

CS 224S / Linguist 285

Spoken Language Processing

Tolúlopé Ògúnremí | Stanford University | Spring 2024

Lecture 12: Speech Recognition Beyond English

Sound check

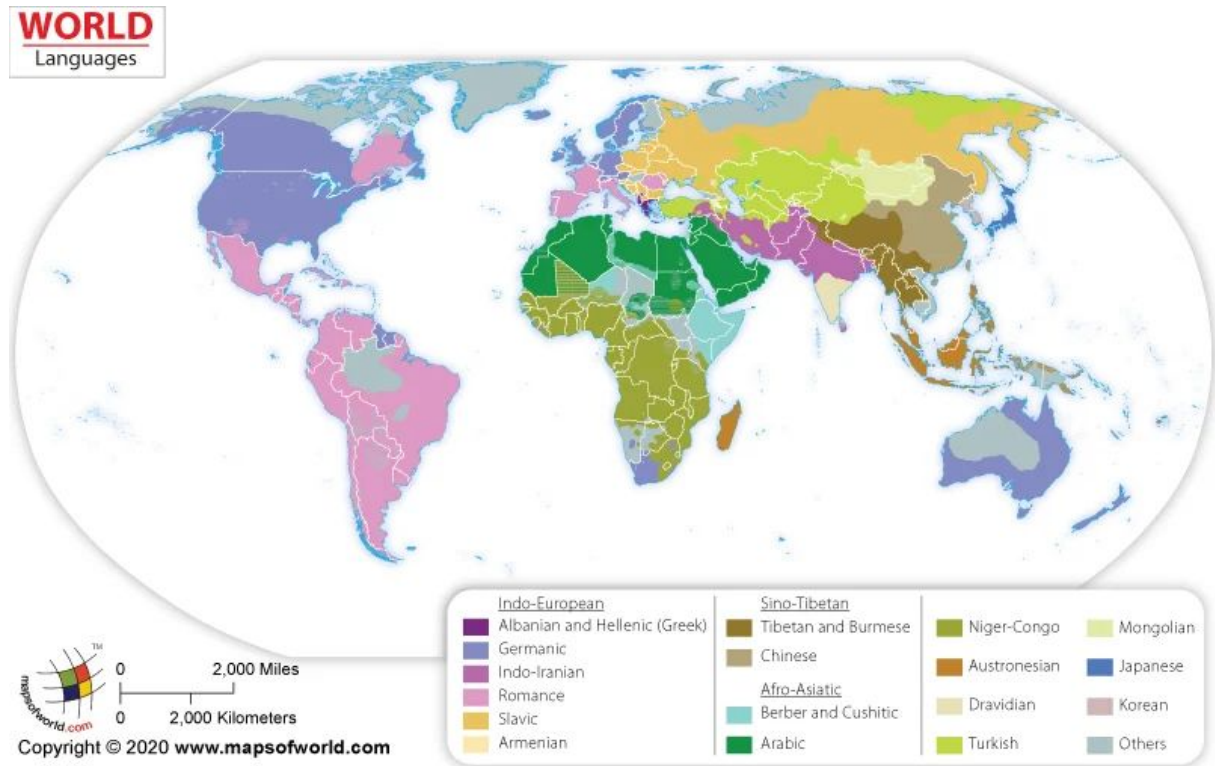


Project Check-Ins

Outline

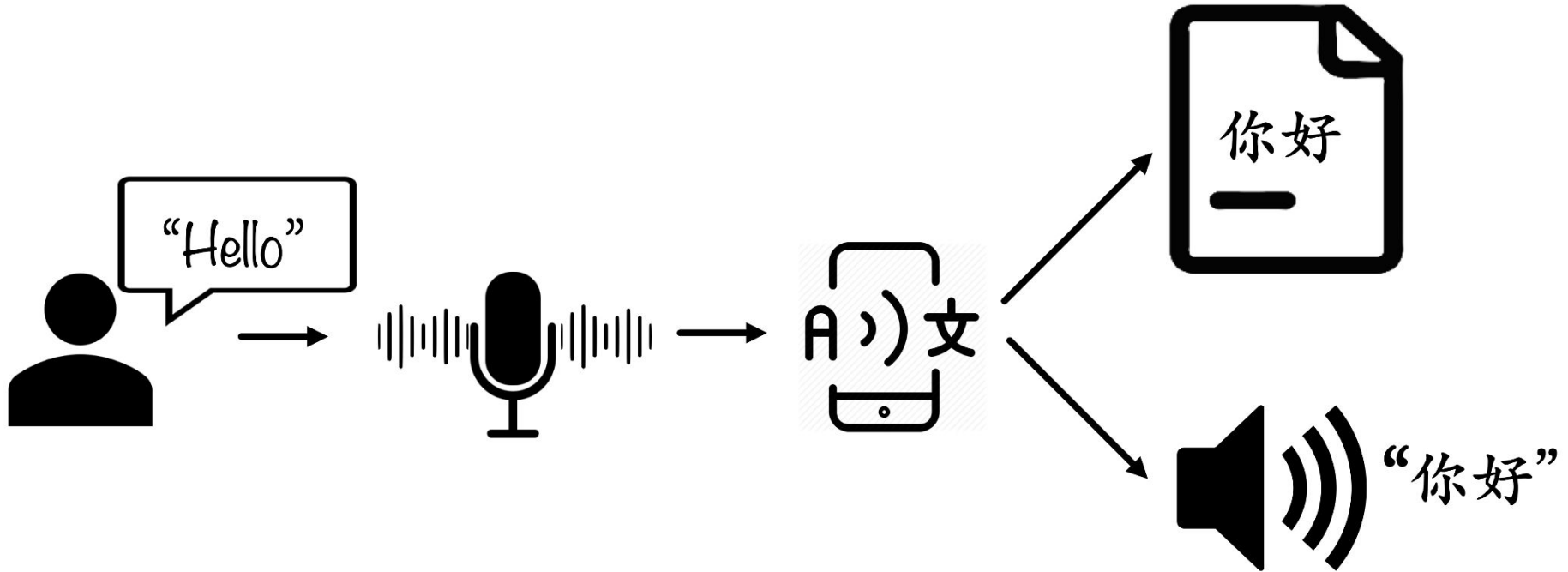
- How languages can differ from English
- Multilingual large pre-trained models
- Datasets and Benchmarks
- Language-specific ASR techniques

