## Multilingual Language Models: NLP Beyond English

Berkeley

N L P

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

---

## Translate Baseline

español → Google Translate → english

---

## Translate Baseline

español → Google Translate → english

Pros:
- Straightforward to implement
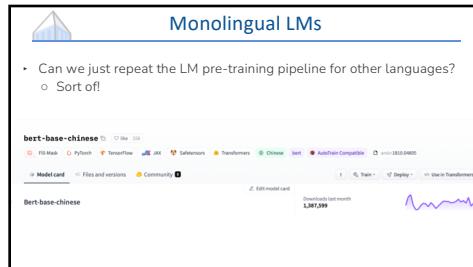- Strong baseline, especially for classification tasks

---

## Translate Baseline

español → Google Translate → english
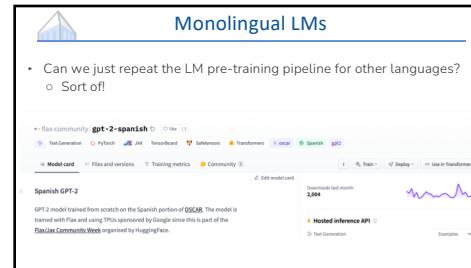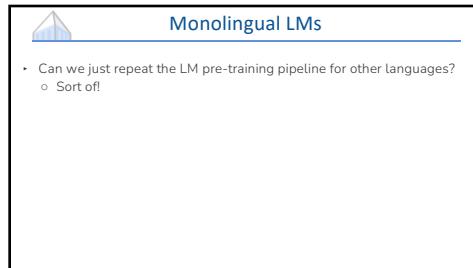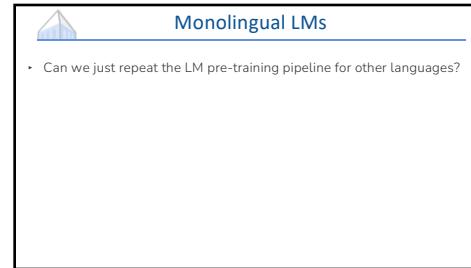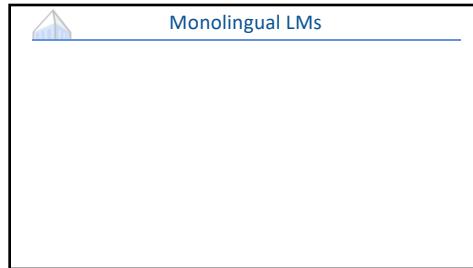
Pros:
- Straightforward to implement
- Strong baseline, especially for classification tasks

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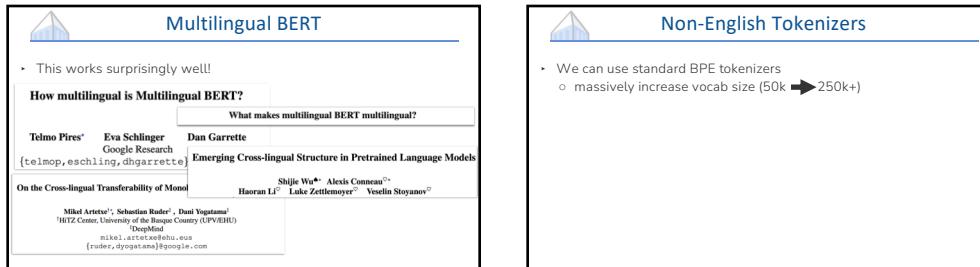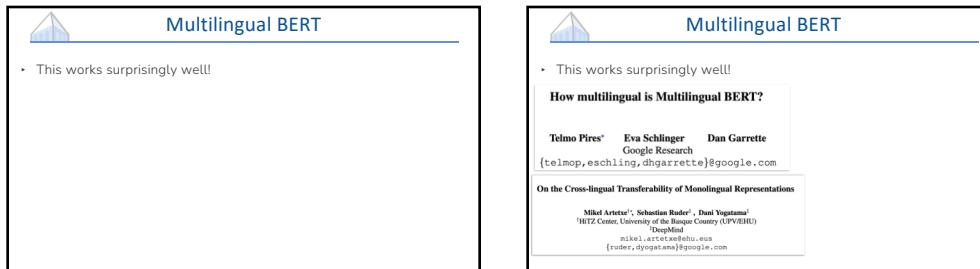- Translation is fundamentally lossy?

## Monolingual LMs

## Monolingual LMs

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

## 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!



flax-community **gpt-2-spanish**

**Spanish GPT-2**

GPT-2 model trained from scratch on the Spanish portion of OSCAR. The model is trained with Flax and using TPUs sponsored by Google since this is part of the Flax/Jax Community Week organised by HuggingFace.

