

Multilingual Language Models: NLP Beyond English



Eric Wallace
CS 288

NLP Beyond English

- An overwhelming majority of NLP research focuses on English!

How to build non-English NLP systems?

- Translate baseline
- Monolingual LMs for each language
- Multilingual LMs

Translate Baseline



Pros:

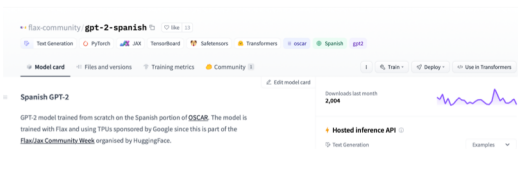
- Straightforward to implement
- Surprisingly strong baseline, especially for classification tasks

Cons:

- Suffers from cascading errors
- Limited to languages that translation systems support
- Can be slow and computationally expensive
- Translation is fundamentally lossy?

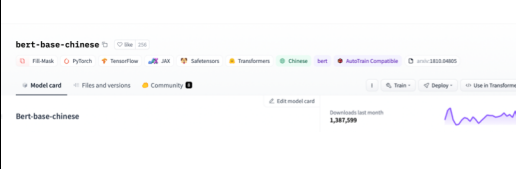
Monolingual LMs

- Can we just repeat the LM pre-training pipeline for other languages?
 - Sort of!



Monolingual LMs

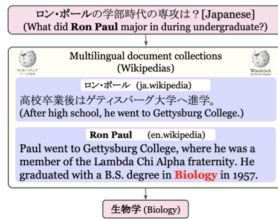
- Can we just repeat the LM pre-training pipeline for other languages?
 - Sort of!



Few-shot Learning in Other Languages

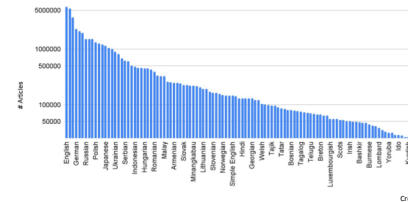
Language	Premise / Hypothesis	Genre	Label
English	You don't have to stay there. You can leave.	Face-To-Face	Entailment
French	La figure 4 montre la courbe d'offre des services de partage de travaux. Les services de partage de travaux ont une offre variable.	Government	Entailment
Spanish	Y se estremeció con el recuerdo. El pensamiento sobre el acontecimiento hizo su estremecimiento.	Fiction	Entailment
German	Während der Depression war es die ärmste Gegend, kurz vor dem Hungertod. Die Weltwirtschaftskrise dauerte mehr als zehn Jahre an.	Travel	Neutral
Swahili	Ni silaha ya plastiki ya moja kwa moja inayopigwa risasi. Inadumu zaidi kuliko silaha ya chuma.	Telephone	Neutral
Russian	Их ниса музонаменес рени уне на спреченне 85 лет. Мах тринае маршале рени музонаменес.	Letters	Contradiction
Chinese	让我告诉你，美国人最终如何看待你作为独立顾问的表现。 美国人完全不知道您是独立律师。	Slack	Contradiction
Arabic	لنحتاج التوكيد ان نكون قادرا على قيام مسؤوليات المنهج في مستقبلنا كوكلاء في الشرف مع (أ) كادنا لاجلنا (ب)	Nine-Eleven	Contradiction

Few-shot Learning in Other Languages



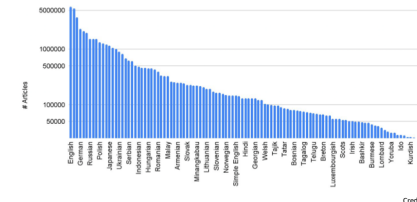
Challenges with Monolingual LMs

- There is not enough unlabeled data for each language

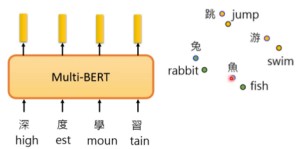


Challenges with Monolingual LMs

- Compute and complexity of serving 100-1000s of different models

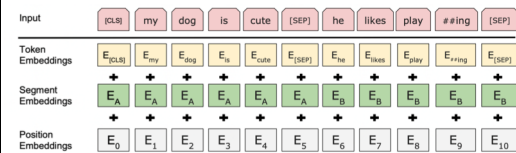


Multilingual Language Models?



Multilingual BERT

- Simply rerun BERT, except use 100+ Wikipedias and new BPE



Non-English Tokenizers

- We can use either standard BPE tokenizers or byte-level models
 - massively increase BPE size (50k \rightarrow 250k+)