There are over
7,000 known
languages in
the world.



We need to process (as many of) the languages of the world (as we can).

Example: Speech Translation



Most of the models we have seen in this class have been trained with only English data.

Languages vary

Languages can have different scripts

Writing system	Scripts	
Alphabet	Roman	napenda utambuzi wa hotuba
	Greek	Λατρεύω την αναγνώριση ομιλίας
	Cyrillic	Би яриа таних дуртай
	Korean	나는 음성 인식을 좋아해요
Semanto-Phonetic	Chinese	我喜欢语音识别
Syllabic Alphabet	Devanāgarī	मलाई बोली पहिचान मन पछ
	Thai	ฉันชอบการรู้จำคำพูด
	Tamil	நான் பேச்சு அங்கீகாரத்தை விரும்புகிறேன்
Abjad	Arabic	أنا أحب التعرف على الكلام
	Hebrew	אני אוהב זיהוי דיבור

Lecture 12:

Speech Recognition Beyond English

Adapted from [Tan et al.](#), 2010

Languages can have **lexical** tone

The pitch of the word changes the
meaning of the word

wá



wò



wà



wù



Languages can have different dialects

English	I don't know what to do
Jordanian Arabic	مش عارف شو اعمل
Palestinian Arabic	شو بدني اعمل
Emirati Arabic	معرف شو اسوي
Modern Arabic	لا اعلم ماذا افعل
Egyptian Arabic	مش عارف اعمل ايه
Tunisian Arabic	منعرفش
Algerian Arabic	ما على بالي
Kuwaiti Arabic	ما ادري شو اسوي

Image from [Bani-Hani et al.](#)
2017

Languages can have codeswitching

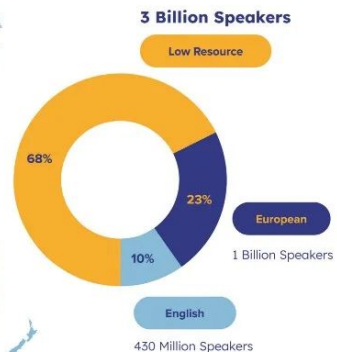


Languages can have little data available to train models

NLP Solutions by Language



Population Size of Languages



Multilingual large pretrained speech models

Multilingual versions of English-only models: wav2vec 2.0 XLSR

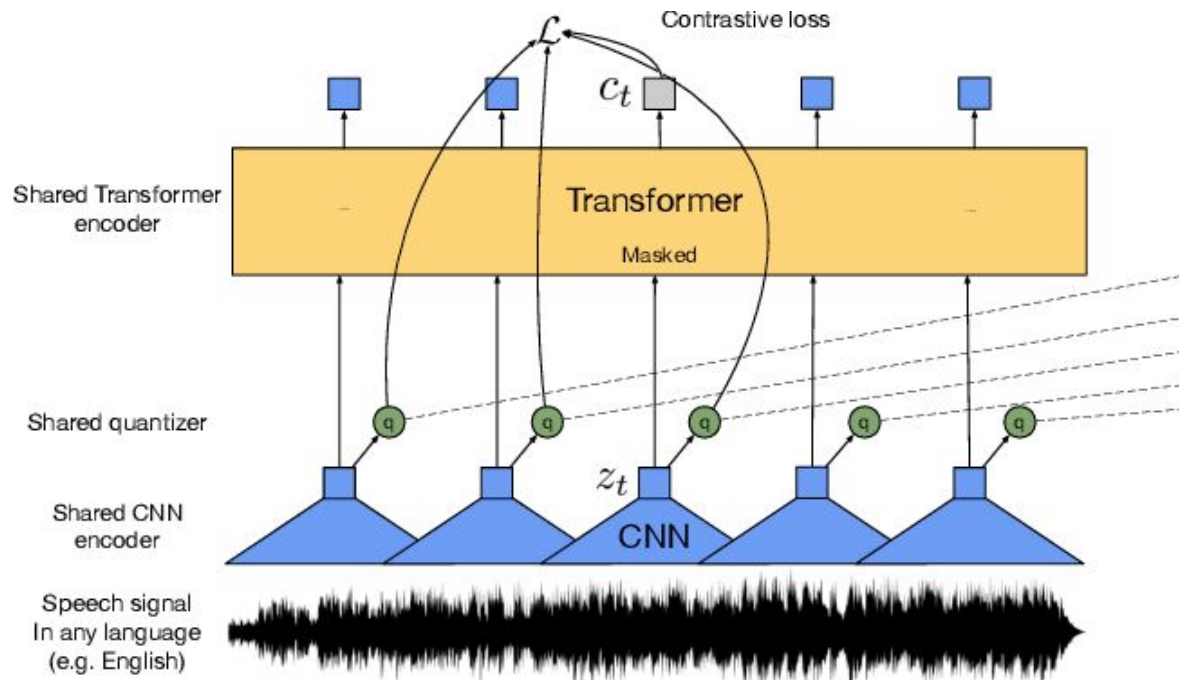


Image from [Conneau et.al, 2020](#)

