

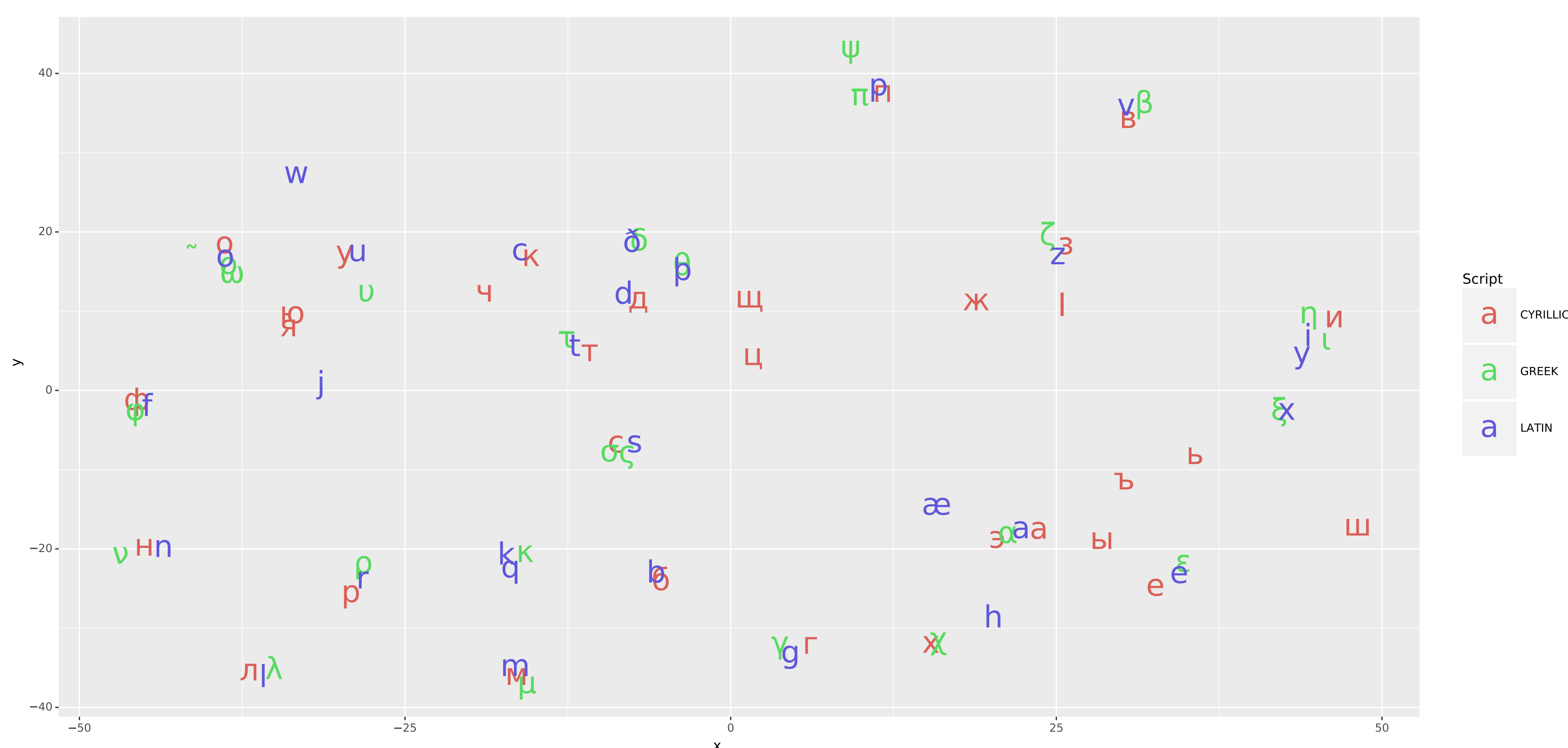
# You don't need language-specific tuning.

# Multilingual, sparse models are very competitive.

## DeepSPIN @ SIGMORPHON

Ben Peters<sup>†</sup> and André F.T. Martins<sup>†‡</sup>

<sup>†</sup>Instituto de Telecomunicações <sup>‡</sup>Unbabel



### One size (and model) fits all

- Per-language hyperparameter tuning is expensive.
- Small train sets require extra (often artificial) data.
- Multilingual training **solves both problems**.

### Sparse seq2seq

- Just replace softmax with entmax everywhere.
- Interpretable **sparse** attention.
- Sparse output distributions can make decoding **exact**.
- Requires **no other changes** to architecture.

