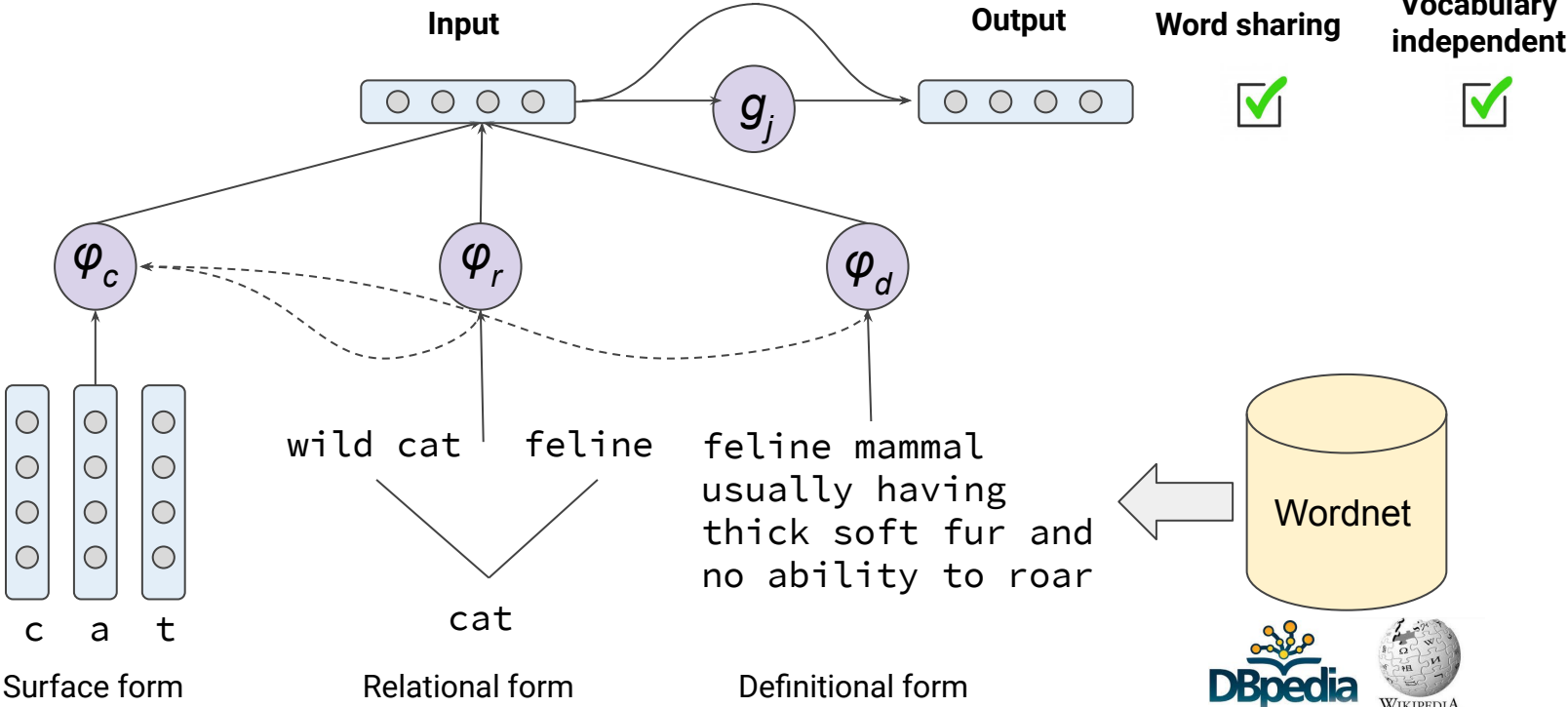
Grounded compositional forms



Grounded compositional forms



Grounded compositional forms



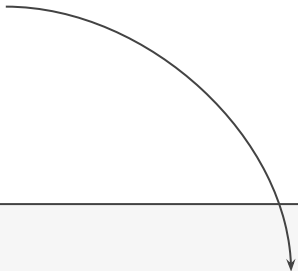
Adapting to any vocabulary

$$p(X_t = x_t | \mathbf{h}_{t-1}) \propto \exp(\mathbf{E}^{out} \mathbf{h}_{t-1} + \mathbf{b})$$

Adapting to any vocabulary

- We first represent the vocabulary with GroC

$$\mathbf{E}^{out} = \text{GroC}(\mathcal{V}^*)$$

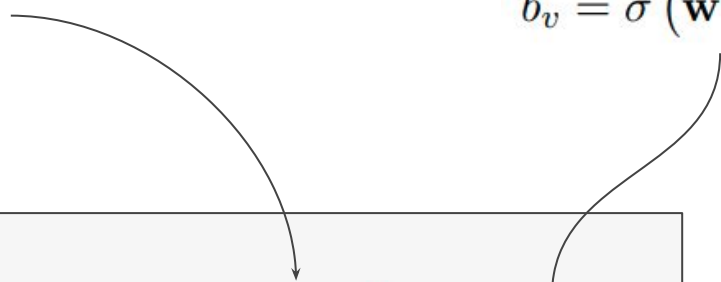

$$p(X_t = x_t | \mathbf{h}_{t-1}) \propto \exp(\mathbf{E}^{out} \mathbf{h}_{t-1} + \mathbf{b})$$

Adapting to any vocabulary

- We first represent the vocabulary with GroC
- Then we estimate the bias for each word u

$$\mathbf{E}^{out} = \text{GroC}(\mathcal{V}^*)$$

$$b_v = \sigma(\mathbf{w} \cdot \mathbf{e}_v^{out} + a)$$


$$p(X_t = x_t | \mathbf{h}_{t-1}) \propto \exp(\mathbf{E}^{out} \mathbf{h}_{t-1} + \mathbf{b})$$

GroC summary

- Creates a compact representation of any vocabulary
- Grounds the language model predictions to prior knowledge
- Enables the decoupling of training and test vocabularies

Experiments

Conventional language modeling

How GroC compares to previous output embedding methods?