- Trained on Multilingual LibriSpeech, Common Voice and BABEL
- (56,000 hours)
- 53 languages: XLSR-53

Multilingual versions of English-only models: wav2vec 2.0 XLSR

Multilingual quantized latent speech representations

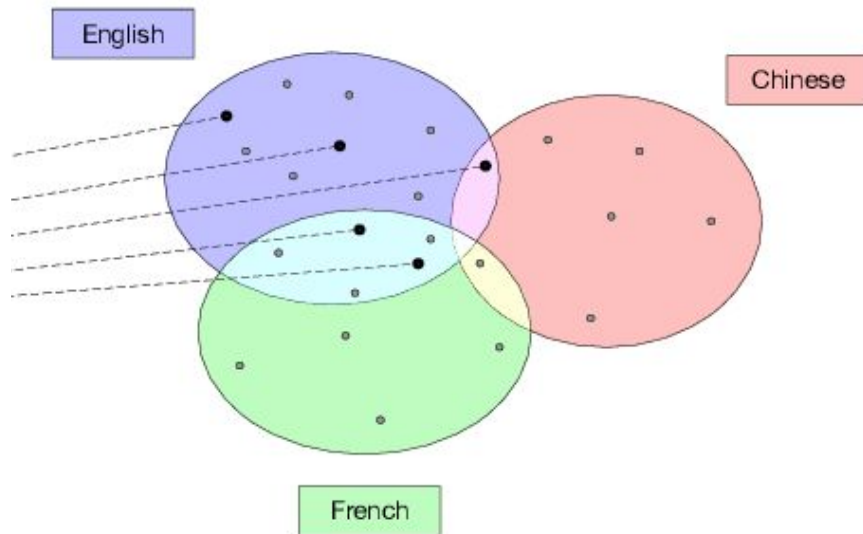
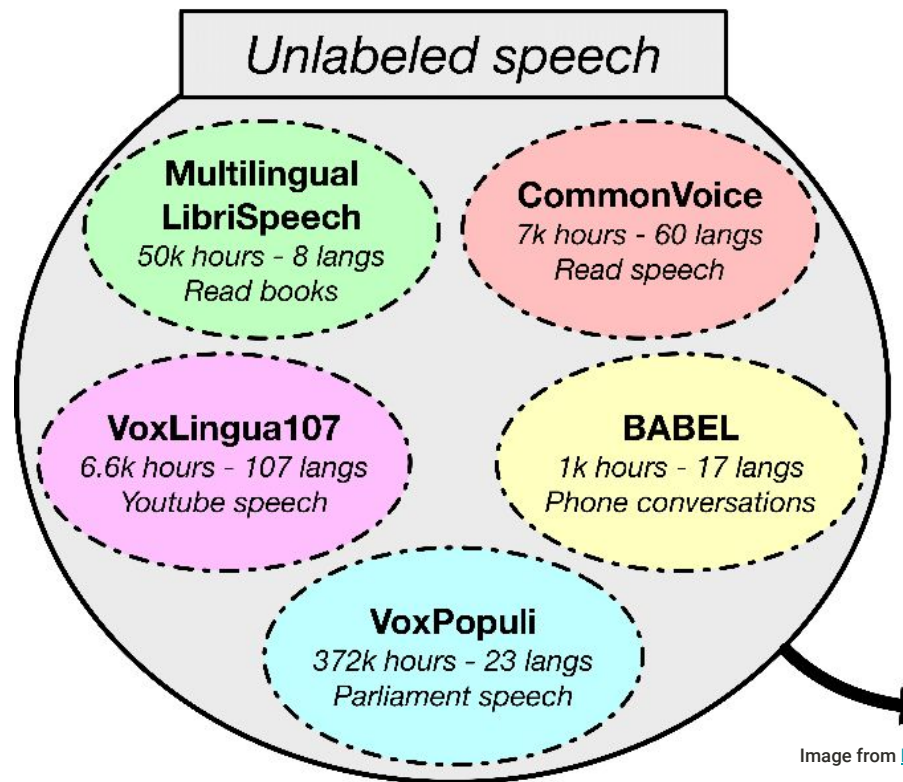


Image from [Conneau et al., 2020](#)

