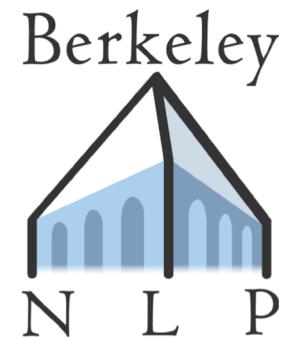


Multilingual Language Models: NLP Beyond English



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 computational expensive
- Translation is fundamentally lossy?



Monolingual LMs

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

The screenshot shows the Hugging Face Model Hub interface for the "gpt-2-spanish" model. At the top, there's a navigation bar with links for "Text Generation", "PyTorch", "JAX", "TensorBoard", "Safetensors", "Transformers", "oscar", "Spanish", and "gpt2". Below the navigation bar, there are tabs for "Model card", "Files and versions", "Training metrics", and "Community". The "Model card" tab is currently selected. On the right side, there are buttons for "Edit model card", "Train", "Deploy", and "Use in Transformers". A chart shows "Downloads last month" at 2,004. Below the chart, there's a section for "Hosted inference API" with a "Text Generation" button and an "Examples" dropdown.

flax-community/gpt-2-spanish □ like 13

Text Generation PyTorch JAX TensorBoard Safetensors Transformers oscar Spanish gpt2

Model card Files and versions Training metrics Community Edit model card

Downloads last month
2,004

Hosted inference API Examples

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.



Monolingual LMs

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

bert-base-chinese like 256

Fill-Mask PyTorch TensorFlow JAX Safetensors Transformers Chinese bert AutoTrain Compatible arxiv:1810.04805

Model card Files and versions Community 8 Edit model card

Bert-base-chinese Downloads last month 1,387,599





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 inayopiga risasi. Inadumu zaidi kuliko silaha ya chuma.	Telephone	Neutral
Russian	И мы занимаемся этим уже на протяжении 85 лет. Мы только начали этим заниматься.	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。	Slate	Contradiction
Arabic	تحتاج الوكالات لأن تكون قادرة على قياس مستويات النجاح . لا يمكن للوكالات أن تعرف ما إذا كانت ناجحة أم لا	Nine-Eleven	Contradiction



Few-shot Learning in Other Languages

ロン・ポールの学部時代の専攻は？ [Japanese]
(What did **Ron Paul** major in during undergraduate?)



Multilingual document collections
(Wikipedias)

ロン・ポール (ja.wikipedia)

高校卒業後はゲティスバーグ大学へ進学。
(After high school, he went to Gettysburg College.)

Ron Paul (en.wikipedia)

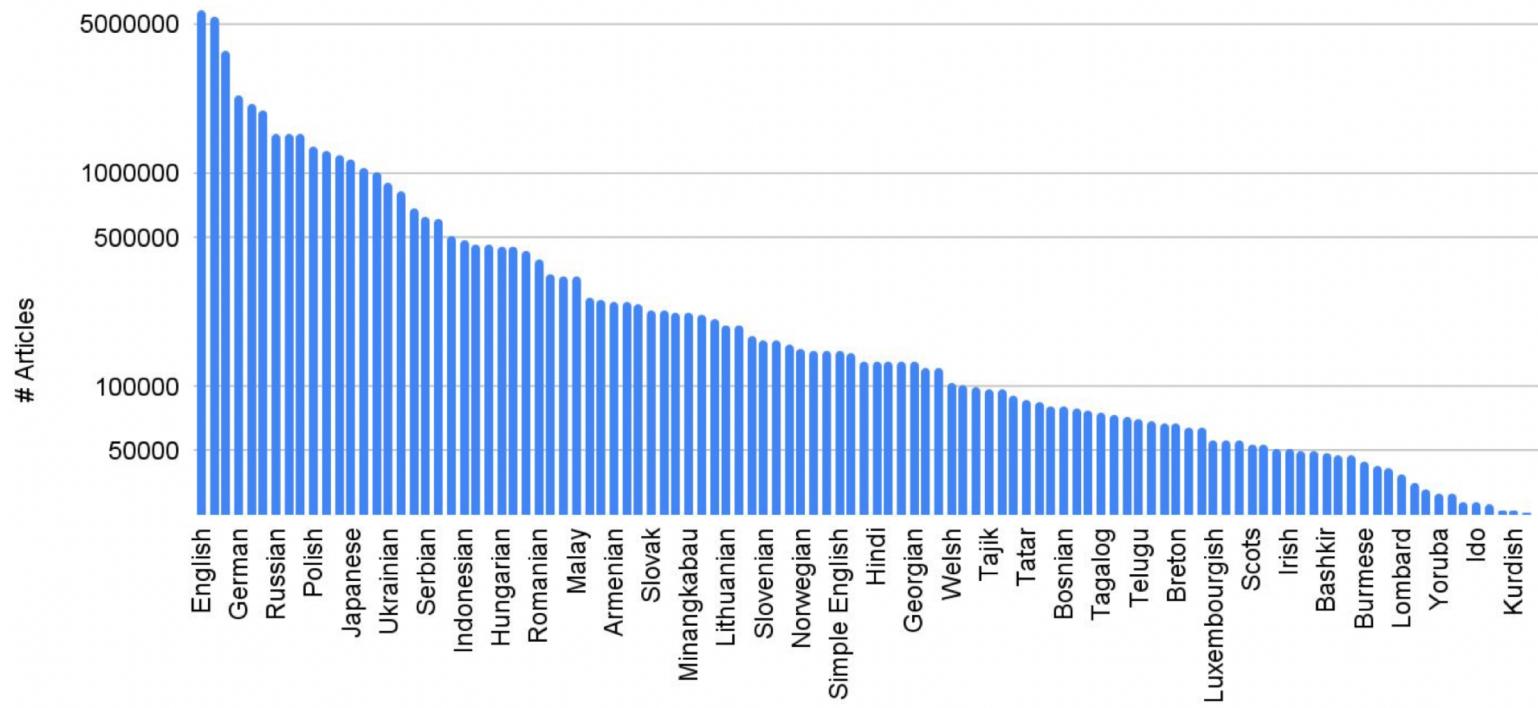
Paul went to Gettysburg College, where he was a member of the Lambda Chi Alpha fraternity. He graduated with a B.S. degree in **Biology** in 1957.