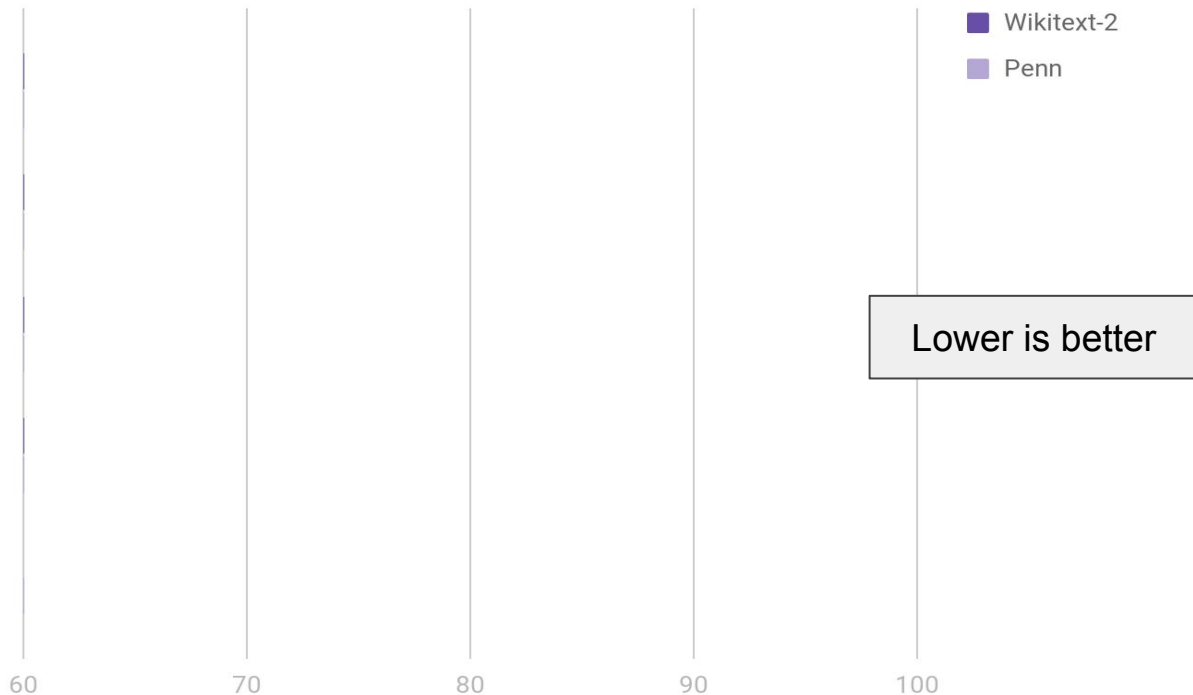
Conventional language modeling

Perplexity

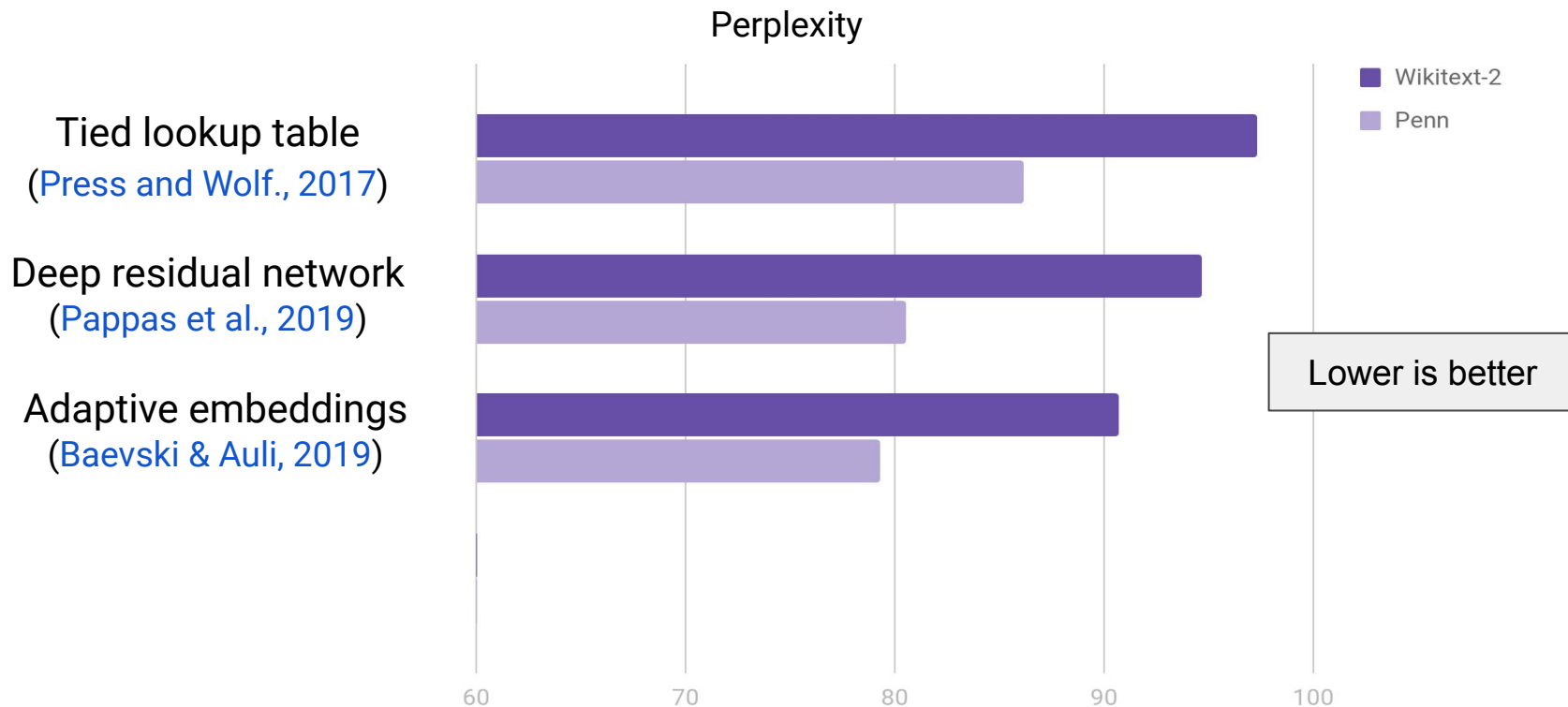
Tied lookup table
([Press and Wolf., 2017](#))

Deep residual network
([Pappas et al., 2019](#))