- Latent multilingual speech representations are theorised

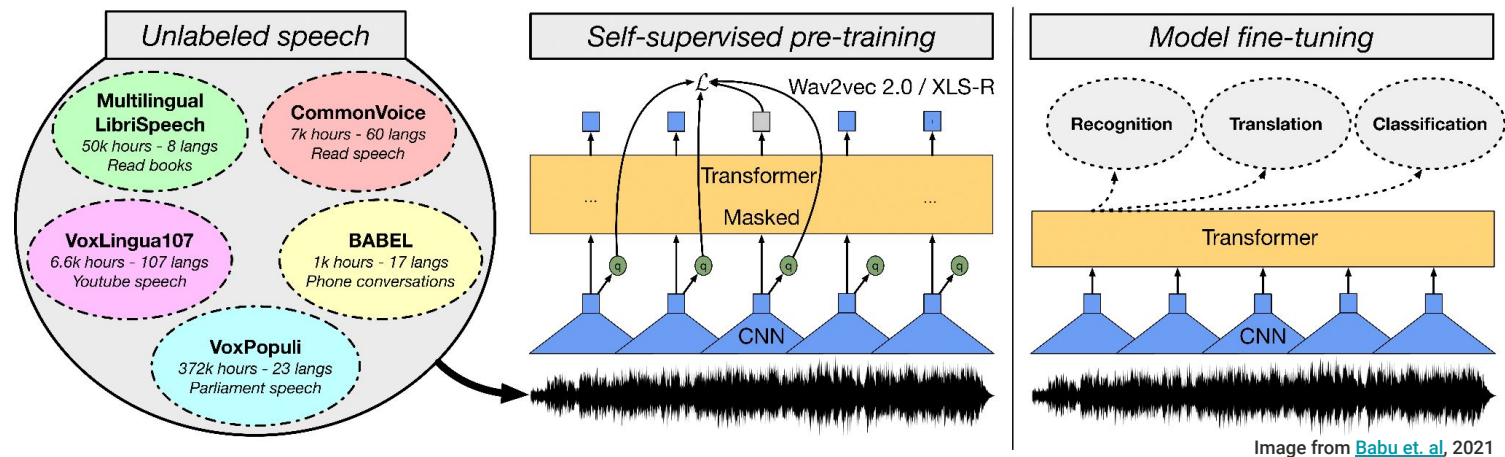
Multilingual versions of English-only models: wav2vec 2.0 XLS-R



- Trained on XLSR datasets and Vox Lingua 107 and Vox Populi, totalling 436,000 hours

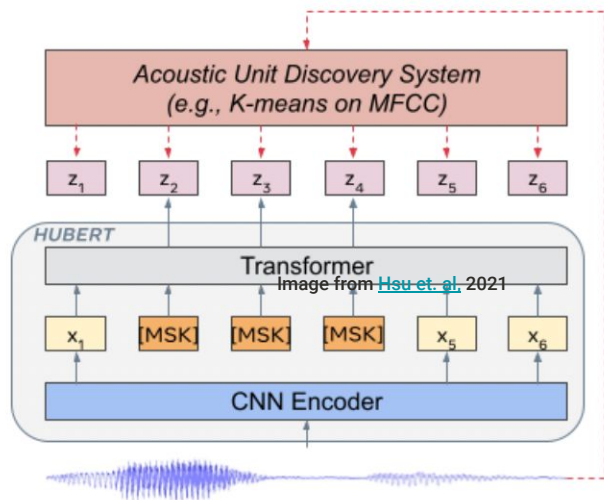
Image from [Babu et. al, 2021](#)

Multilingual versions of English-only models: wav2vec 2.0 XLS-R



- Tested on ASR and AST (Automatic Speech Translation)

Multilingual versions of English-only models: mHuBERT



- Trained specifically for speech translation in “[Textless Speech-to-Speech Translation on Real Data](#)” (Lee et. al, 2022)
- Trained with the 100,000 hour subset of Vox Populi

