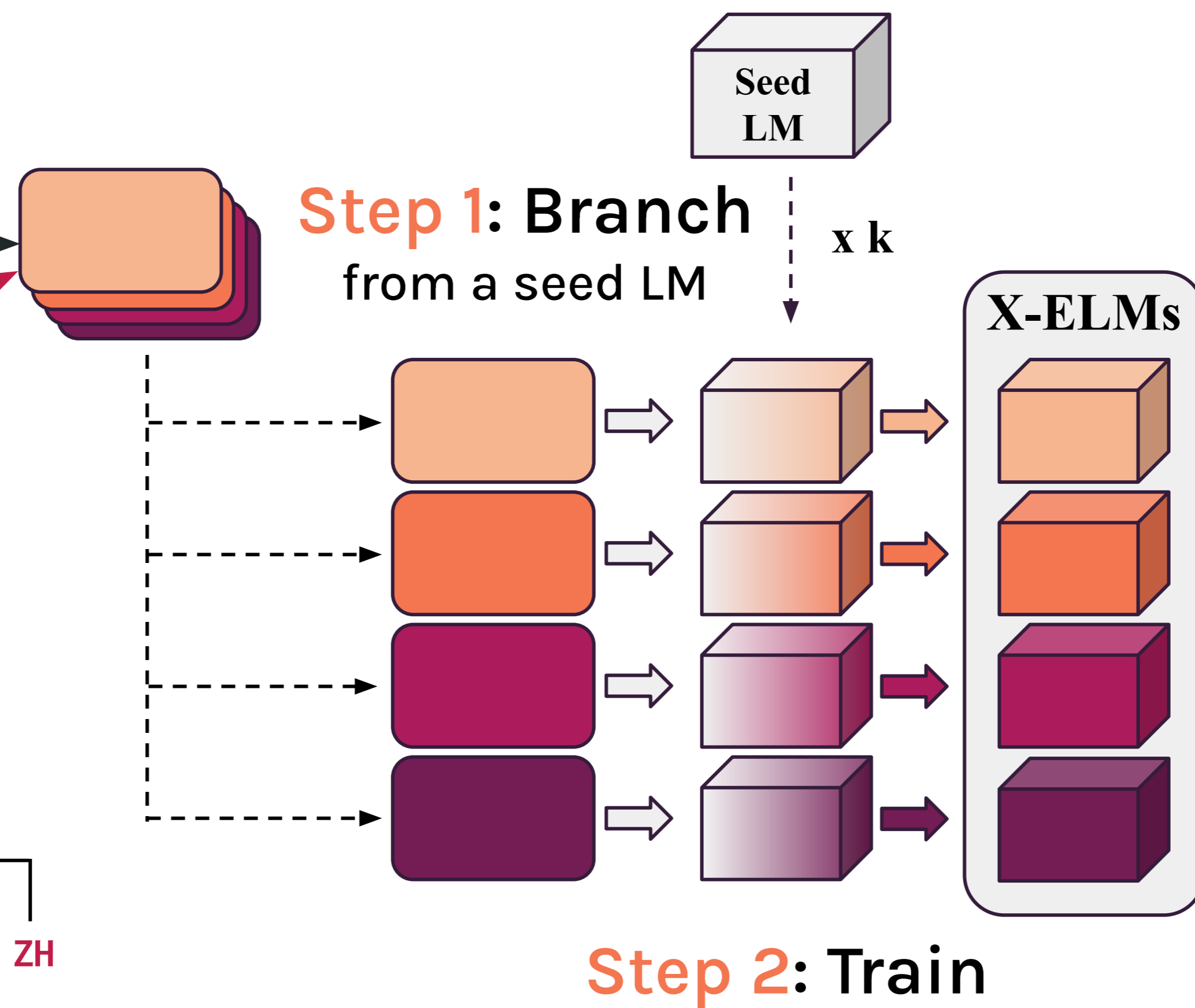
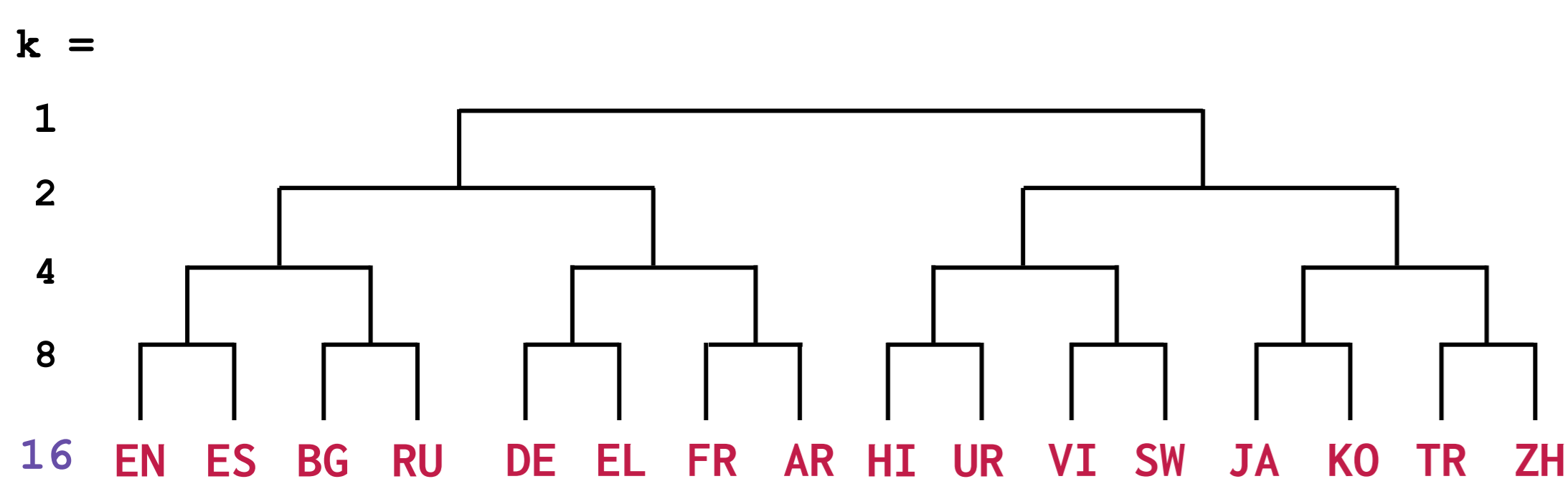