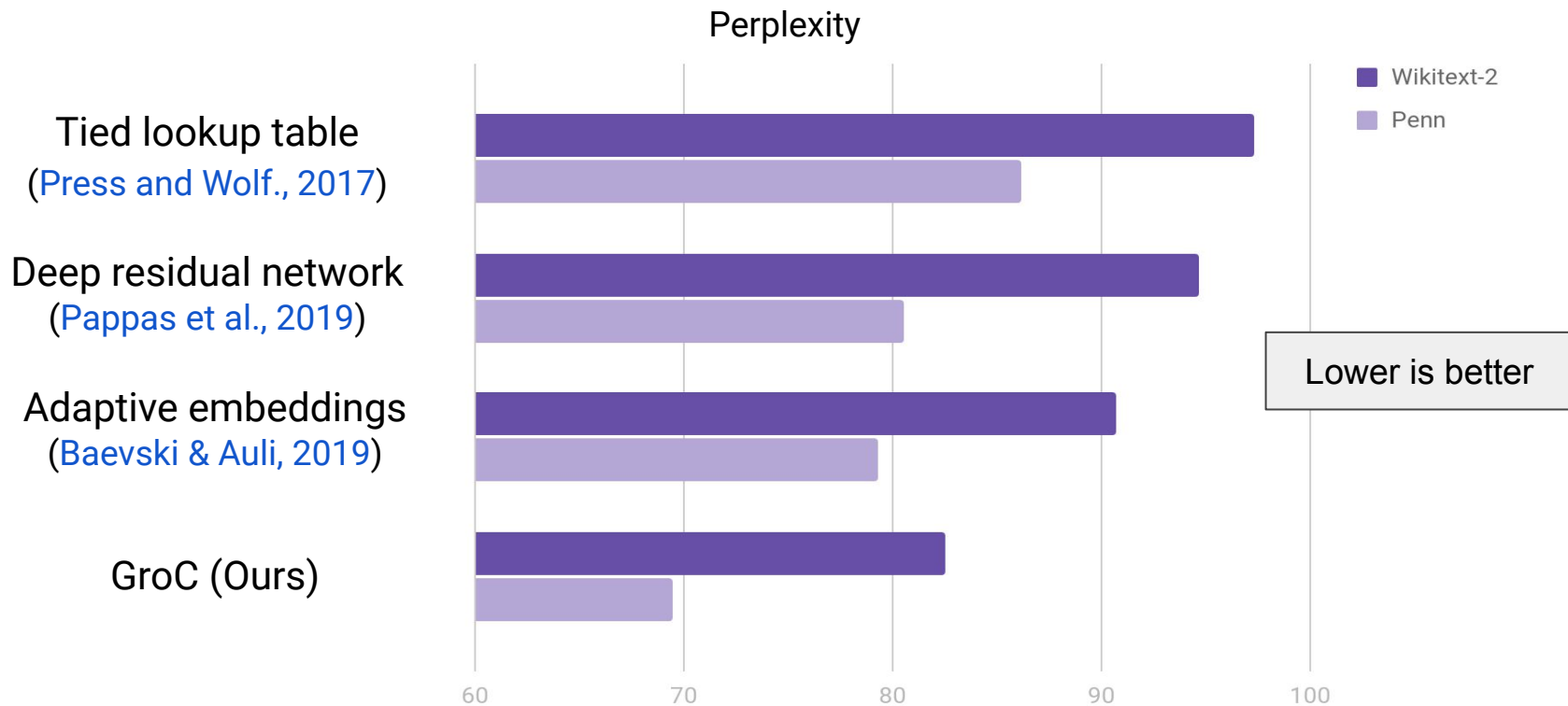
Adaptive embeddings
([Baevski & Auli, 2019](#))