Multilingual from the start: Whisper

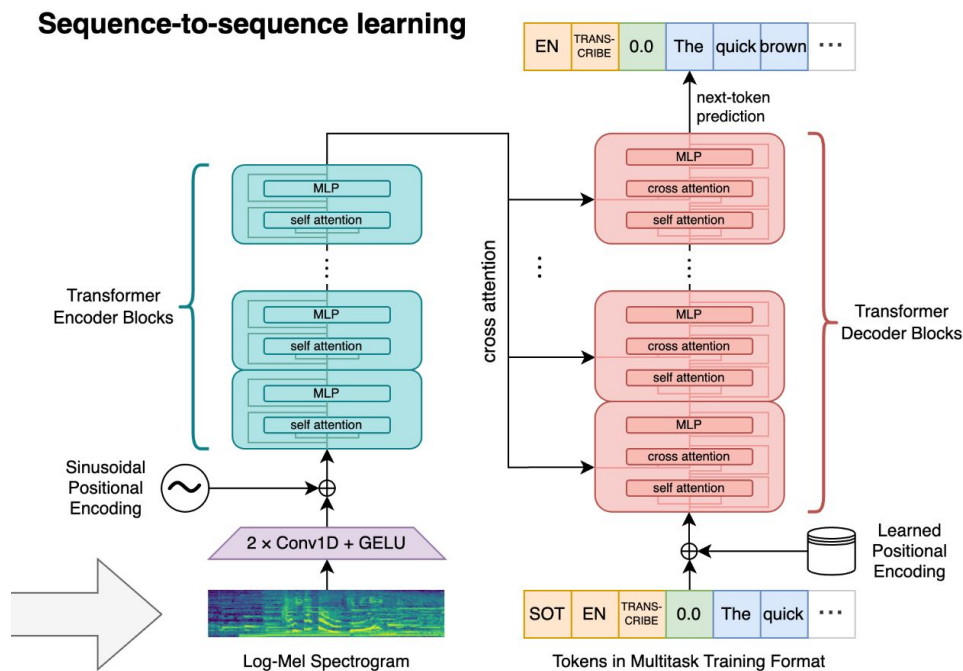


Image from [Radford et. al](#), 2022

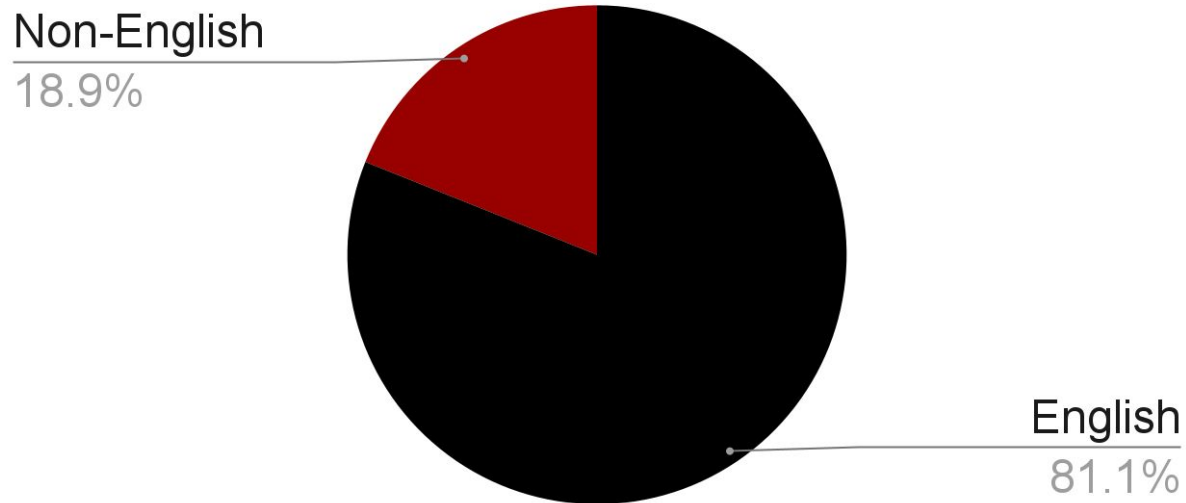
- “Multilingual and multitask”
- Trained with 680,000 hours of data
- Training data is not publicly available.

What is the data distribution?

How multilingual are these models?

wav2vec 2.0
XLSR

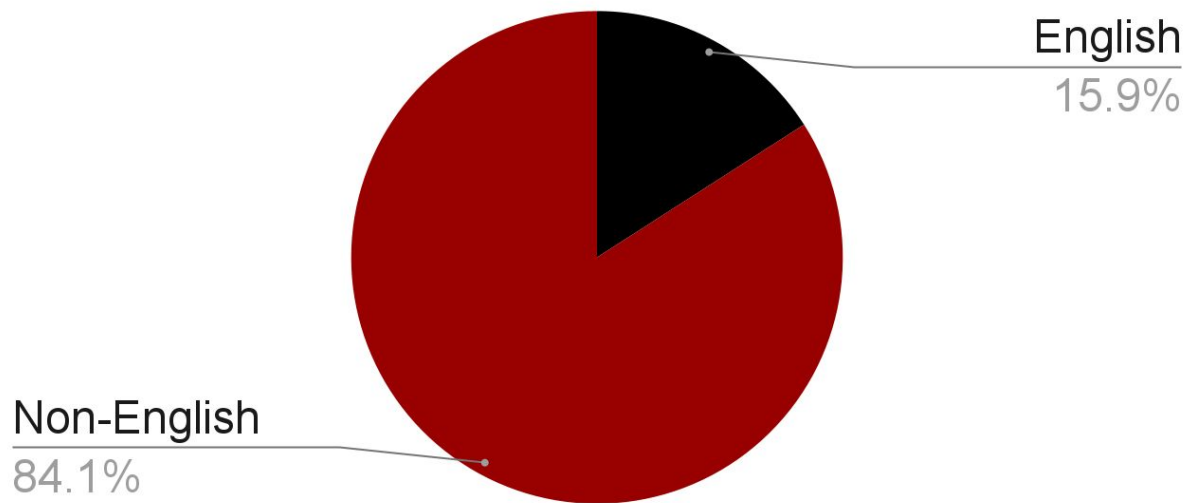
Languages in wav2vec 2.0 XLSR



How multilingual are these models?