生物学 (Biology)



Challenges with Monolingual LMs

- There is not enough unlabeled data for each language

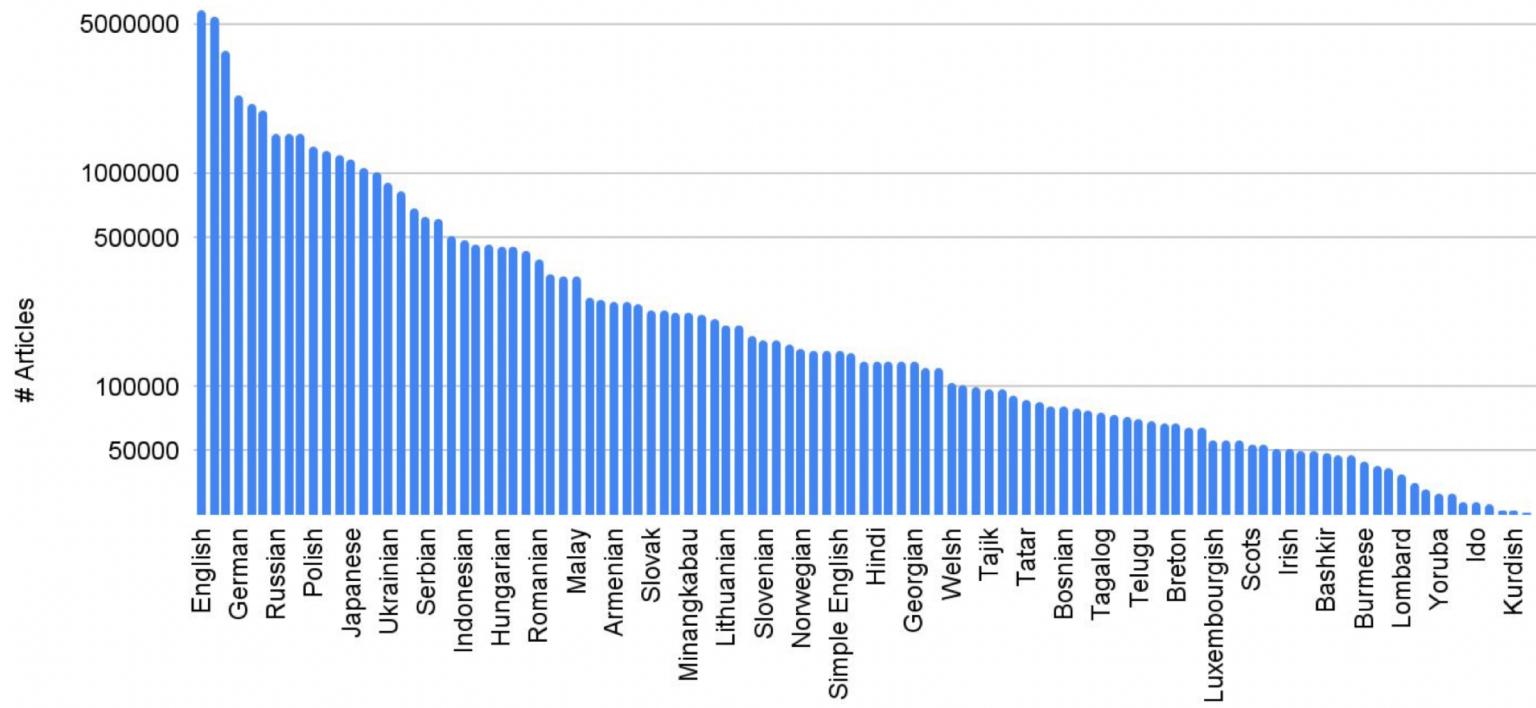


Credit: Graham Neubig



Challenges with Monolingual LMs