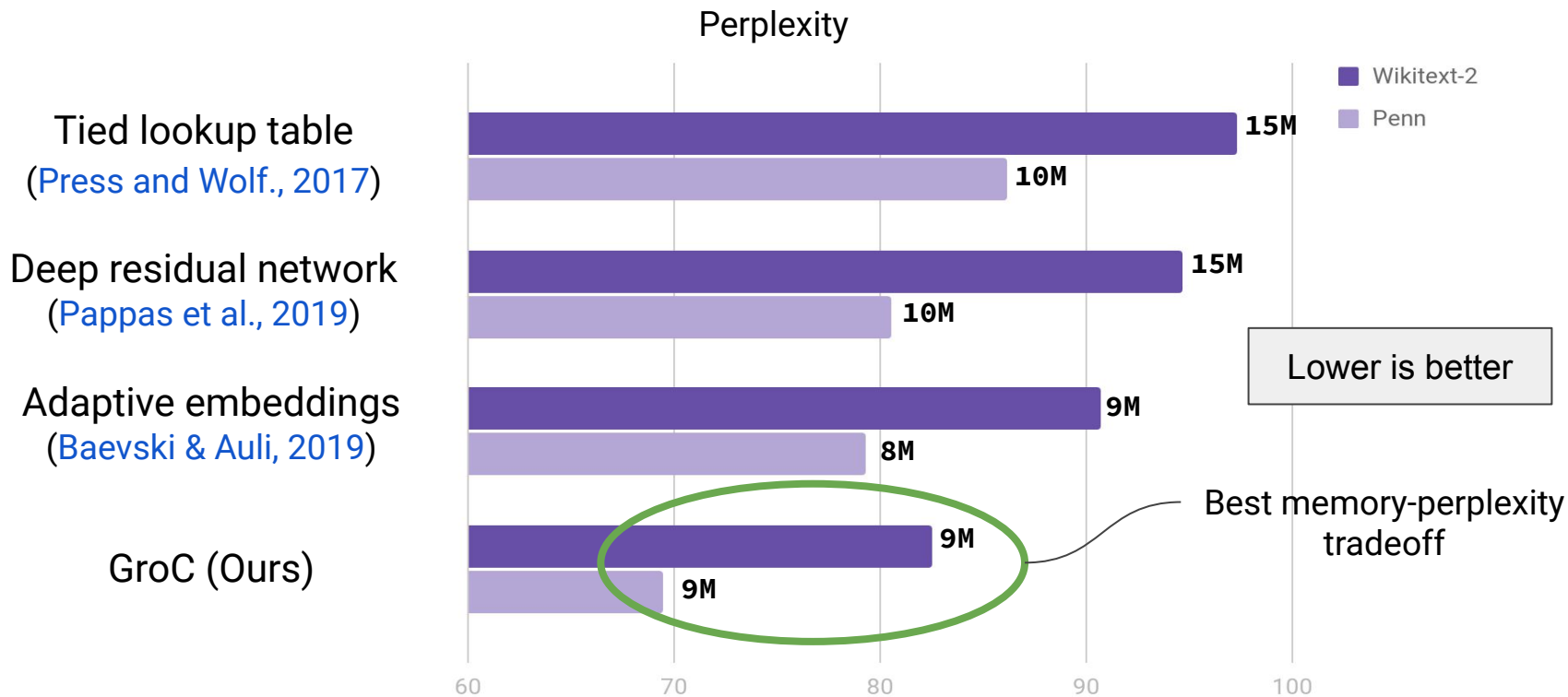
Conventional language modeling



Conventional language modeling



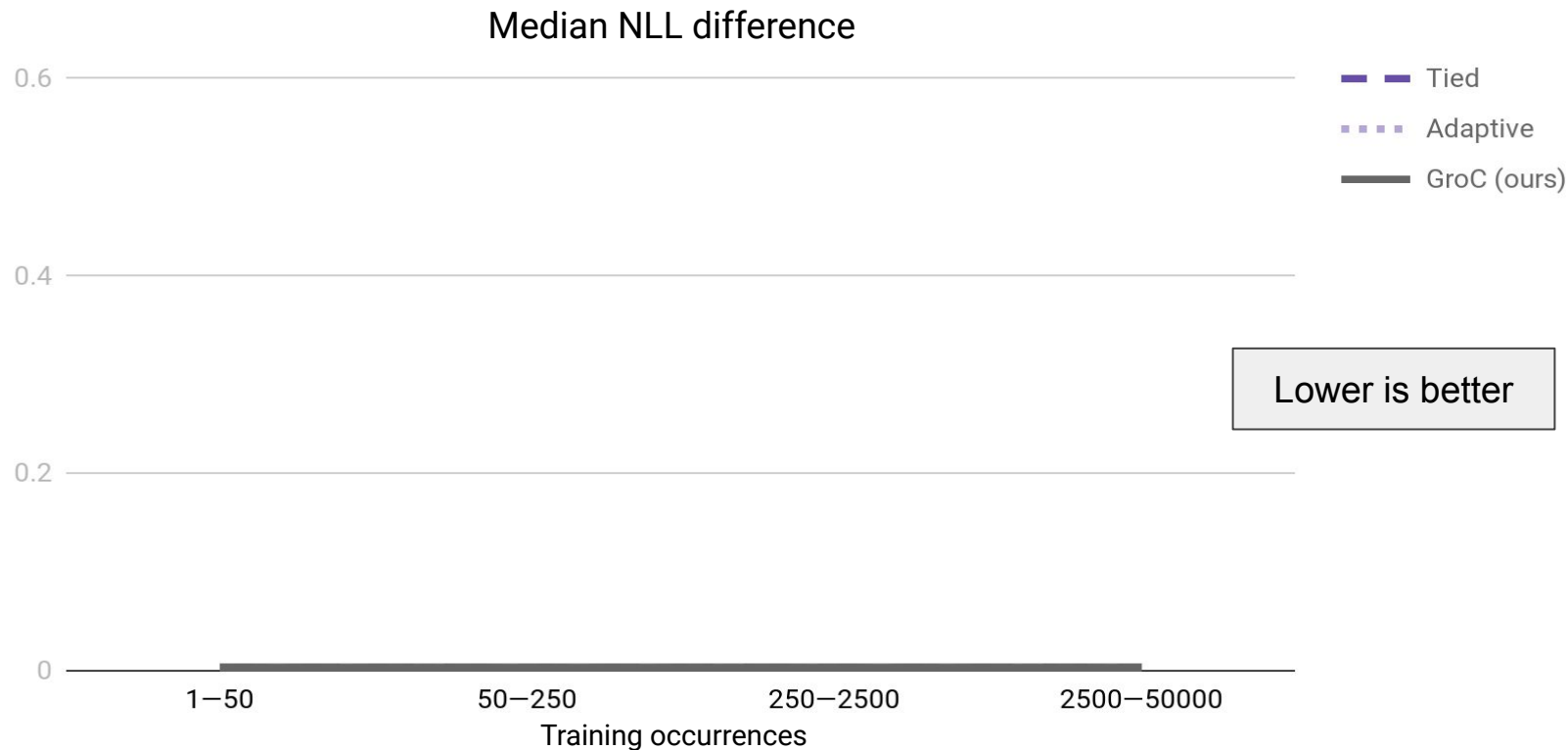