wav2vec 2.0
XLS-R

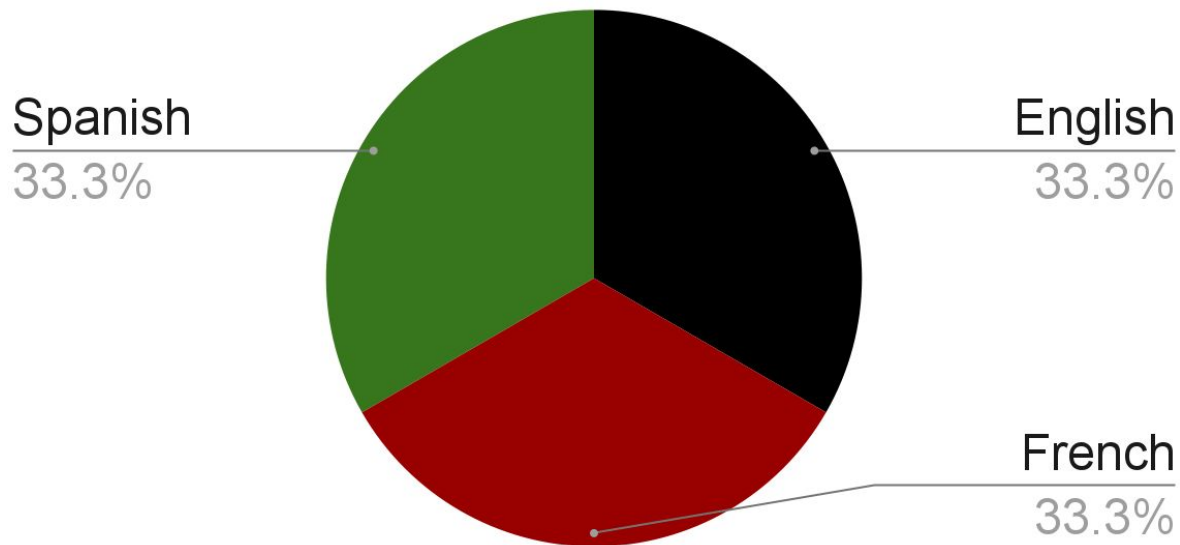
Languages in wav2vec 2.0 XLS-R



How multilingual are these models?

mHuBERT

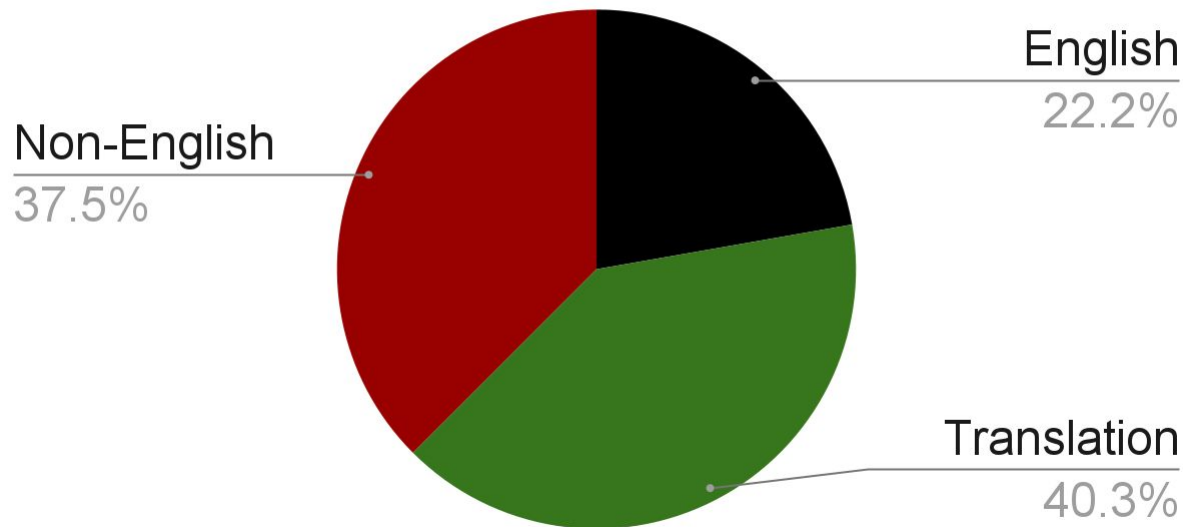
Languages in mHuBERT



How multilingual are these models?

Whisper

Languages in Whisper



Open Source Multilingual Datasets: CommonVoice and Yodas

Common Voice

English 🕒 Hours: 3485 🗣️ Speakers: 93140 📊 Validation Progress: 73% 📄 Sentences: 1673433 CONTRIBUTE	Catalan 🕒 Hours: 3838 🗣️ Speakers: 35957 📊 Validation Progress: 73% 📄 Sentences: 1214229 COL-LABOREU-HI	Kinyarwanda 🕒 Hours: 2388 🗣️ Speakers: 1134 📊 Validation Progress: 84% 📄 Sentences: 1404853 FASHA, TANGA UMUSANZU
Belarusian 🕒 Hours: 1794 🗣️ Speakers: 8363 📊 Validation Progress: 97% 📄 Sentences: 379505 ЗРАБІЦЬ УНЁСАК	Esperanto 🕒 Hours: 1936 🗣️ Speakers: 1758 📊 Validation Progress: 75% 📄 Sentences: 180562 KONTRIBUI	German 🕒 Hours: 1412 🗣️ Speakers: 19151 📊 Validation Progress: 94% 📄 Sentences: 2056443 MITMACHEN
French 🕒 Hours: 1145 🗣️ Speakers: 19584 📊 Validation Progress: 88% 📄 Sentences: 1646292 CONTRIBUER	Kabyle 🕒 Hours: 699 🗣️ Speakers: 1560 📊 Validation Progress: 81% 📄 Sentences: 182716 TTEKKI	Spanish 🕒 Hours: 2326 🗣️ Speakers: 26107 📊 Validation Progress: 20% 📄 Sentences: 1080695 COLABORAR
Luganda 🕒 Hours: 583 🗣️ Speakers: 660 📊 Validation Progress: 76% 📄 Sentences: 191407 YAMBAKO	Swahili 🕒 Hours: 1081 🗣️ Speakers: 1459 📊 Validation Progress: 37% 📄 Sentences: 134669 CHANGIA	Persian 🕒 Hours: 426 🗣️ Speakers: 4453 📊 Validation Progress: 87% 📄 Sentences: 55667 مشارکت

