# A bit of Progress and Stronger n-gram LM Baselines

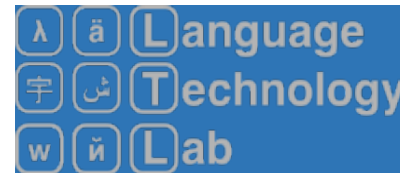**Ehsan Shareghi***      Daniela Gerz      Ivan Vulić      Anna Korhonen
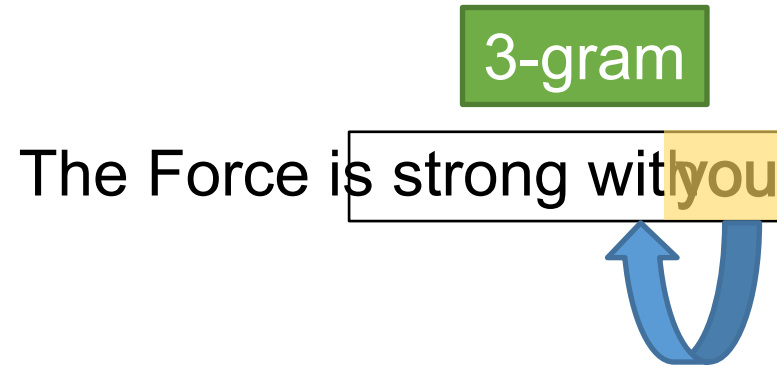
# n-gram LM

$$P(w_1^N) = \prod_{i=1}^{N} P(w_i | w_{i-n+1}^{i-1})$$

3-gram

The Force is strong with you

## Questions:

- Do I know what is Kneser-Ney (KN), and what is Modified KN (MKN)?

- Has there been any progress in n-gram smoothing?

- Are NLM always superior to n-grams? When are they likely to fall short?

# Smoothing in general

Donald Trump is a X

$$P(x = \text{politician} \mid \text{Donald Trump is a}) = \frac{\text{count(Donald Trump is a politician)}}{\text{count(Donald Trump is a)}}$$

count(Donald Trump is a politician) = 0
count(Trump is a politician) = 100

P(politician|Donald Trump is a) ≈ P(politician|Trump is a)

# Smoothing- something old something new

$$\beta(w_i|w_{i-n+1}^{i-1},\Theta)+\gamma(w_{i-n+1}^{i-1},\Theta)P(w_i|w_{i-n+2}^{i-1},\Theta)$$

# Smoothing- something old something new

$$\boxed{\beta(w_i|w_{i-n+1}^{i-1},\Theta)+\gamma(w_{i-n+1}^{i-1},\Theta)P(w_i|w_{i-n+2}^{i-1},\Theta)}$$

KN

$$\dfrac{\beta(w_i|w_{i-n+1}^{i-1},\Theta)}{\dfrac{c(w_{i-n+1}^{i})-D_n}{c(w_{i-n+1}^{i-1})}} \qquad \dfrac{\Theta}{D_n}$$

# Smoothing- something old something new

$$\boxed{\beta(w_i|w_{i-n+1}^{i-1},\Theta)+\gamma(w_{i-n+1}^{i-1},\Theta)P(w_i|w_{i-n+2}^{i-1},\Theta)}$$

| | $\beta(w_i|w_{i-n+1}^{i-1},\Theta)$ | $\Theta$ |
|---|---|---|
| KN | $\dfrac{c(w_{i-n+1}^i)-D_n}{c(w_{i-n+1}^{i-1})}$ | $D_n$ |
| MKN | $\dfrac{c(w_{i-n+1}^i)-D_n^{c(w_{i-n+1}^i)}}{c(w_{i-n+1}^{i-1})}$ | $D_n^{i\in\{1,2,3+\}}$ |