Breaking the Curse of Multilinguality with Cross-lingual Expert Language Models

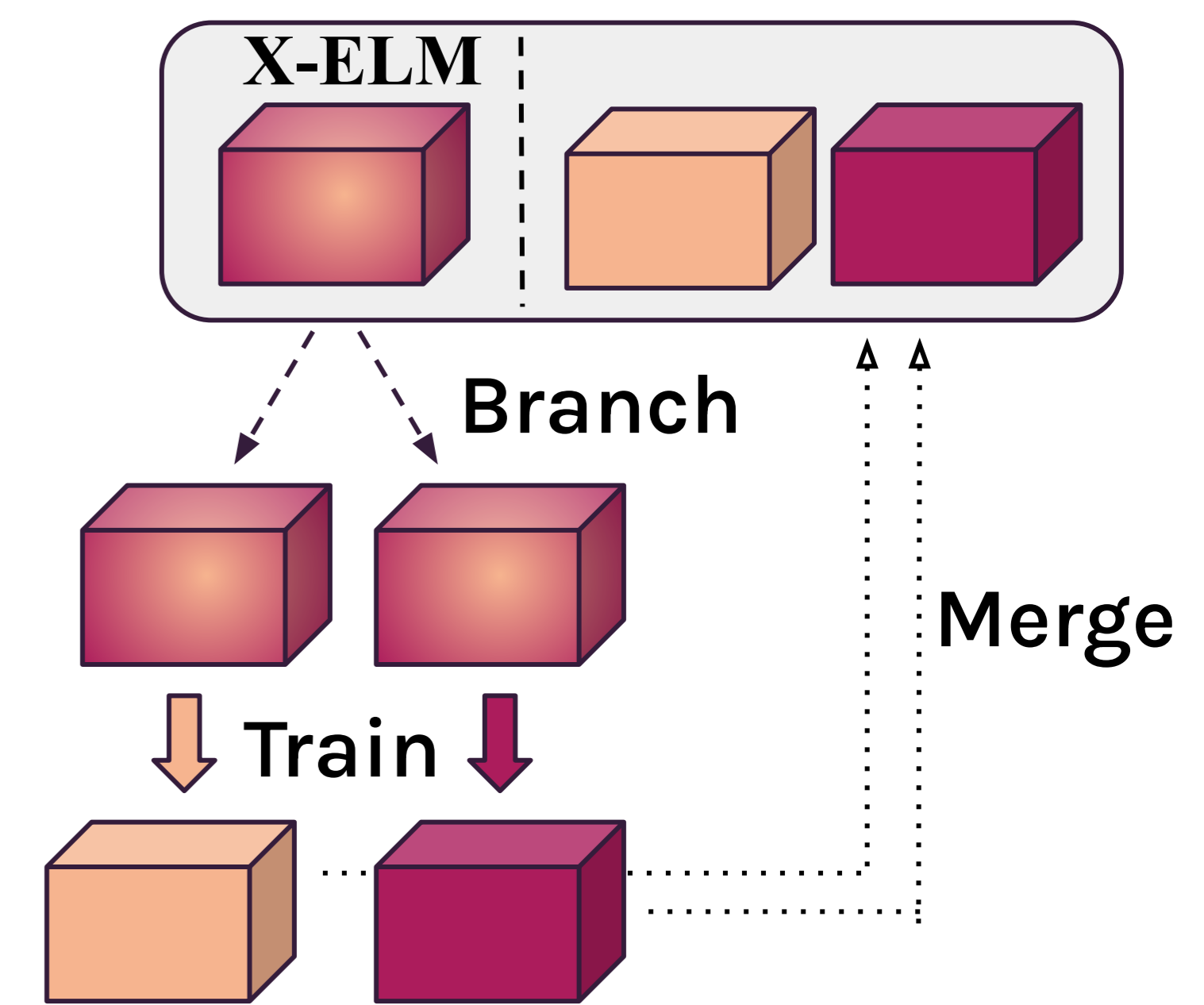
Terra Blevins¹ Tomasz Limisiewicz² Suchin Gururangan¹
Margaret Li¹ Hila Gonen¹ Noah A. Smith^{1,3} Luke Zettlemoyer¹