- Multilingual living dataset
- 30,000 recorded hours covering 124 languages
- Anyone can set up a Common Voice page for their language
- Anyone can record utterances for the dataset
- Dataset is noisier than LibriSpeech due to less controlled recording environments

Common Voice

Marathi Common Voice

Kinyarwanda Common Voice

Attempting to open source datasets: Yodas

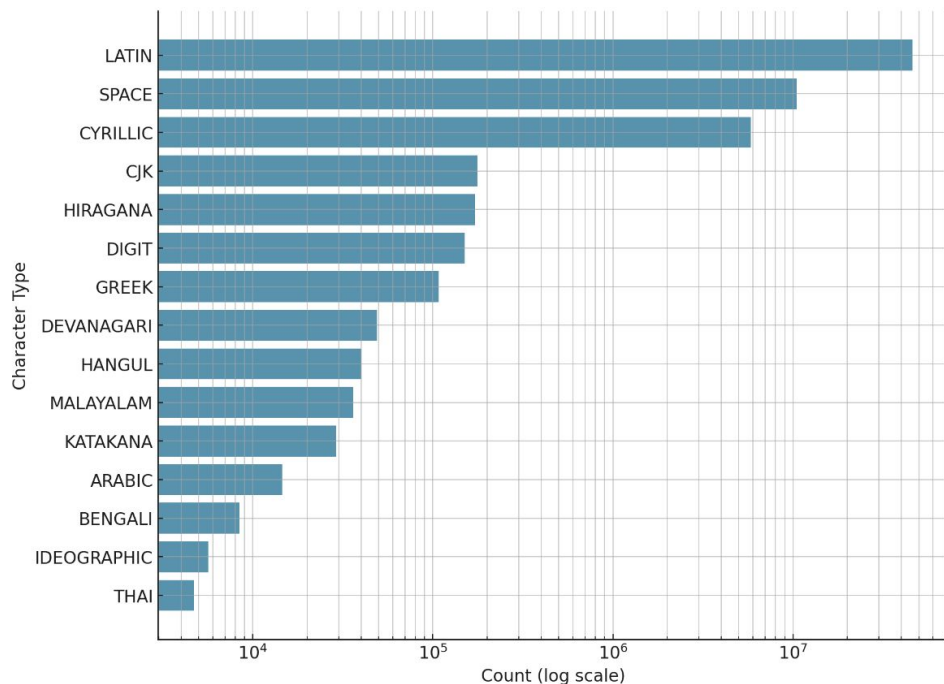


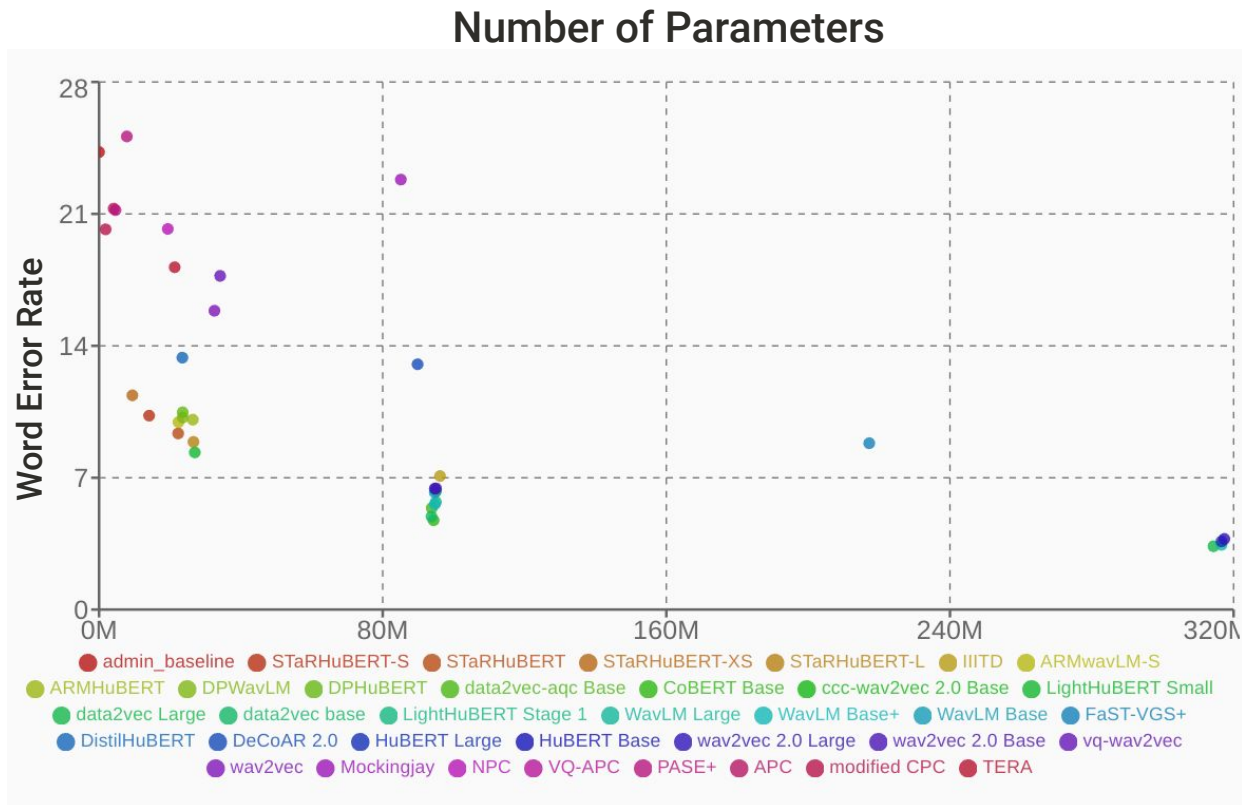
Image from [WAVLab post](#)