# Smoothing- something old something new

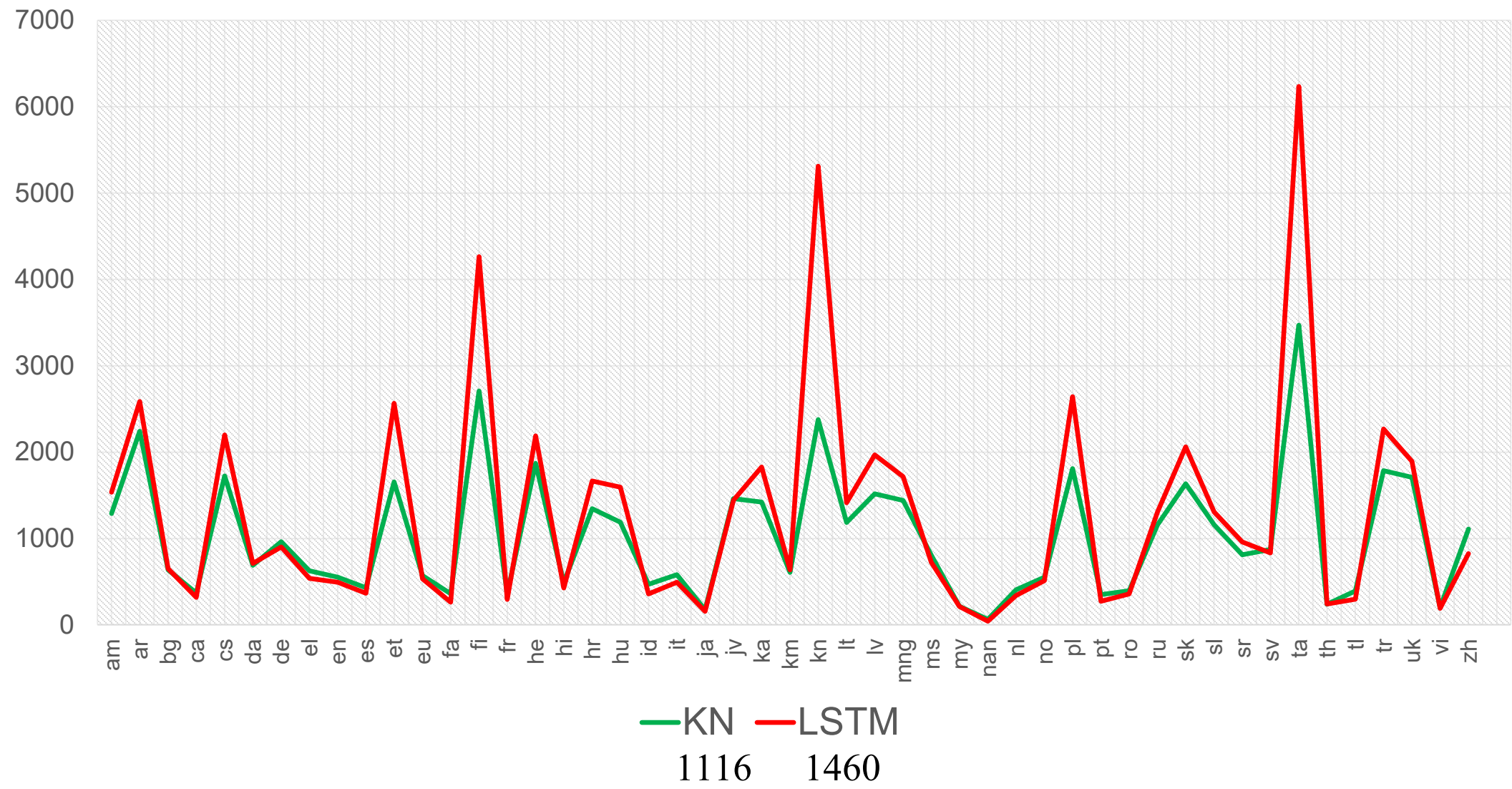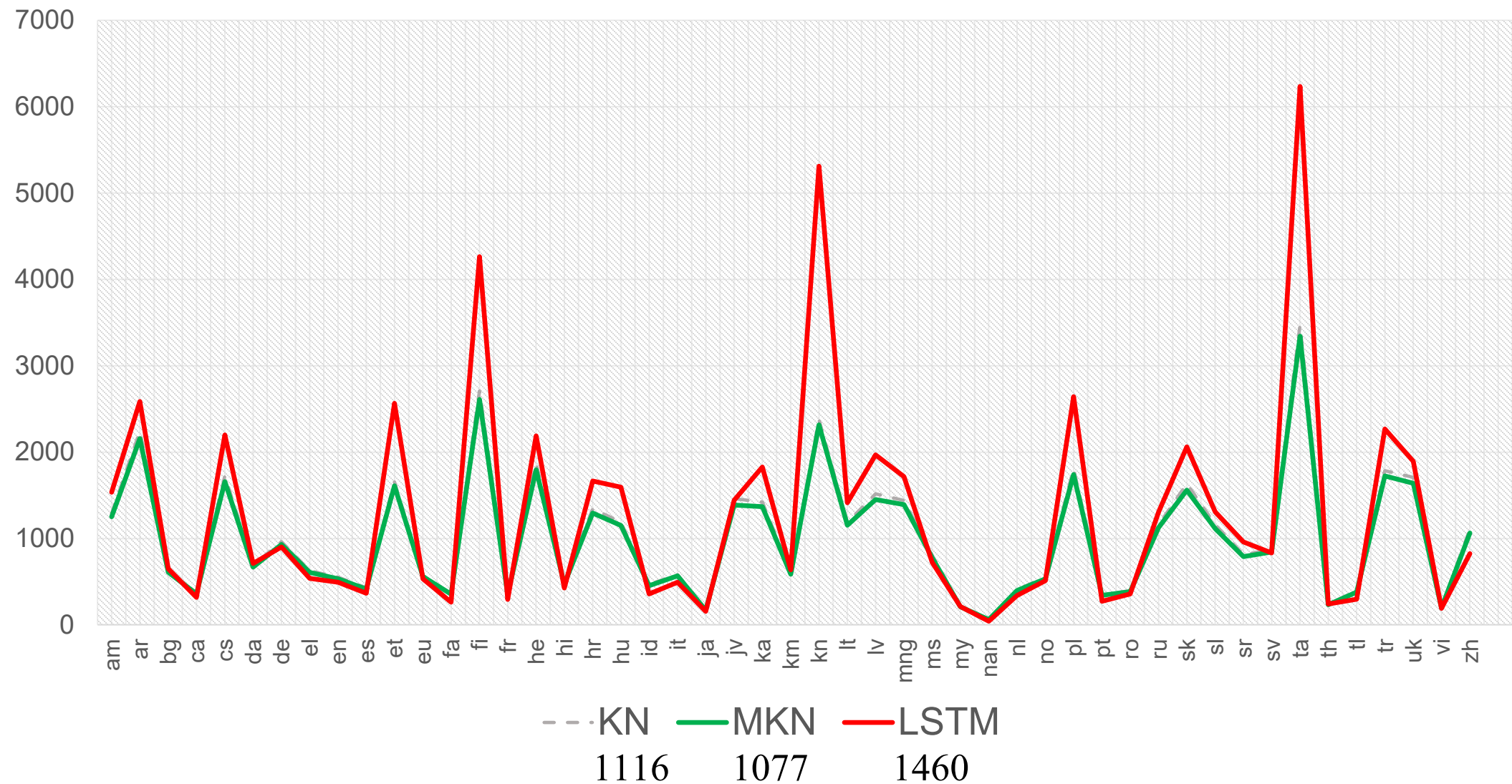$$\beta(w_i | w_{i-n+1}^{i-1}, \Theta) + \gamma(w_{i-n+1}^{i-1}, \Theta) P(w_i | w_{i-n+2}^{i-1}, \Theta)$$