Step 0: Data Allocation

- ❖ TF-IDF Clustering over training documents
- ❖ **Typological Clustering** of training languages



Hierarchical Multi-round (HMR) Training



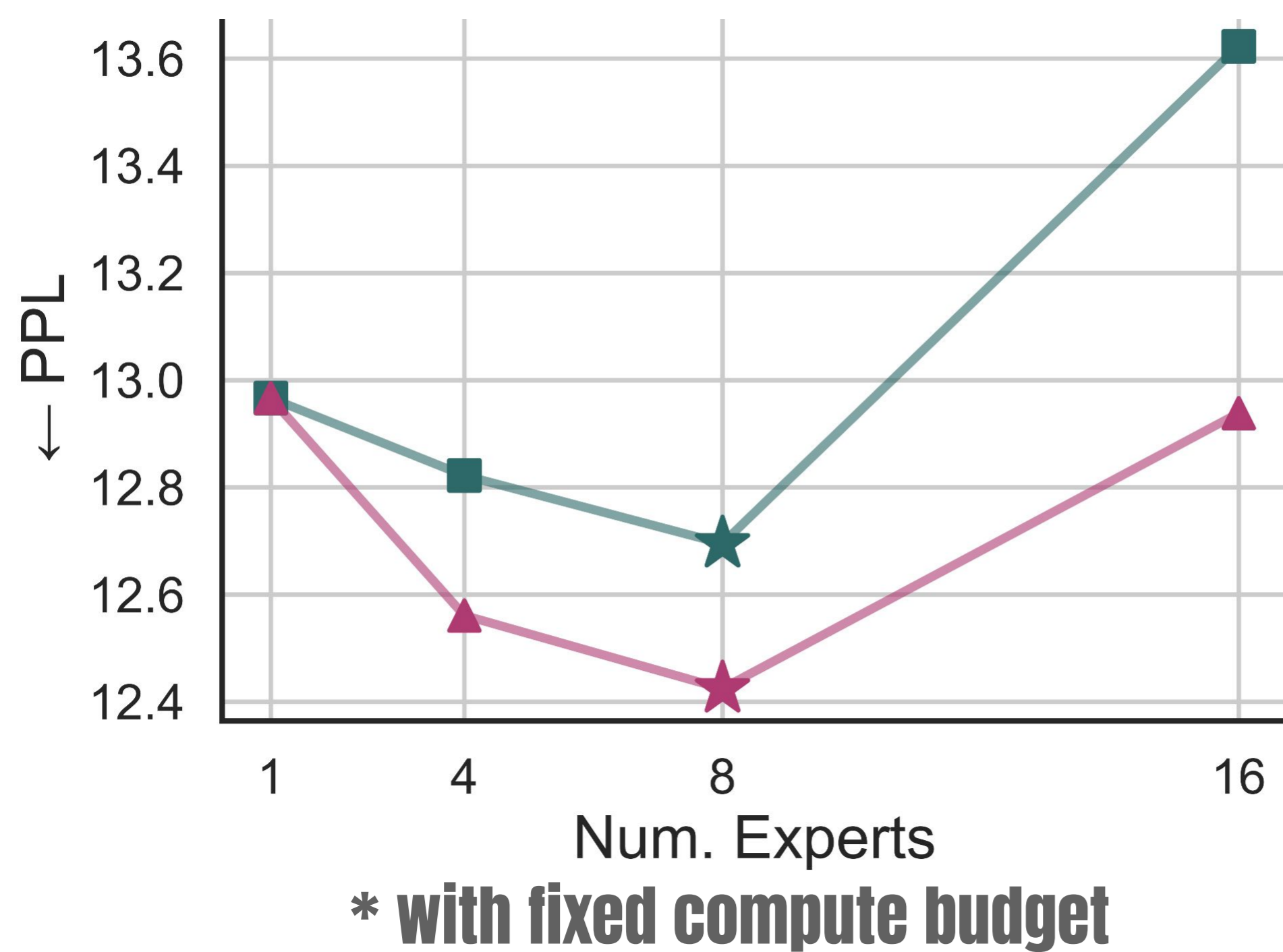
Goal → build better LM for **all** languages while maintaining **cross-linguality** of model