Conventional language modeling



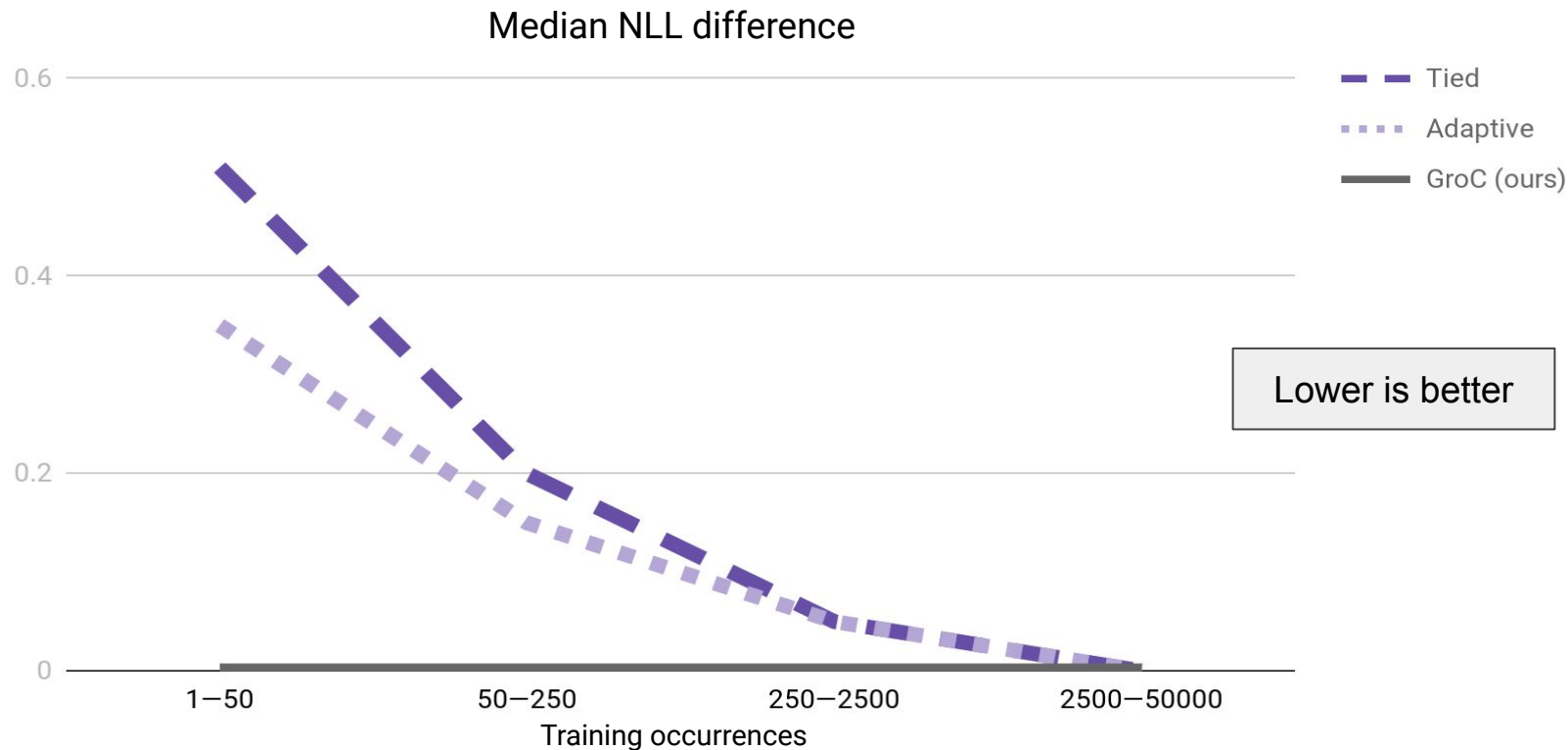
Conventional language modeling

Where does the improvement come from?