### Multilinguality

`martins` + `<en>` → `martinz`

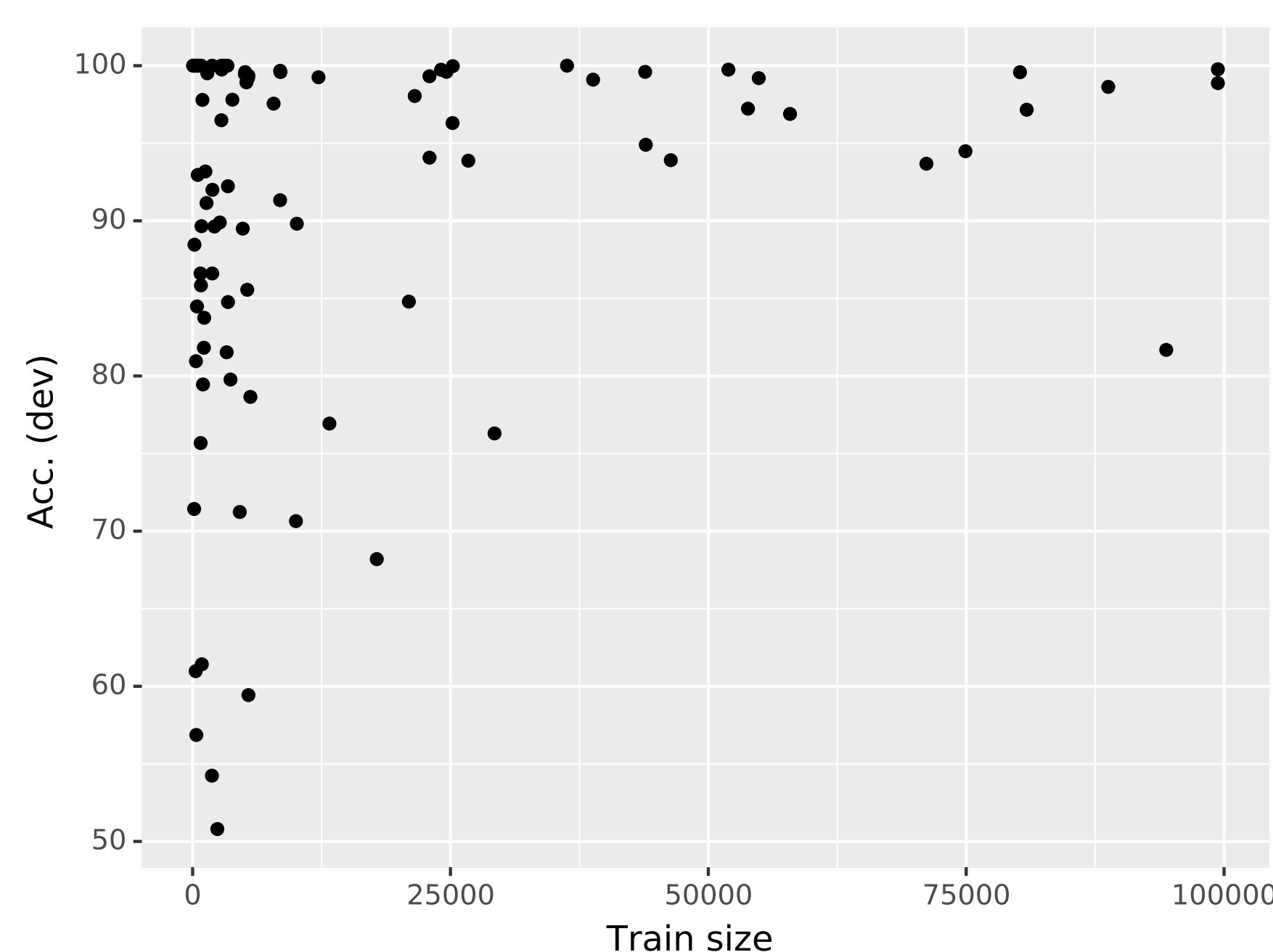
`martins` + `<pt>` → `mertiŃ`

- Label each sample with its language.
- Learn a **language embedding** for each label.
- Concatenate to the character embedding at each step.

### Phonology-aware char embeddings

- Multilingual g2p maps from disjoint scripts to shared IPA.
- Grapheme embeddings cluster by **phonological** similarity.
- Applications to transliteration?

### How data-hungry is inflection?



- Inflection-sparsemax dev results; each point is a language.
- Bigger is better, but other factors make a huge difference.
- Typology** matters, but it doesn't explain the sometimes large differences between related languages.

### Inflection results

Model	Acc. ↑	Lev. Dist. ↓
Inflection-entmax-1.5	90.5	0.217
Inflection-sparsemax	90.9	0.211
Baseline	90.6	0.215

- Tied for first place!
- All models use multi-encoder RNNs
- Same model as DeepSPIN's 2019 submission

### g2p results

Model	WER ↓	PER ↓
RNN-entmax-1.5	14.47	2.85
RNN-sparsemax	14.19	2.78
Transformer-entmax-1.5	14.15	2.92
Transformer-sparsemax	14.53	2.92
Baseline (RNN)	16.84	3.99

- Third place
- Best-performing transformer

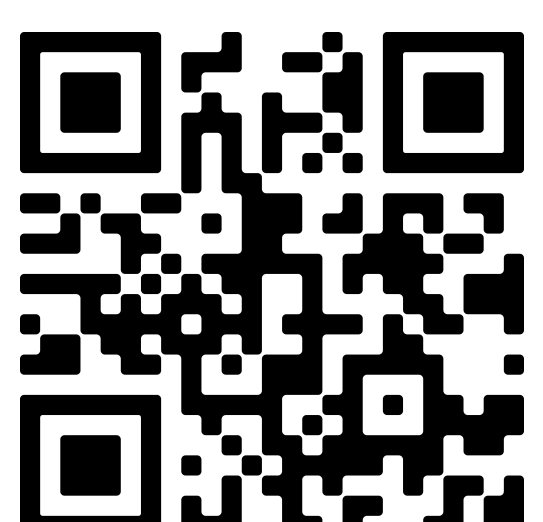
### What is entmax?

$$\begin{aligned} \alpha\text{-entmax}(z) &:= \operatorname{argmax}_{p \in \Delta^d} p^\top z + H_\alpha^\top(p) \\ &= [(\alpha - 1)z - \tau \mathbf{1}]_+^{1/\alpha-1} \end{aligned}$$

where

$$H_\alpha^\top(p) := \begin{cases} \frac{1}{\alpha(\alpha-1)} \sum_j (p_j - p_j^\alpha), & \alpha \neq 1, \\ H^S(p), & \alpha = 1 \end{cases}$$

- $\alpha = 1$  → softmax
- $\alpha > 1$  → sparsity possible
- $\alpha = 2$  → sparsemax
- $\alpha = \infty$  → argmax
- Sparse, **differentiable** softmax replacement.



← paper