Train **expert language models** with multilingual **Branch-Train-Merge**

Improves language model **performance, efficiency, & adaptation** to new settings

How many (k) experts are ideal for this data setting? Using...

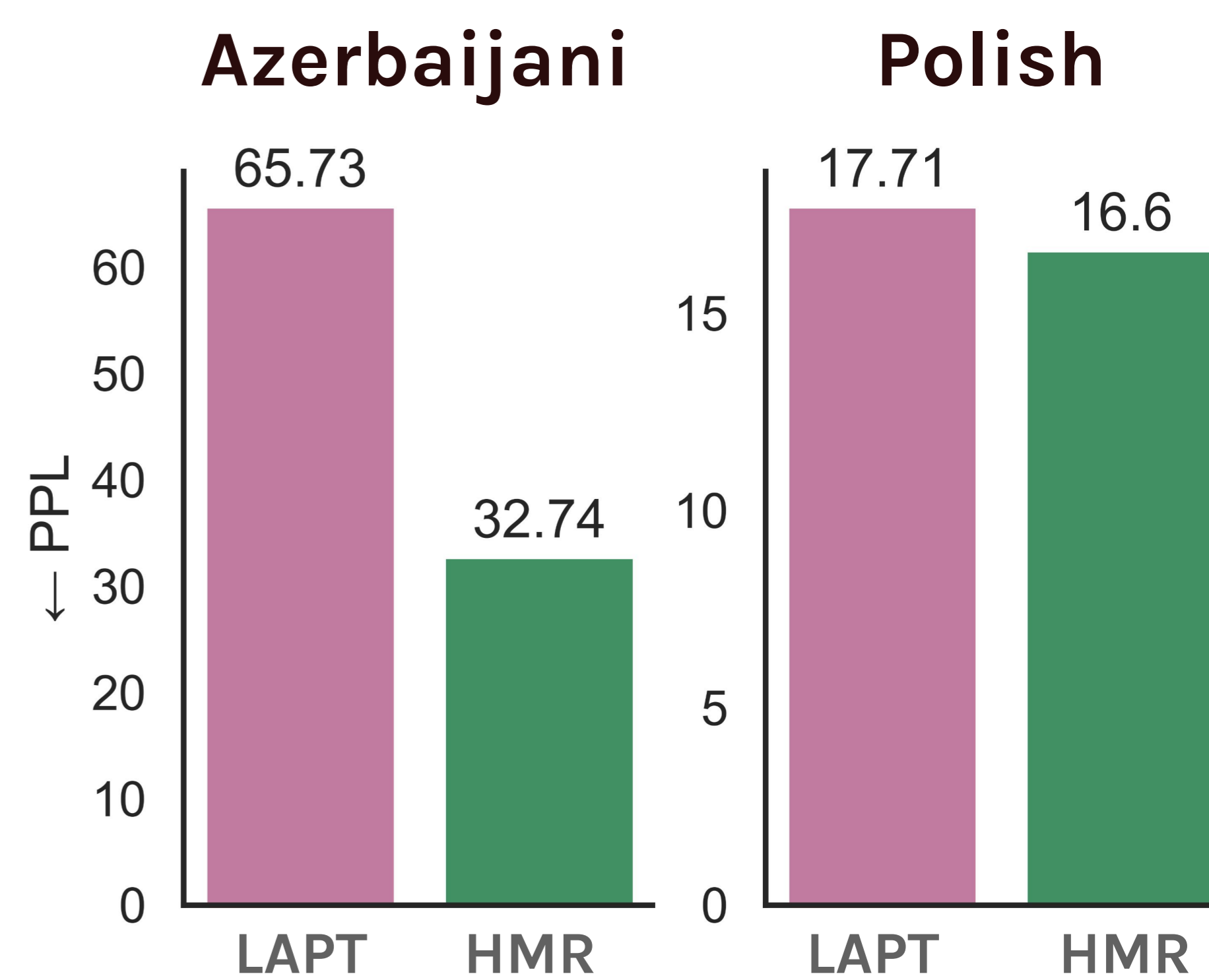
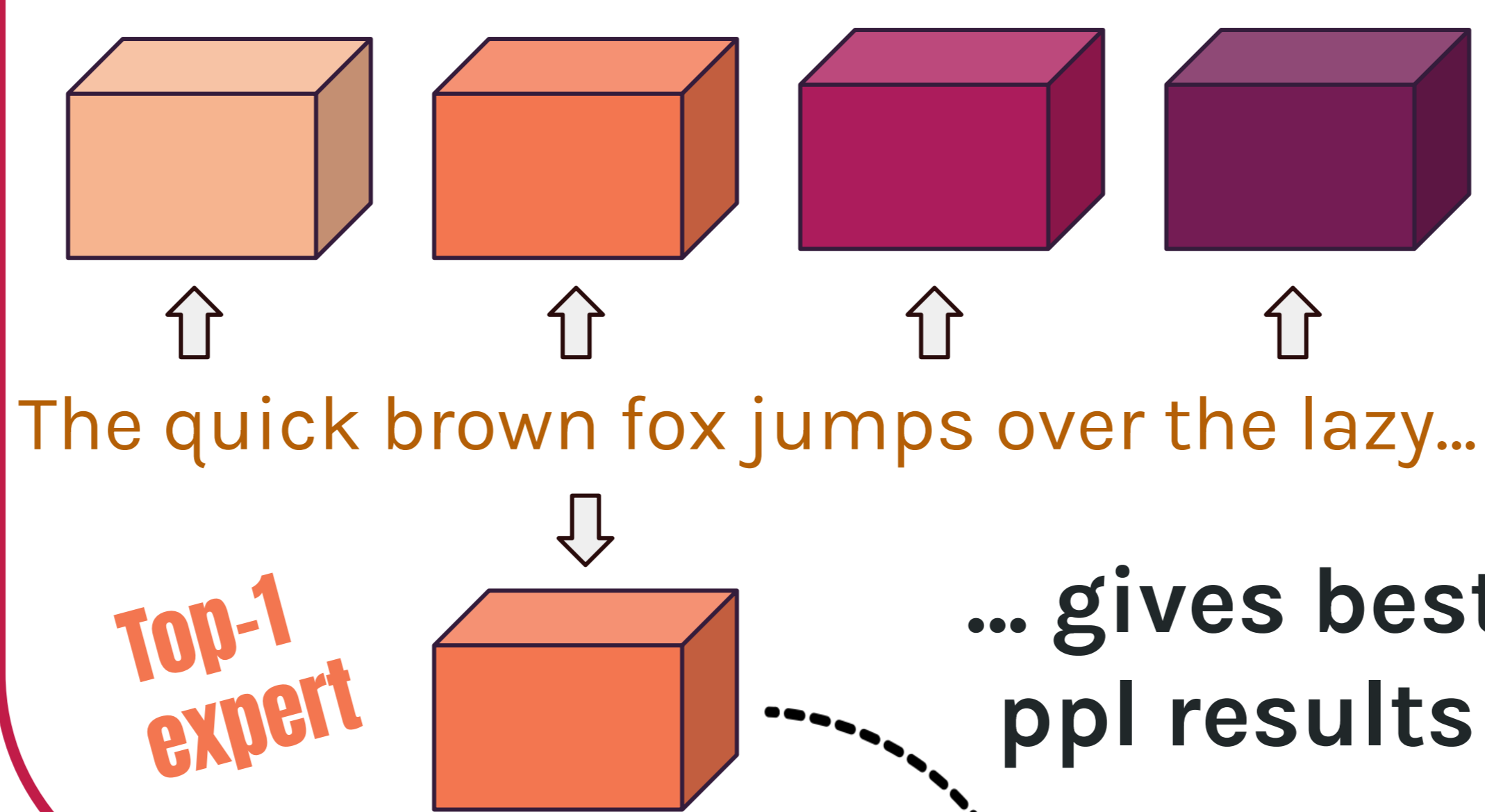
- TF-IDF Clusters
- ▲ **Typological Clusters**