Break down by frequency

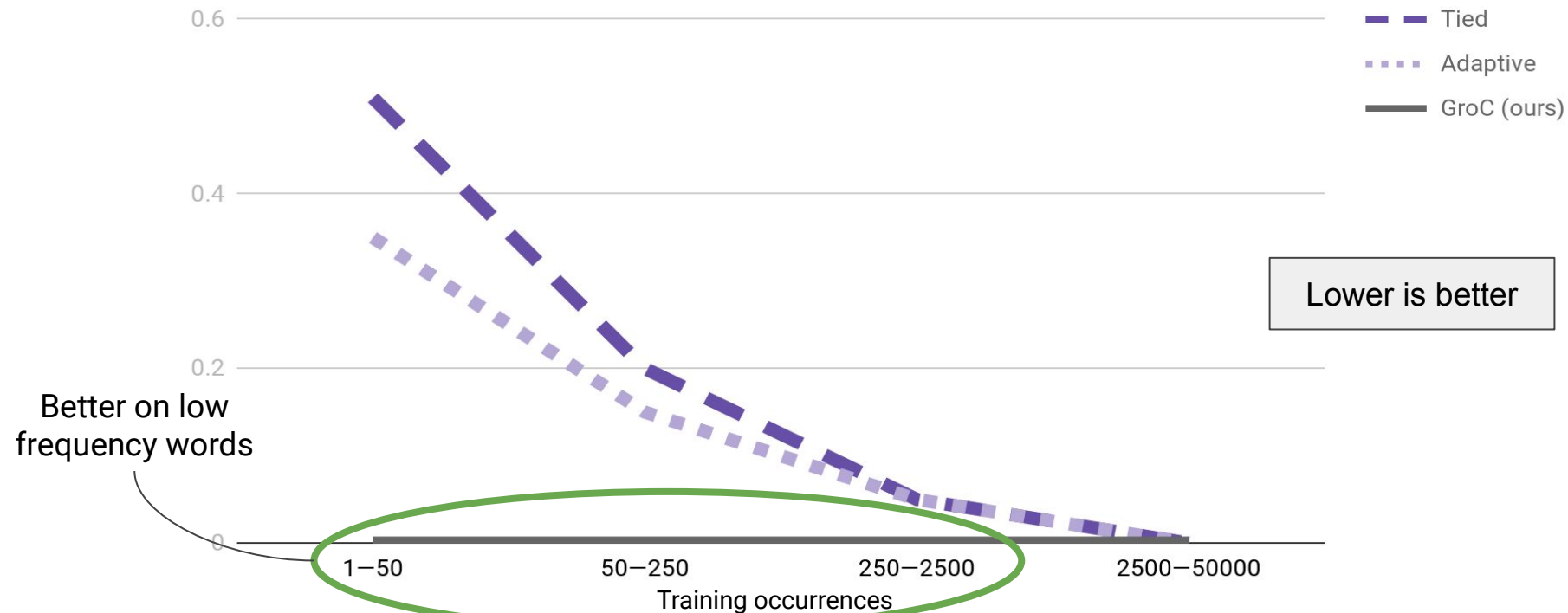


Break down by frequency



Break down by frequency

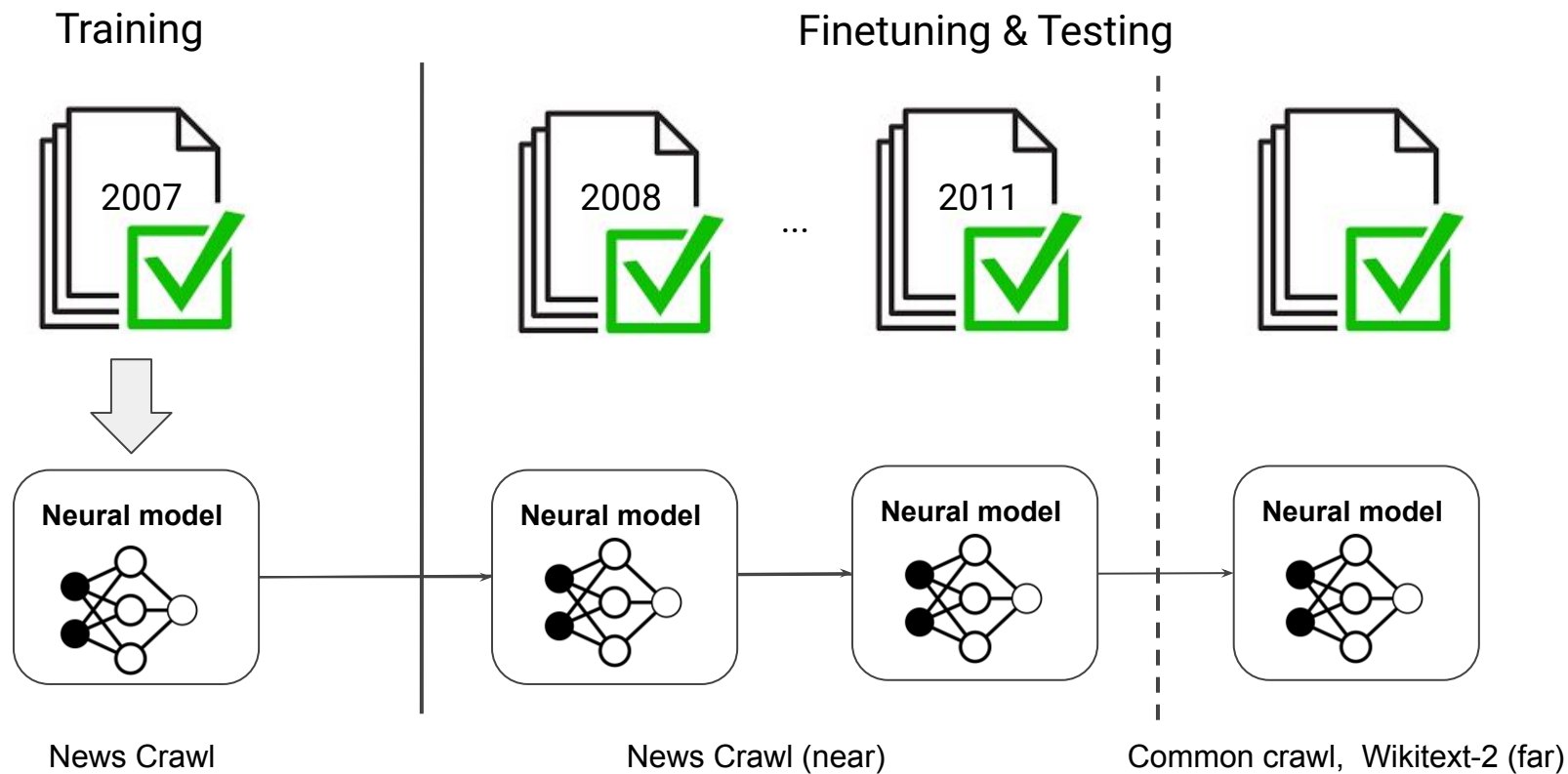
Median NLL difference



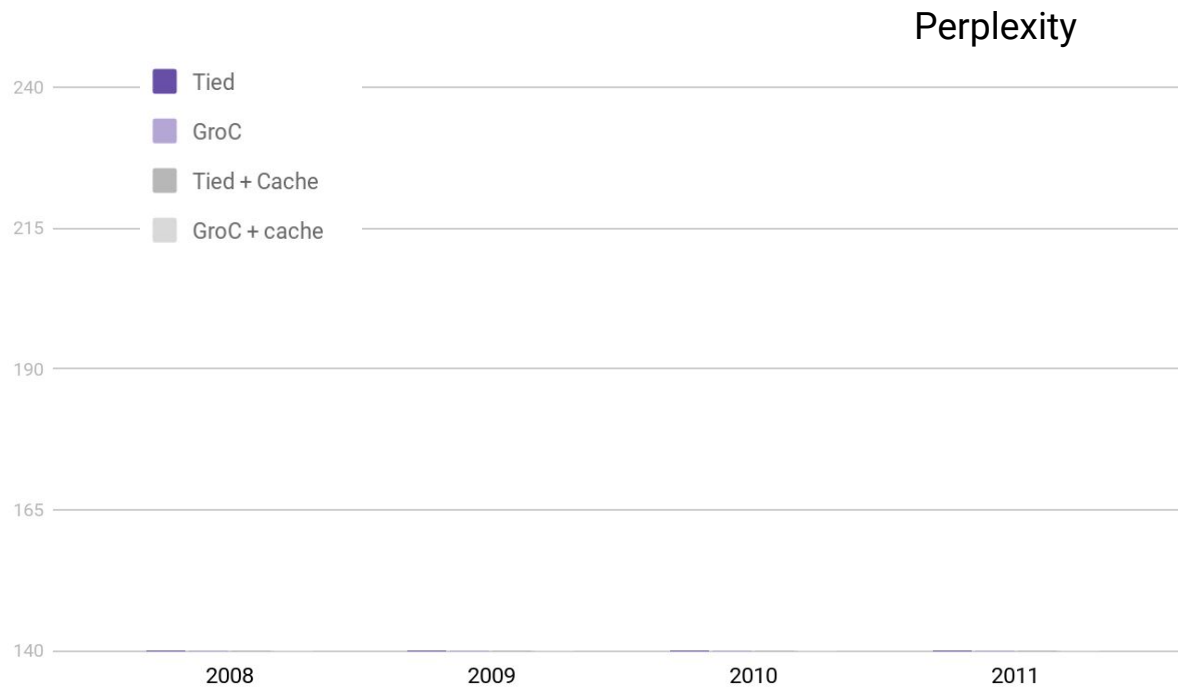
Cross-domain modeling: Zero resource

Does GroC generalize on zero resource adaptation settings?