| | $\beta(w_i \| w_{i-n+1}^{i-1}, \Theta)$ | $\Theta$ |
|---|---|---|
| KN | $\dfrac{c(w_{i-n+1}^i) - D_n}{c(w_{i-n+1}^{i-1})}$ | $D_n$ |
| MKN | $\dfrac{c(w_{i-n+1}^i) - D_n^{c(w_{i-n+1}^i)}}{c(w_{i-n+1}^{i-1})}$ | $D_n^{i \in \{1,2,3+\}}$ |
| GKN | $\dfrac{c(w_{i-n+1}^i) - D_n^{c(w_{i-n+1}^i)}}{c(w_{i-n+1}^{i-1})}$ | $D_n^{i \in \{1,\dots,10+\}}$ |

# Smoothing- something old something new

$$\beta(w_i|w_{i-n+1}^{i-1},\Theta)+\gamma(w_{i-n+1}^{i-1},\Theta)P(w_i|w_{i-n+2}^{i-1},\Theta)$$

| | $\beta(w_i|w_{i-n+1}^{i-1},\Theta)$ | $\Theta$ |
|---|---|---|
| KN | $\dfrac{c(w_{i-n+1}^i)-D_n}{c(w_{i-n+1}^{i-1})}$ | $D_n$ |
| MKN | $\dfrac{c(w_{i-n+1}^i)-D_n^{c(w_{i-n+1}^i)}}{c(w_{i-n+1}^{i-1})}$ | $D_n^{i\in\{1,2,3+\}}$ |
| GKN | $\dfrac{c(w_{i-n+1}^i)-D_n^{c(w_{i-n+1}^i)}}{c(w_{i-n+1}^{i-1})}$ | $D_n^{i\in\{1,...,10+\}}$ |
| BKN | $\dfrac{c(w_{i-n+1}^i)-D_{w_{i-n+1}}t_{w_{i-n+1}}^{w_i}}{c(w_{i-n+1}^{i-1})+\theta_{w_{i-n+1}}}$ | $D_{w_{i-n+1}},\theta_{w_{i-n+1}},t_{w_{i-n+1}}^{w_i}$ |