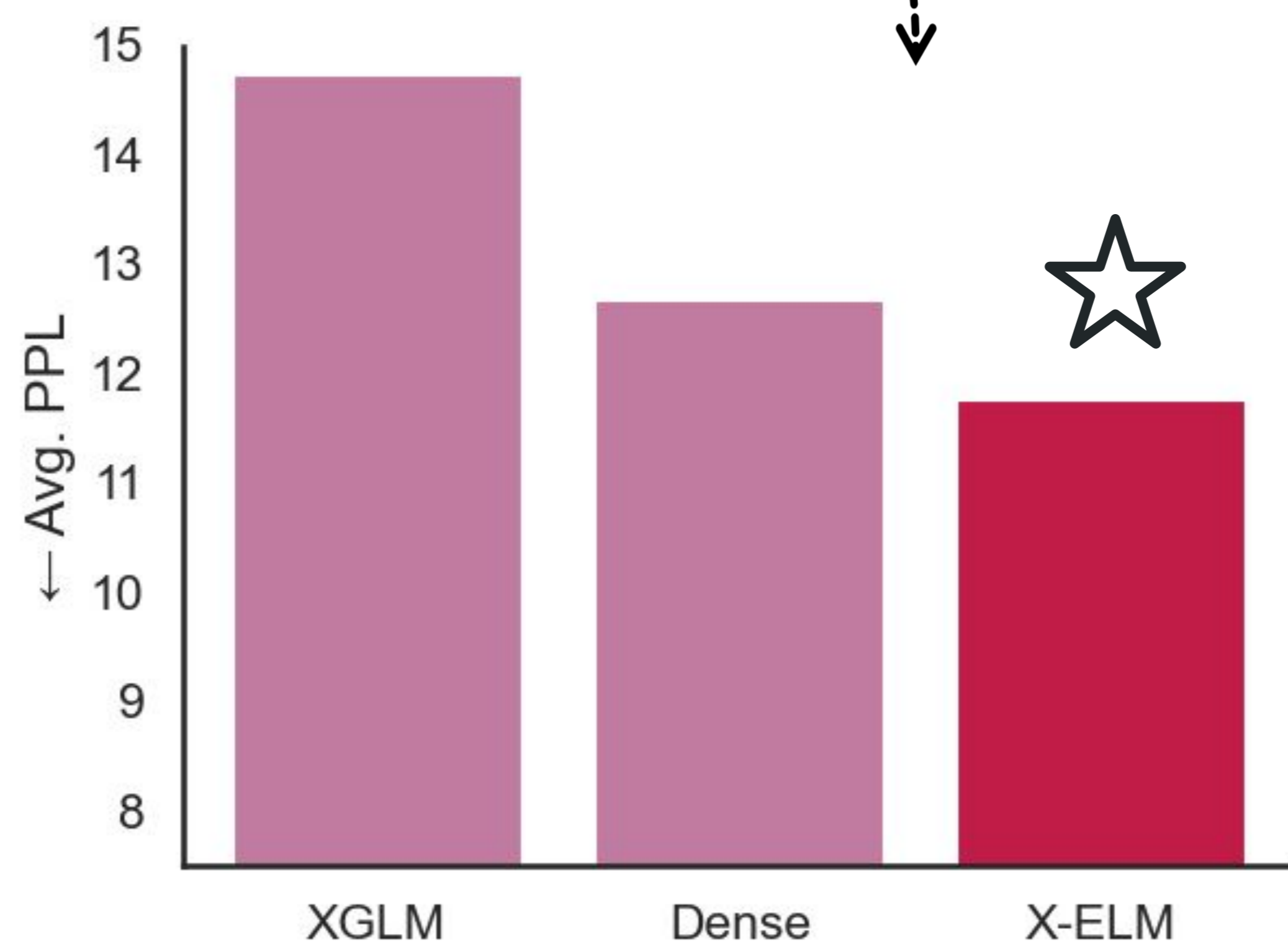