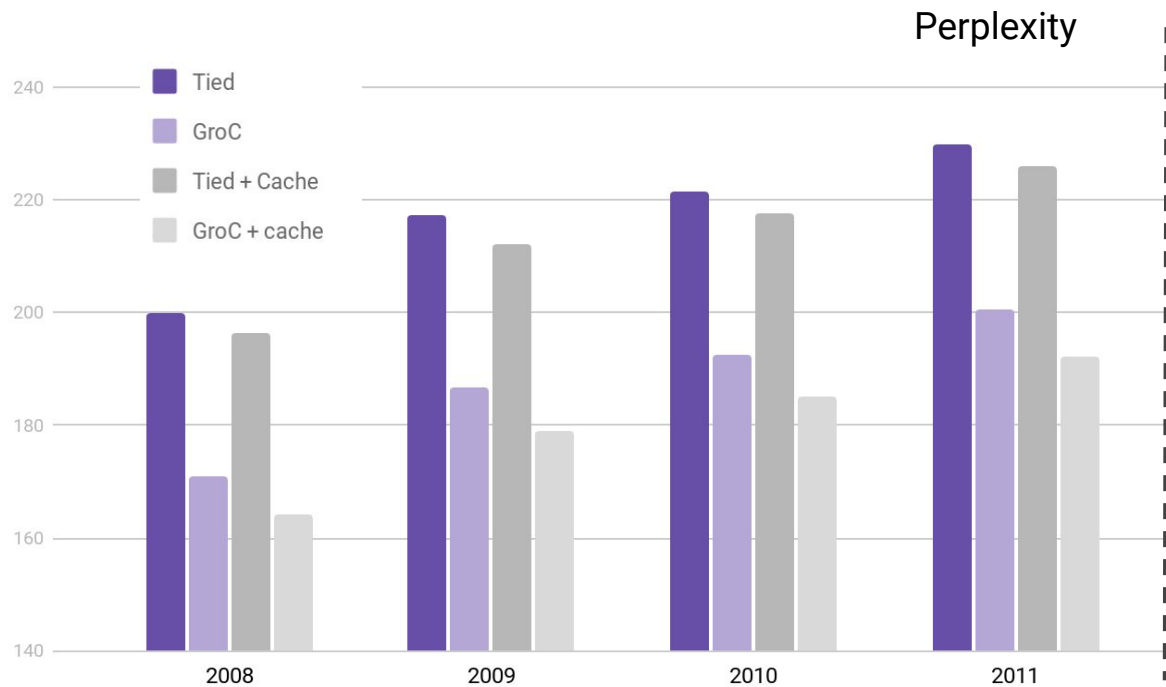
Cross-domain modeling: Zero resource



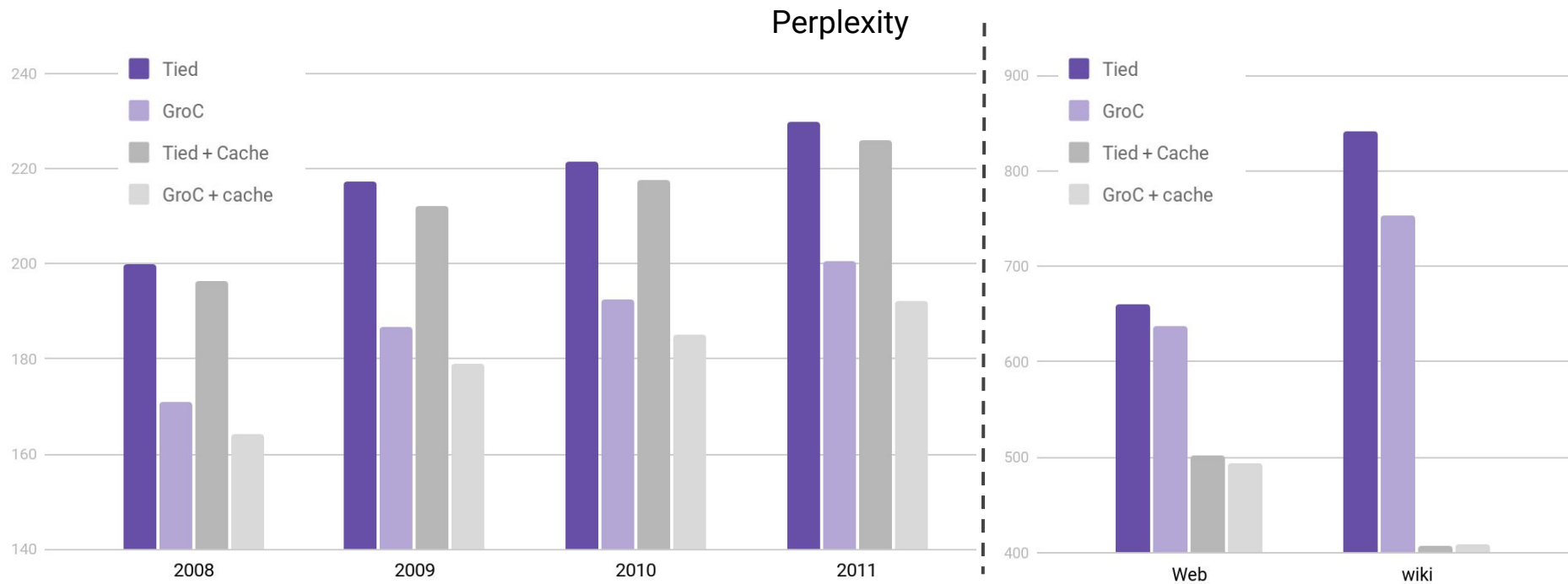
Cross-domain modeling: Zero resource



Cross-domain modeling: Zero resource



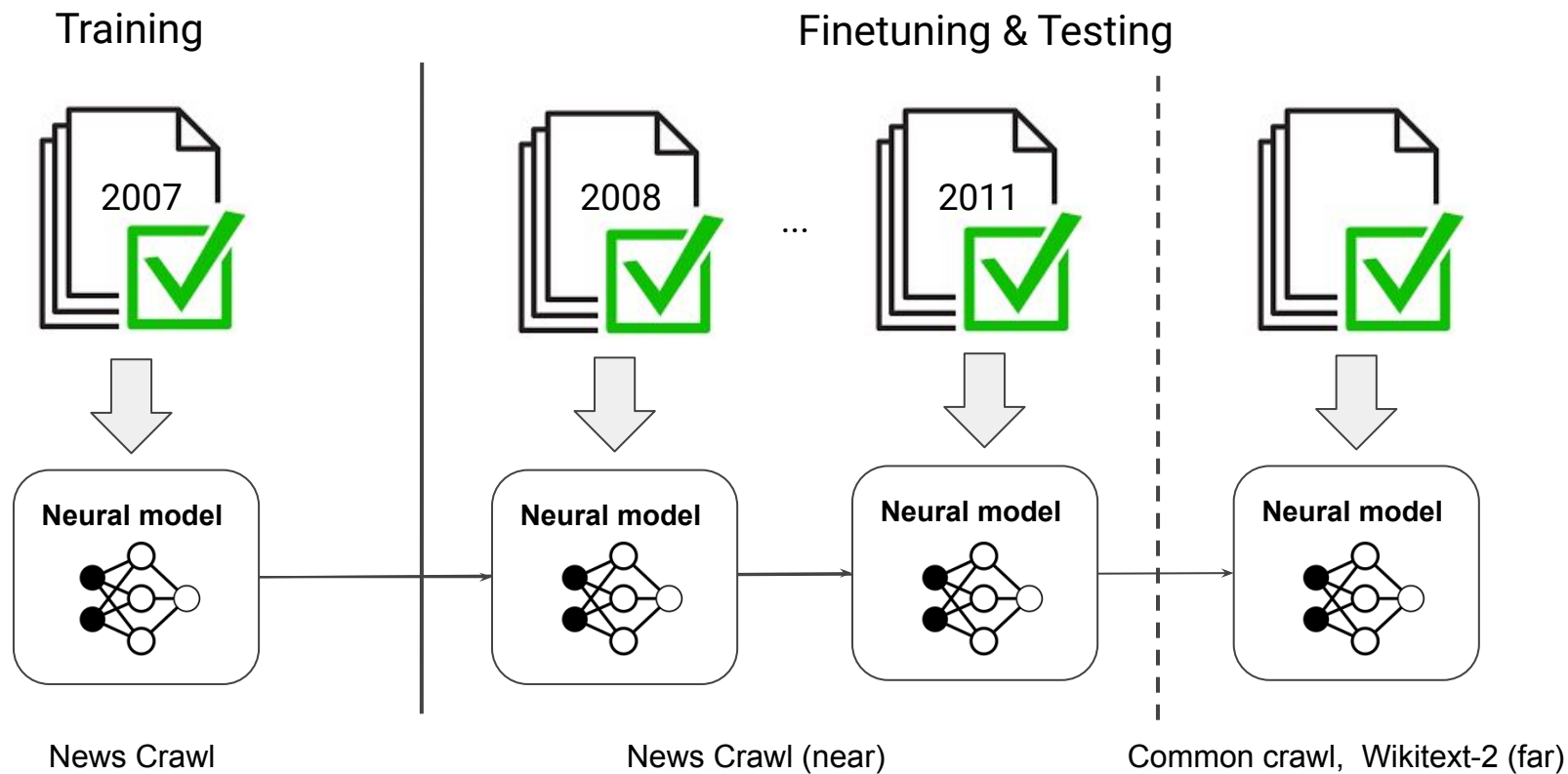
Cross-domain modeling: Zero resource



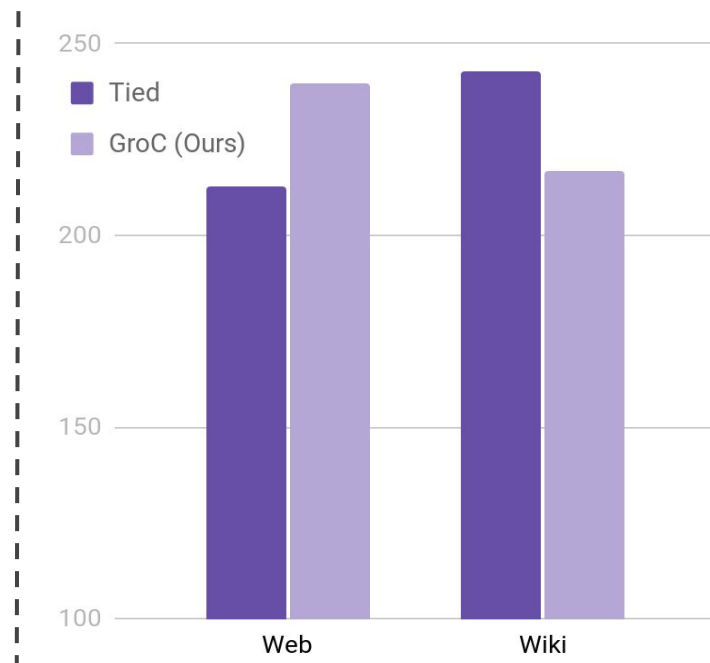
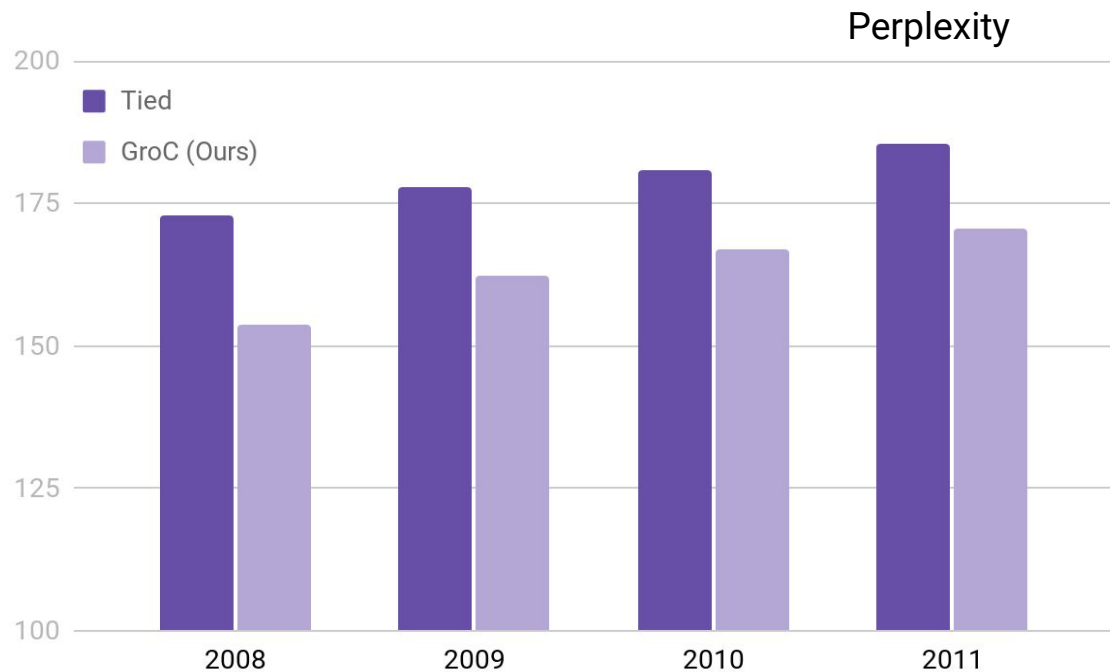
Cross-domain modeling: Low resource

Does GroC help on low resource adaptation settings?