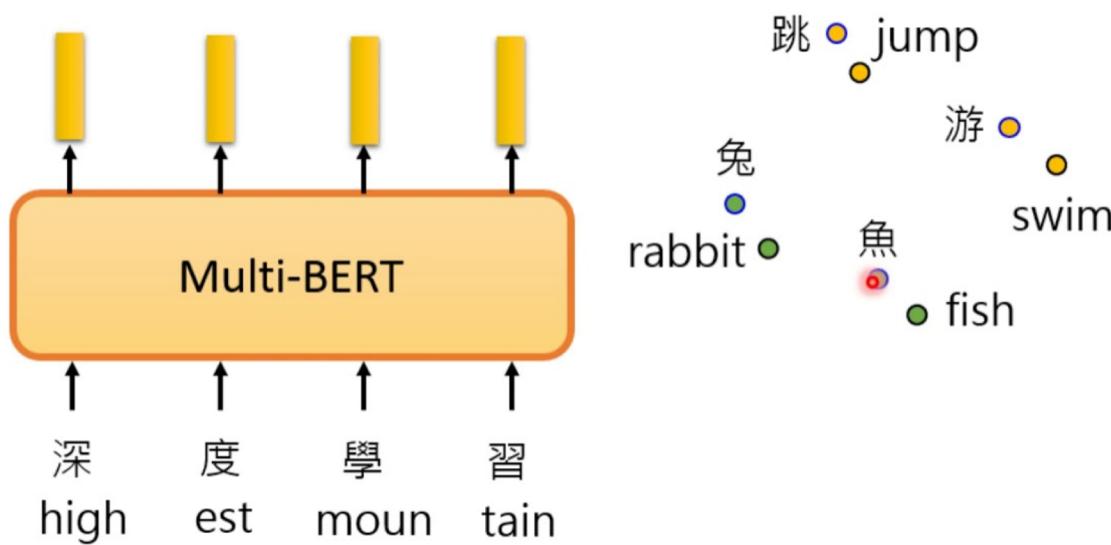
- ▶ Compute and complexity of serving 100-1000s of different models



Credit: Graham Neubig



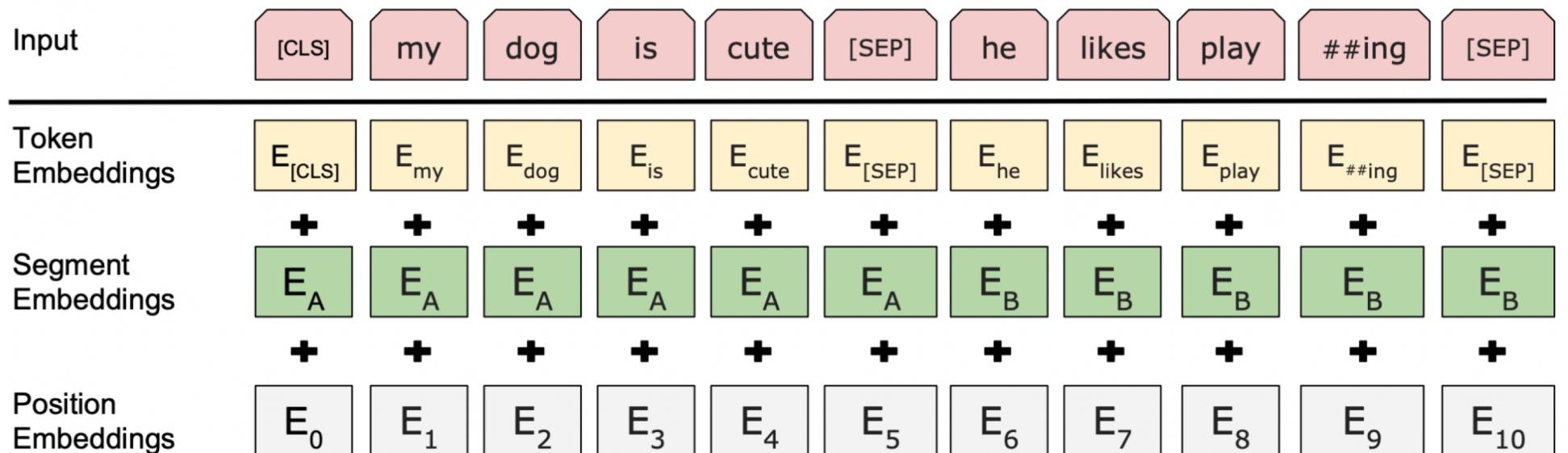
Multilingual Language Models?





