

Q1: Are all languages equally hard to model?

Not with current models, all models we test have very different performance on different languages.

Q2: So what makes a language hard to model?

Hypothesis: inflectional morphology.

- \rightarrow LM performance negatively correlated with morphological counting complexity (MCC; Sagot, 2013)
- \rightarrow Correlation disappears when modeling lemmata (obtained using UDPipe (Straka et al., 2016)) instead of forms

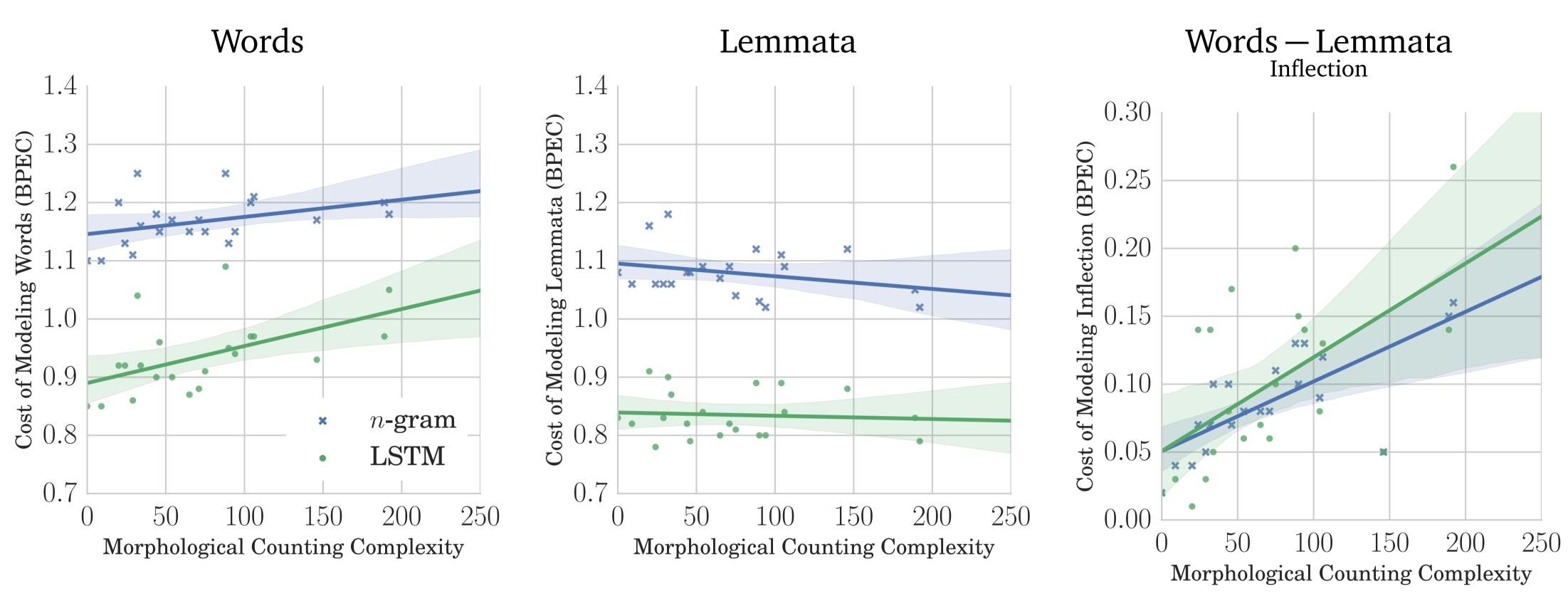
Differing corpora are unfair \rightarrow Multi-text

Domain differences impact estimates, so use aligned parallel text, 21 usable languages in Europarl (Koehn, 2005):

bg/cs/da/de/el/en/es/fr/it lt/lv/nl/pl/pt/ro/sk/sl/sv	et/fi/hu
Indo-European	Uralic

Translationese (Baker, 1993, translations is stylistically different from "native" text) only underestimates the difficulty of non-English languages—and we find the opposite!

Are All Languages Equally Hard to Language-Model? Ryan Cotterell¹, Sebastian J. Mielke¹, Jason Eisner¹, Brian Roark² ¹Department of Computer Science, Johns Hopkins University, Baltimore, USA ²Google



Each point is a language, the cost of modeling is plotted against the MCC of a language.

Closed-vocab is unfair \rightarrow Open-vocab LMs

Replacing rare words with UNK leads to unfairly good scores for languages with many word types (e.g., morphologically rich languages). So instead use *open-vocabulary* LMs like:

• Kneser-Ney smoothed **n-gram** LM over "flat" hybrid representations (Bisani and Ney, 2005):

A character-level LSTM LM (Sundermeyer et al., 2012; Zaremba et al., 2014):

> $p(c_t | \mathbf{c}_{< i}) = \operatorname{softmax}(Wh_i + b)$ $h_i = \text{LSTM}(h_{i-1}, c_{i-1})$

Information contained in...

BPC is unfair \rightarrow Bits per English character

BPC (i.e., the information contained in one character) values depend on a language and cannot be compared. Example: these three strings that contain equal information (total cross-entropy is around 6 bits for each of them), but different BPC:

Sentence/word/lemma	Σ bits
EN 🗮 c o u p	≈6
DE <mark>E</mark> Putsch	≈ 6
CZ 🚬 puč	≈ 6

We normalize the total number of bits (i.e., information) for length, arbitrarily choosing the number of English characters in the utterance, obtaining bits per **English** character (BPEC).