# Perplexity
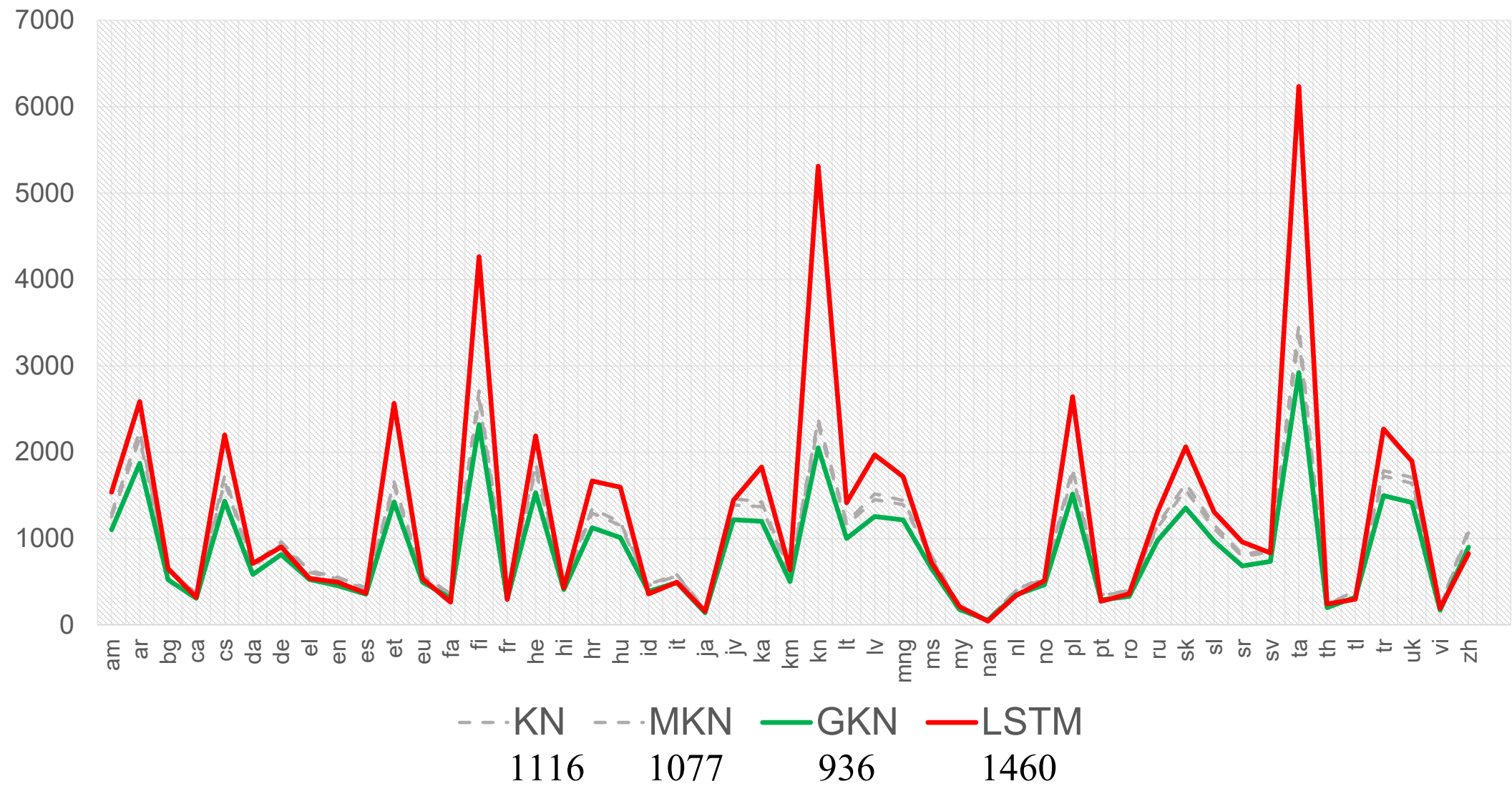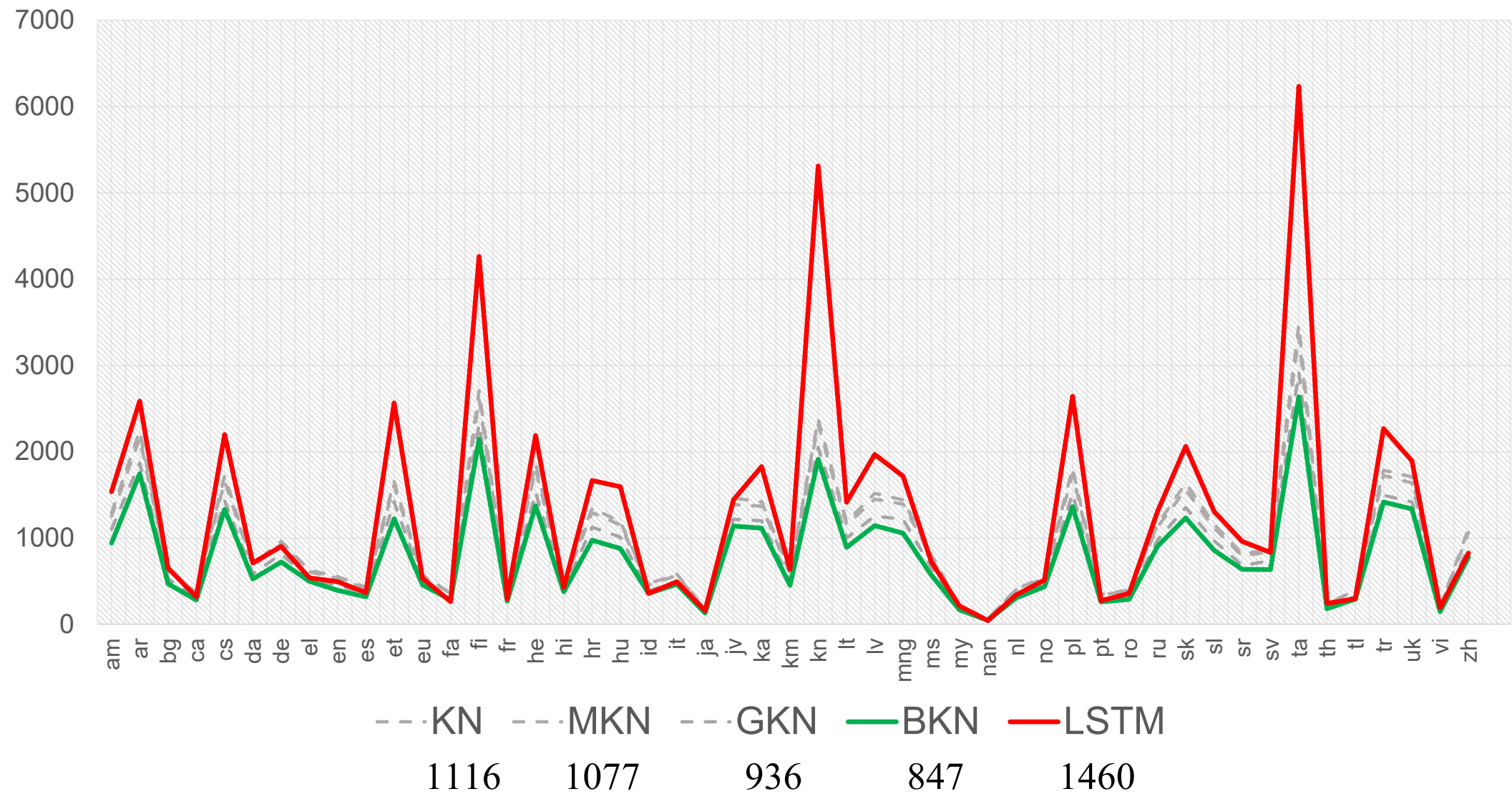


KN — LSTM
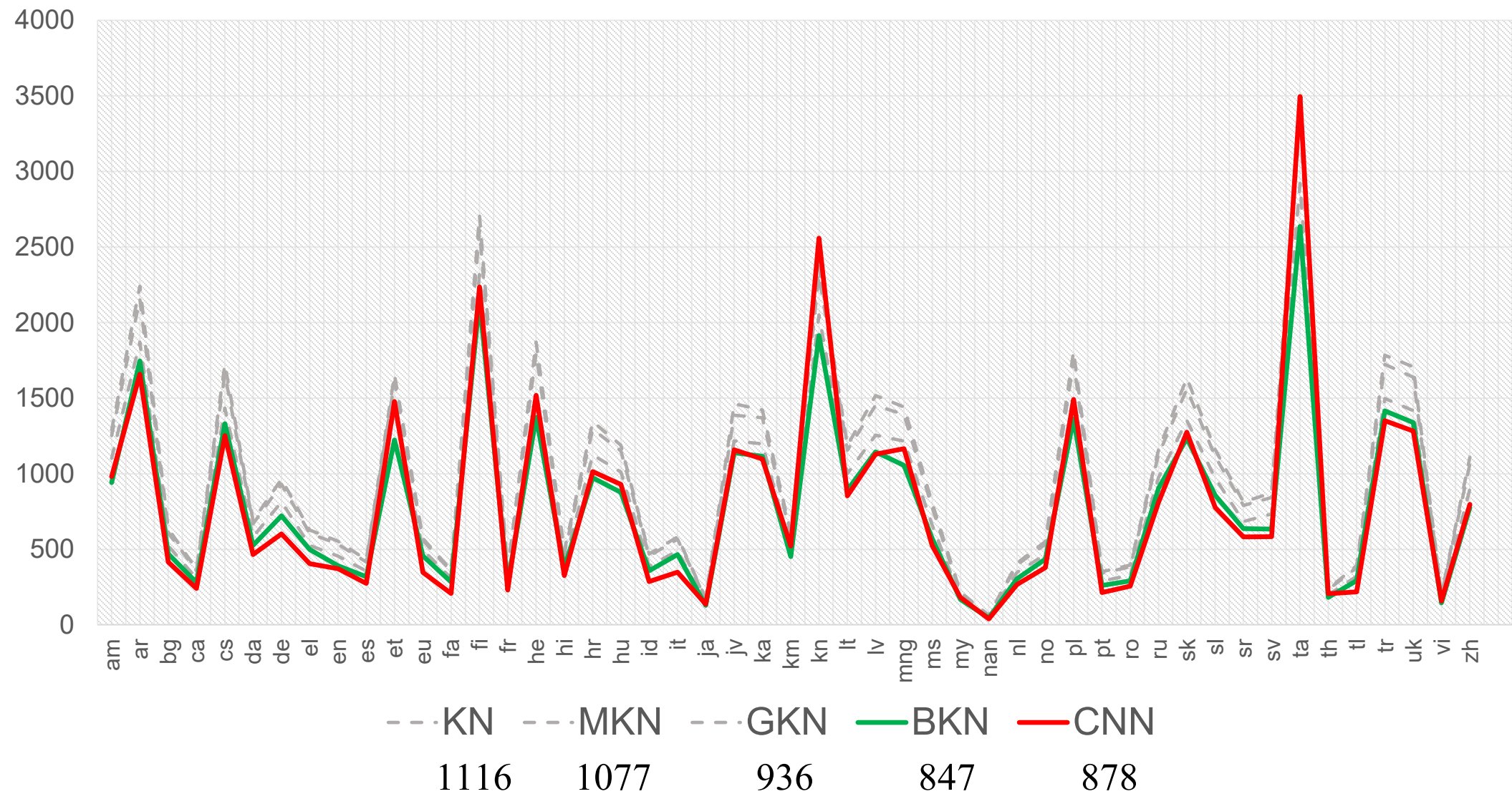1116    1460

# Perplexity



KN — MKN — LSTM
1116   1077   1460

# Perplexity

# Perplexity



| | KN | MKN | GKN | BKN | LSTM |