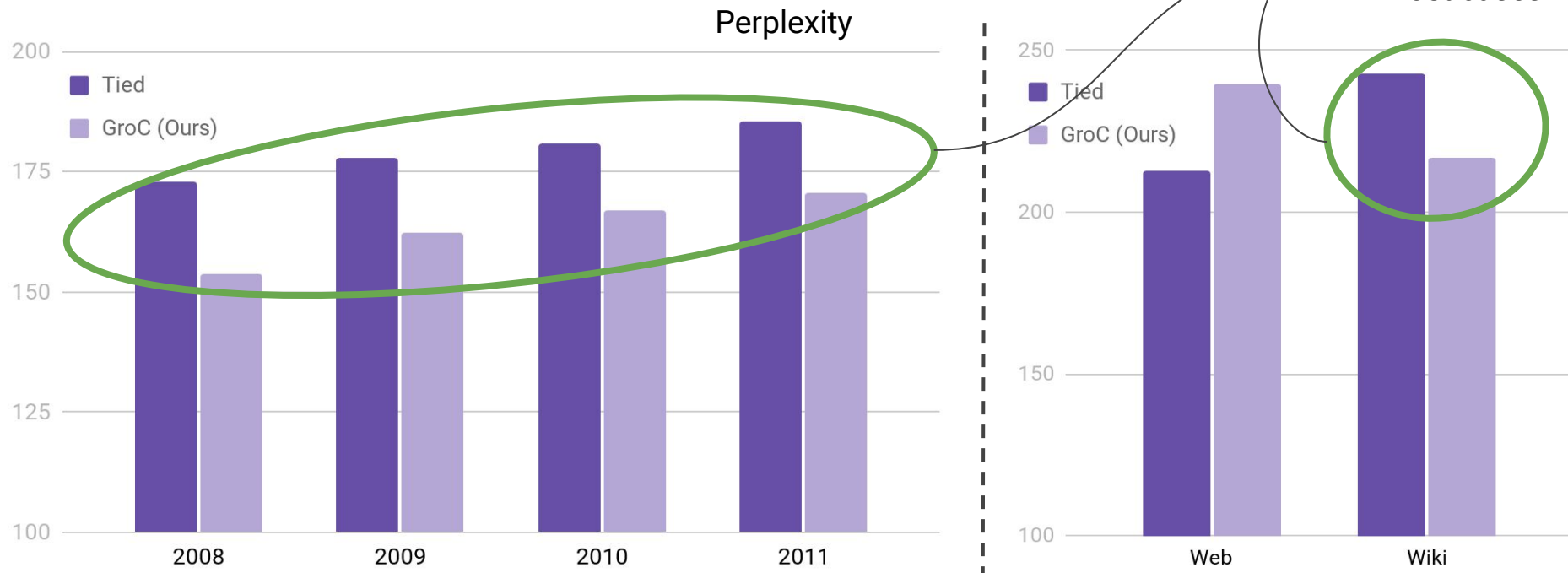
Cross-domain modeling: Low resource



Cross-domain modeling: Low resource



Cross-domain modeling: Low resource



Conclusion

Grounded compositional outputs for language models

- Outperform previous methods on conventional settings
- Achieve low perplexity on rare words
- Generalize well to previously unseen domains

Thank you



<https://github.com/Noahs-ARK/groc>