Multilingual BERT

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





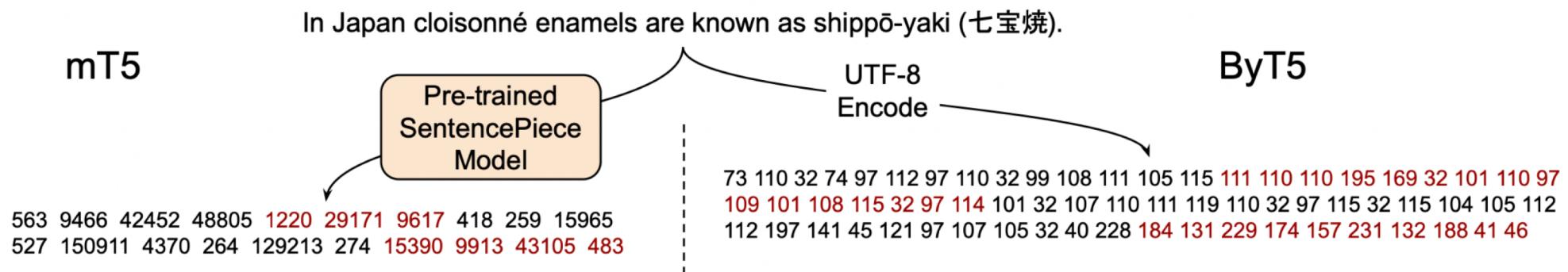
Non-English Tokenizers

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



Non-English Tokenizers

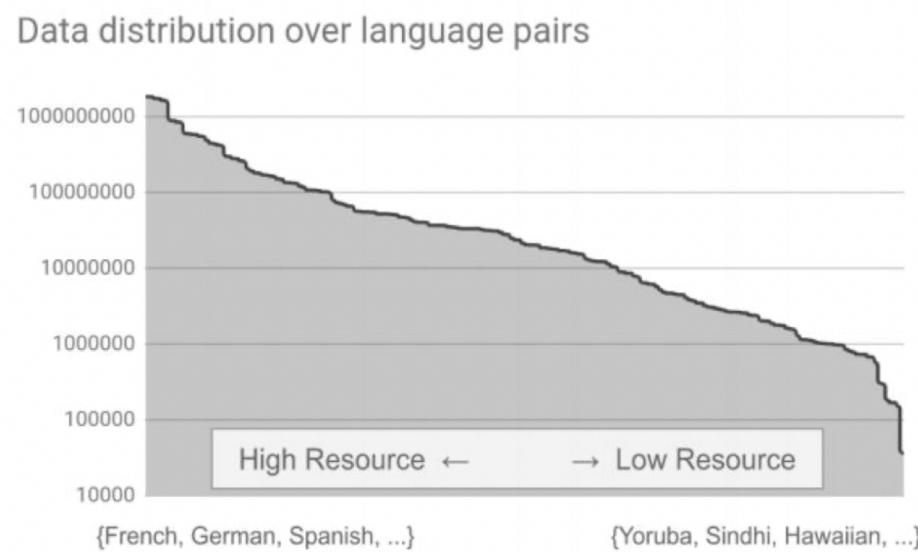
- We can use either standard BPE tokenizers or byte-level models
 - massively increase BPE size (50k → 250k+)