|---|---|---|---|---|---|
| | 1116 | 1077 | 936 | 847 | 1460 |

# Perplexity

# Perplexity



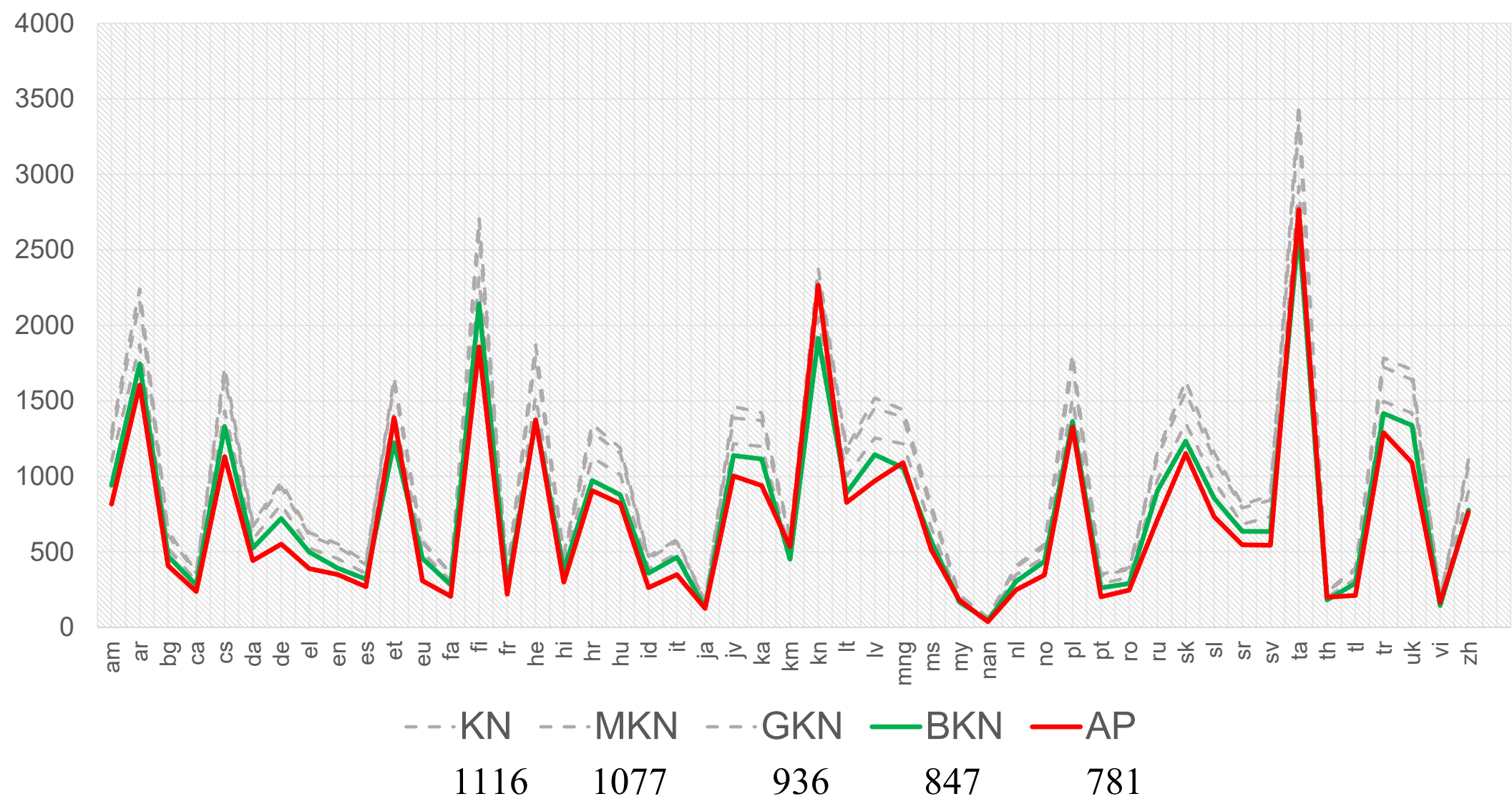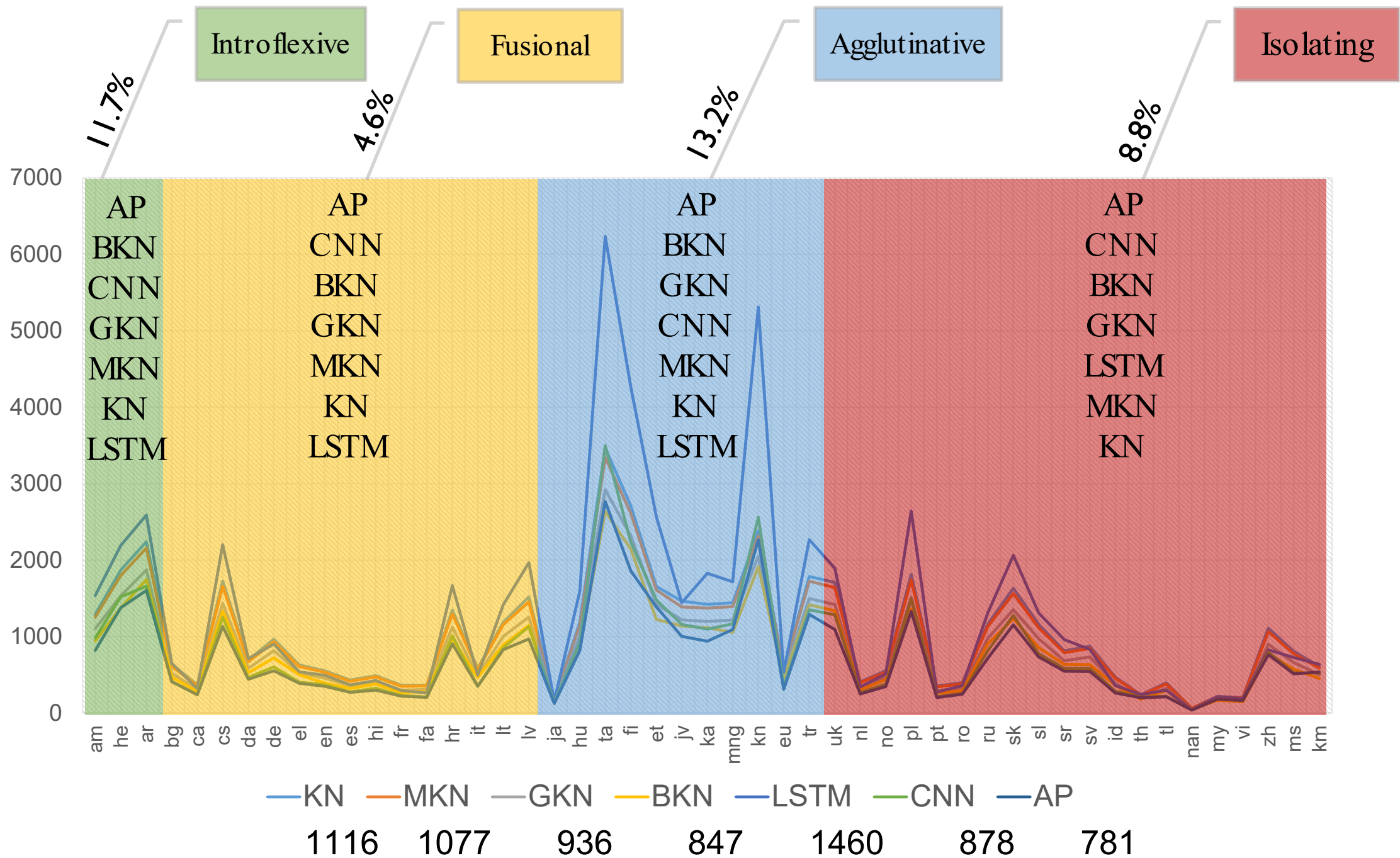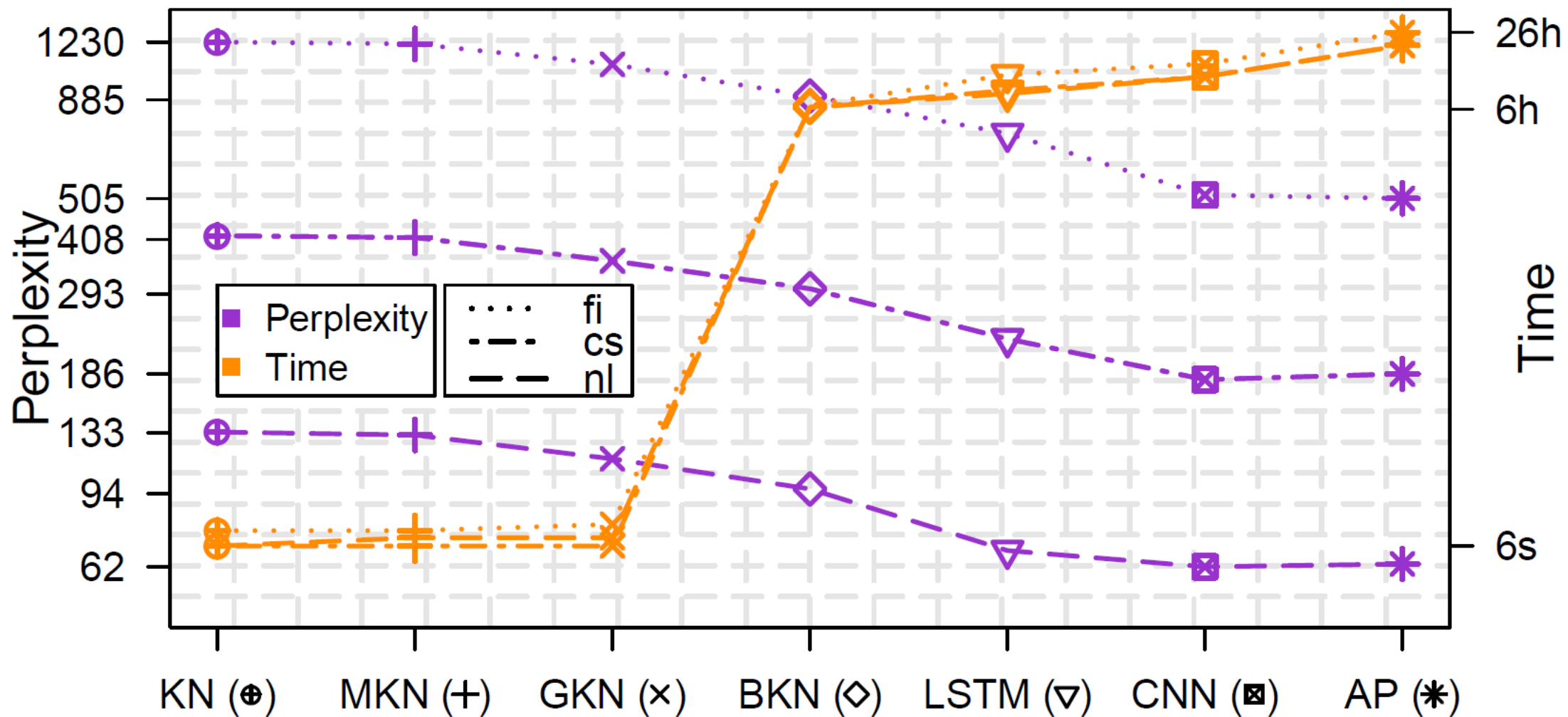| KN | MKN | GKN | BKN | AP |
|----|-----|-----|-----|-----|
| 1116 | 1077 | 936 | 847 | 781 |

# Training Time

# Wrapping up …

- n-grams are highly competitive with neural LMs for low-resource setting, high OOV ratio, or for languages with high type-token ratio

- Recent developments in n-gram models permit to lift the finite-order Markov assumption, hence in theory models should be capable of capturing long range dependencies

- The gap between neural and stand-alone n-gram models could be reduced by (somehow) incorporating continuous word representations into n-gram models

- n-gram models have far more attractive computational properties (Memory/Time usage) for both training and inference steps. So invest in improving neural models computational shortcomings.

Thanks! ☺

Contact: Ehsan Shareghi - es776@cam.ac.uk