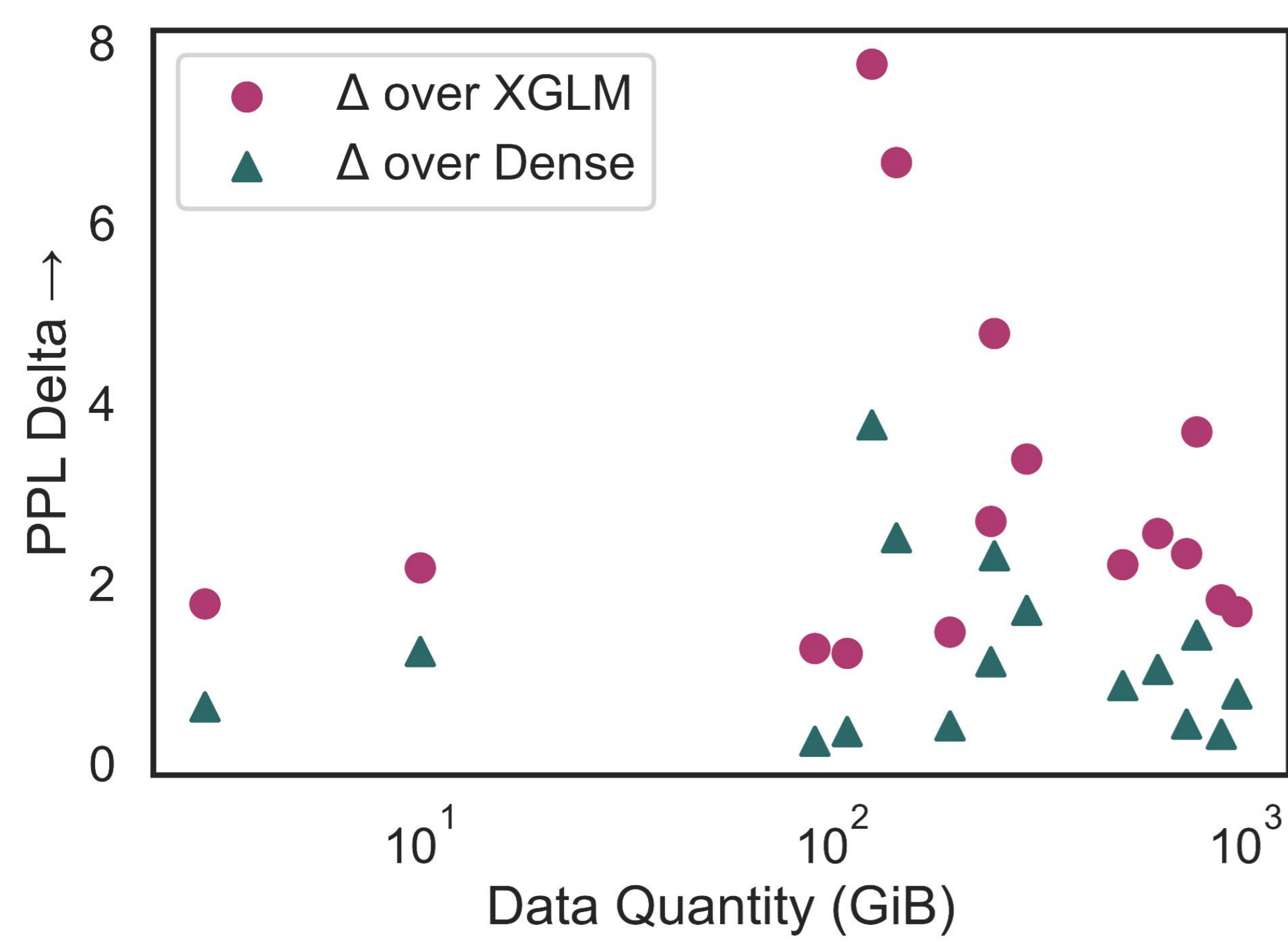
Inference with X-ELM