Multilingual Language Models

Problem: training data highly imbalanced



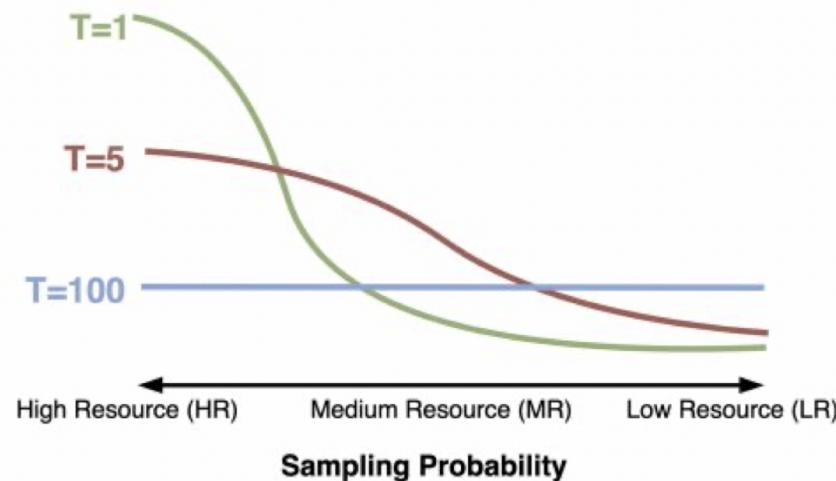
- High resource languages have much more data than low-resource ones
- Important to upsample low-resource data in this setting!

Credit: Graham Neubig



Multilingual Language Models

Problem: training data highly imbalanced



- Sample data based on dataset size scaled by a temperature term
- Easy control of how much to upsample low-resource data

Credit: Graham Neubig



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

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 inayopiga risasi. Inadumu zaidi kuliko silaha ya chuma.	Telephone	Neutral
Russian	И мы занимаемся этим уже на протяжении 85 лет. Мы только начали этим заниматься.	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。	Slate	Contradiction
Arabic	تحتاج الوكالات لأن تكون قادرة على قياس مستويات النجاح . لا يمكن للوكالات أن تعرف ما إذا كانت ناجحة أم لا	Nine-Eleven	Contradiction



Cross-lingual Transfer

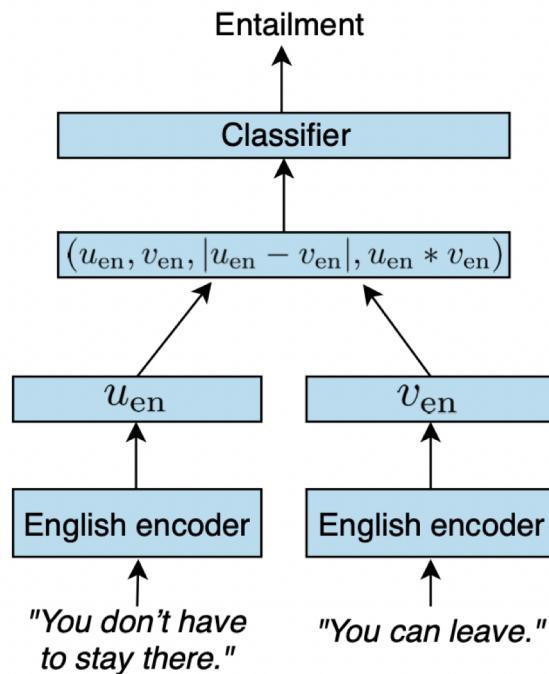
- ▶ If I have English supervised data, can transfer to low-resource language



Cross-lingual Transfer

- ▶ If I have English supervised data, can transfer to low-resource language

A) Learning NLI English encoder and classifier

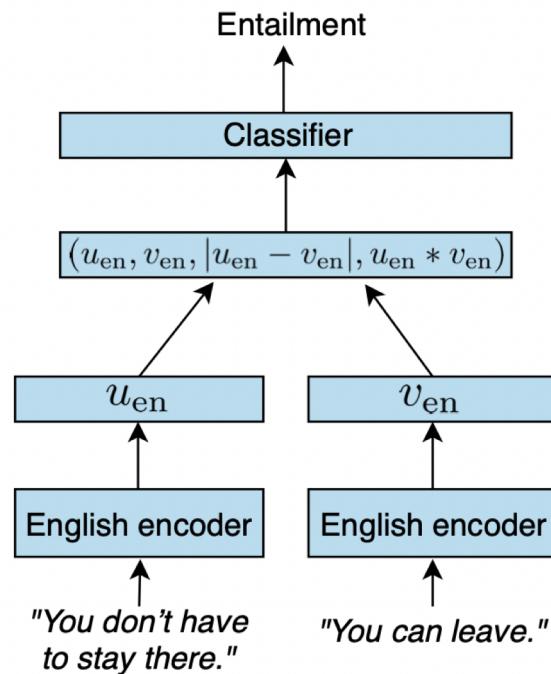




Cross-lingual Transfer

- If I have English supervised data, can transfer to low-resource language

A) Learning NLI English encoder and classifier



C) Inference in the other language