Downloads last month
2,004

## Monolingual LMs

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



**bert-base-chinese**

Bert-base-chinese

Downloads last month
1,307,599

## Evaluating LMs in Other Languages



Datasets: **squad_es**

Datasets: Hello-SimpleAI/**HC3-Chinese**

Datasets: **german_legal_entity_recognition**

## Evaluating LMs in Other Languages

| Language | Premise / Hypothesis | Genre | Label |
|---|---|---|---|
| English | You don't have to stay there.<br>You can leave. | Face-To-Face | Entailment |
| French | La figure 4 montre la courbe d'offre des services de partage de travaux.<br>Les services de partage de travaux ont une offre variable. | Government | Entailment |
| Spanish | Y se estremeció con el recuerdo.<br>El pensamiento sobre el acontecimiento hizo su estremecimiento. | Fiction | Entailment |
| German | Während der Depression war es die ärmste Gegend, kurz vor dem Hungertod.<br>Die Weltwirtschaftskrise dauerte mehr als zehn Jahre an. | Travel | Neutral |
| Swahili | Ni silaha ya plastiki ya moja kwa moja inayopiga risasi.<br>Inadumu zaidi kuliko silaha ya chuma. | Telephone | Neutral |
| Russian | И мы занимаемся этим уже на протяжении 85 лет.<br>Мы только начали этим заниматься. | Letters | Contradiction |
| Chinese | 让我告诉你，美国人最终如何看待作作为独立顾问的表现。<br>美国人完全不知道您是独立律师。 | Slate | Contradiction |
| Arabic | تحلو اعوادك ال تطول قادرة على قياس مطويات التهذب<br>يمكنكالتولكات ا العرف ما إذا كانت تابعة أو د | Nine-Eleven | Contradiction |

## Challenges with Monolingual LMs

‣ There is not enough unlabeled data for each language



Credit: Graham Neubig

## Challenges with Monolingual LMs

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



Credit: Graham Neubig

## Multilingual Language Models?

## Multilingual Language Models?



## Multilingual Language Models?



Promises:
● Share world knowledge, syntax, etc. across languages?
● Use related languages to enhance transfer?

## Multilingual BERT

- Simply rerun BERT with 100+ language's Wikipedia and new BPE

## Multilingual BERT

- Simply rerun BERT with 100+ language's Wikipedia and new BPE



## Multilingual BERT

- This works surprisingly well!

## Multilingual BERT

- This works surprisingly well!

**How multilingual is Multilingual BERT?**

Telmo Pires*    Eva Schlinger    Dan Garrette
Google Research
{telmop,eschling,dhgarrette}@google.com

**On the Cross-lingual Transferability of Monolingual Representations**

Mikel Artetxe[1*], Sebastian Ruder[2] , Dani Yogatama[2]
[1]HiTZ Center, University of the Basque Country (UPV/EHU)
[2]DeepMind
mikel.artetxe@ehu.eus
{ruder,dyogatama}@google.com

## Multilingual BERT

- This works surprisingly well!

**How multilingual is Multilingual BERT?**

**What makes multilingual BERT multilingual?**

Telmo Pires*    Eva Schlinger    Dan Garrette
Google Research
{telmop,eschling,dhgarrette}

**Emerging Cross-lingual Structure in Pretrained Language Models**

**On the Cross-lingual Transferability of Mono**

Shijie Wu[♠*]   Alexis Conneau[♢*]
Haoran Li[♢]   Luke Zettlemoyer[♢]   Veselin Stoyanov[♢]

Mikel Artetxe[1*], Sebastian Ruder[2] , Dani Yogatama[2]
[1]HiTZ Center, University of the Basque Country (UPV/EHU)
[2]DeepMind
mikel.artetxe@ehu.eus
{ruder,dyogatama}@google.com

## Non-English Tokenizers

- We can use standard BPE tokenizers
  - massively increase vocab size (50k ➡ 250k+)

## Non-English Tokenizers

- We can use standard BPE tokenizers
  - massively increase vocab size (50k ➡ 250k+)

- Or use unicode byte-level models



---

## Data Resampling

### Problem: training data highly imbalanced

Data distribution over language pairs



High Resource ←→ Low Resource

(French, German, Spanish, ...)          (Yoruba, Sindhi, Hawaiian, ...)

→ High resource languages have much more data than low-resource ones

→ Important to upsample low-resource data in this setting!

Credit: Graham Neubig

---

## Data Resampling

### Problem: training data highly imbalanced



Sampling Probability

→ Sample data based on dataset size scaled by a temperature term

→ Easy control of how much to upsample low-resource data

Credit: Graham Neubig

---

## Translation MLM

- If I have translation data, can use it to enhance masked LM training



---

## Existing Multilingual Language Models

| Model | Architecture | Parameters | # languages | Data source |
|---|---|---|---|---|
| mBERT (Devlin, 2018) | Encoder-only | 110M | 104 | Wikipedia |
| XLM (Lample and Conneau, 2019) | Encoder-only | 570M | 100 | Wikipedia |
| XLM-R (Conneau et al., 2019) | Encoder-only | 270M − 550M | 100 | Common Crawl (CCNet) |
| mBART (Lewis et al., 2019a) | Encoder-decoder | 680M | 25 | Common Crawl (CC25) |
| MARGE (Lewis et al., 2020) | Encoder-decoder | 960M | 26 | Wikipedia or CC-News |
| mT5 (ours) | Encoder-decoder | 300M − 13B | 101 | Common Crawl (mC4) |

---

## Multilingual Few-shot Learning

- Can use LMs out of the box for few-shot learning in different languages

## Multilingual Few-shot Learning

- Can use LMs out of the box for few-shot learning in different languages



## Cross-lingual Supervised Transfer

- If I have supervised data, can transfer to languages w/o supervised data



## Cross-lingual Supervised Transfer

- If I have supervised data, can transfer to languages w/o supervised data



## Cross-lingual Supervised Transfer

- If I have supervised data, can transfer to languages w/o supervised data