Ensemble

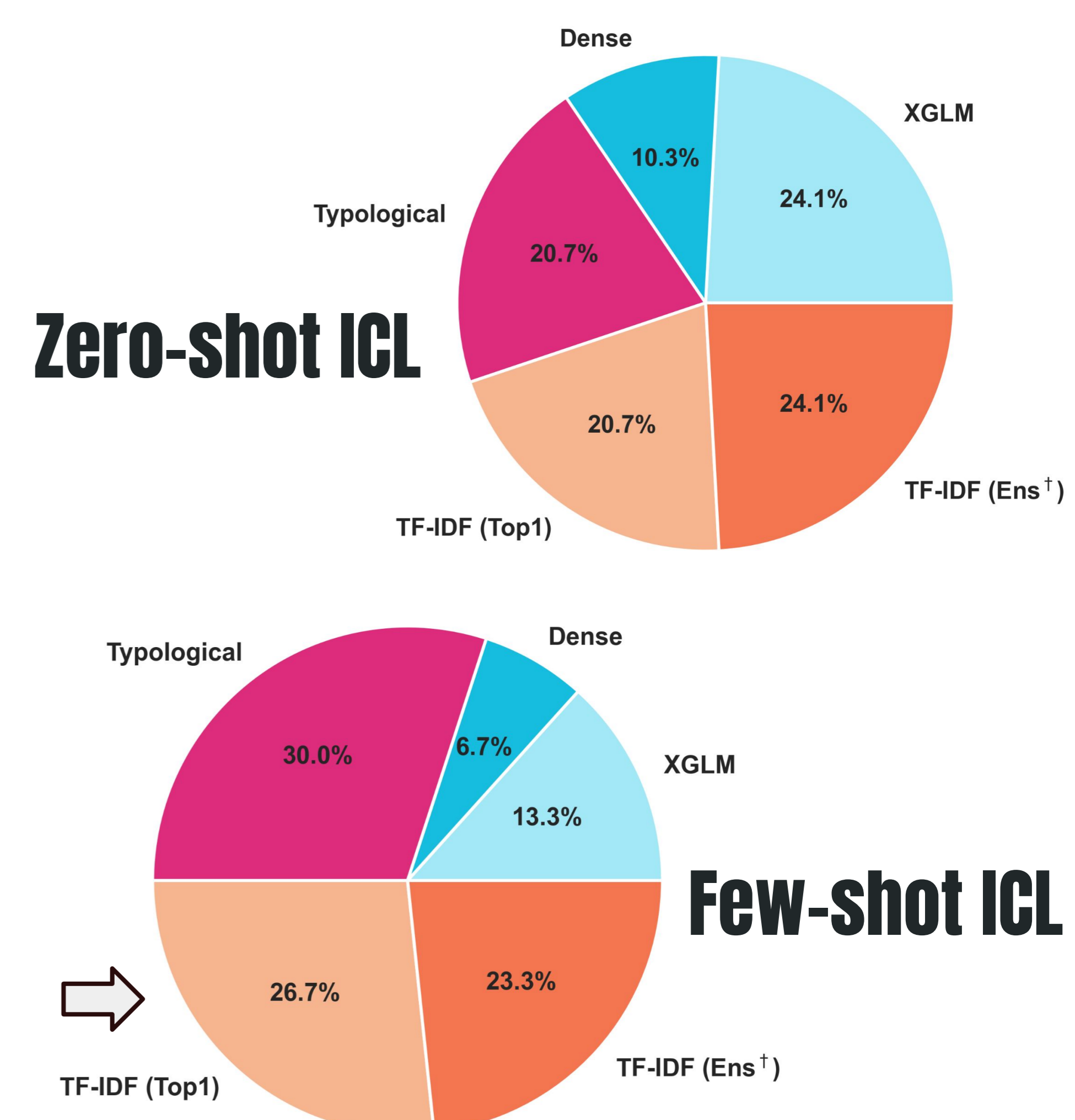
$$p_E(x_t | x_{<t}) = \sum_{e \in E} \alpha_e \cdot p_e(x_t | x_{<t})$$



Adapting to new languages with **HMR training** (vs. LAPT)



← **X-ELM PPL** versus seed (XGLM) model and $k=1$ dense baseline →



Better language modeling for **all** languages

X-ELMs (Typological and TF-IDF) get better **win rates** on ICL than **baselines**

Check out the paper!



blvns.github.io

@terrablvs

¹ University of Washington

² Charles University in Prague

³ Allen Institute for Artificial Intelligence