- [Youtube-Oriented Dataset for Audio and Speech](#)
- Result of a 6-month crawl of YouTube followed by alignment of transcript to audio.
- 500,000 hours of data across 140 languages.
- 420,000 hours of transcribed data

Benchmarking large models multilingually: ML-SUPERB

SUPERB

Speech processing Universal
PERformance Benchmark



ML-SUPERB

Multilingual Speech
processing Universal
PERformance Benchmark

Automatic speech
recognition and language
identification for 143
languages

Method	Mono-ASR ↓	Multi-ASR (Normal) ↓	Multi-ASR (Few-shot) ↓	LID ↑
HuBERT Base	35.3	31.4	42.7	86.1
HuBERT Large	32.2	37.7	43.5	64.1
Mandarin HuBER...	45.6	43.2	46.6	85.3
Mandarin HuBER...	33.7	39.6	45.1	57.3
Robust wav2vec 2...	35.7	31.1	42.2	72.1

Language-specific techniques

Languages can have different scripts

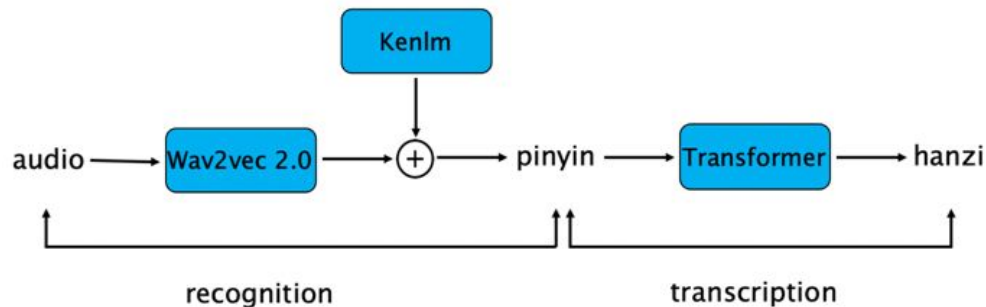
Writing system	Scripts	
Alphabet	Roman	napenda utambuzi wa hotuba
	Greek	Λατρεύω την αναγνώριση ομιλίας
	Cyrillic	Би яриа таних дуртай
	Korean	나는 음성 인식을 좋아해요
Semanto-Phonetic	Chinese	我喜欢语音识别
Syllabic Alphabet	Devanāgarī	मलाई बोली पहिचान मन पर्छ
	Thai	ฉันชอบการรู้จำคำพูด
	Tamil	நான் பேச்சு அங்கீகாரத்தை விரும்புகிறேன்
Abjad	Arabic	أنا أحب التعرف على الكلام
	Hebrew	אני אוהב זיהוי דיבור Adapted from Tan et. al , 2010

Lecture 12:

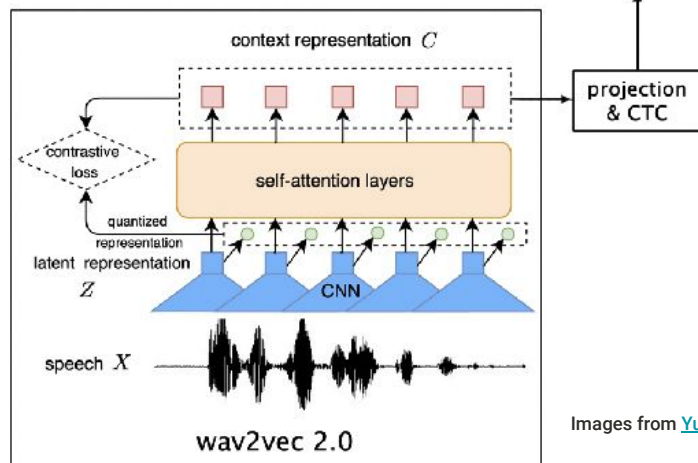
Speech Recognition Beyond English

Using different representations

Incorporating Pinyin for Mandarin Chinese - intermediary phonetic representation



Characters: 她的表现也更加全面
Pinyin+T: ta1 de5 biao3 xian4 ye3 geng4 jia1 quan2 mian4
Pinyin-T: ta de biao xian ye geng jia quan mian



Images from [Yuan et al., 2021](#)

Languages can have **lexical** **tone**

The pitch of the word changes the
meaning of the word

wá



wò



wà



wù



Tonal languages: can we find tones in the representations?

Shen et. al find that models behave similarly to native and non-native human participants in tone and consonant perception studies, but they do not follow the same developmental trajectory.

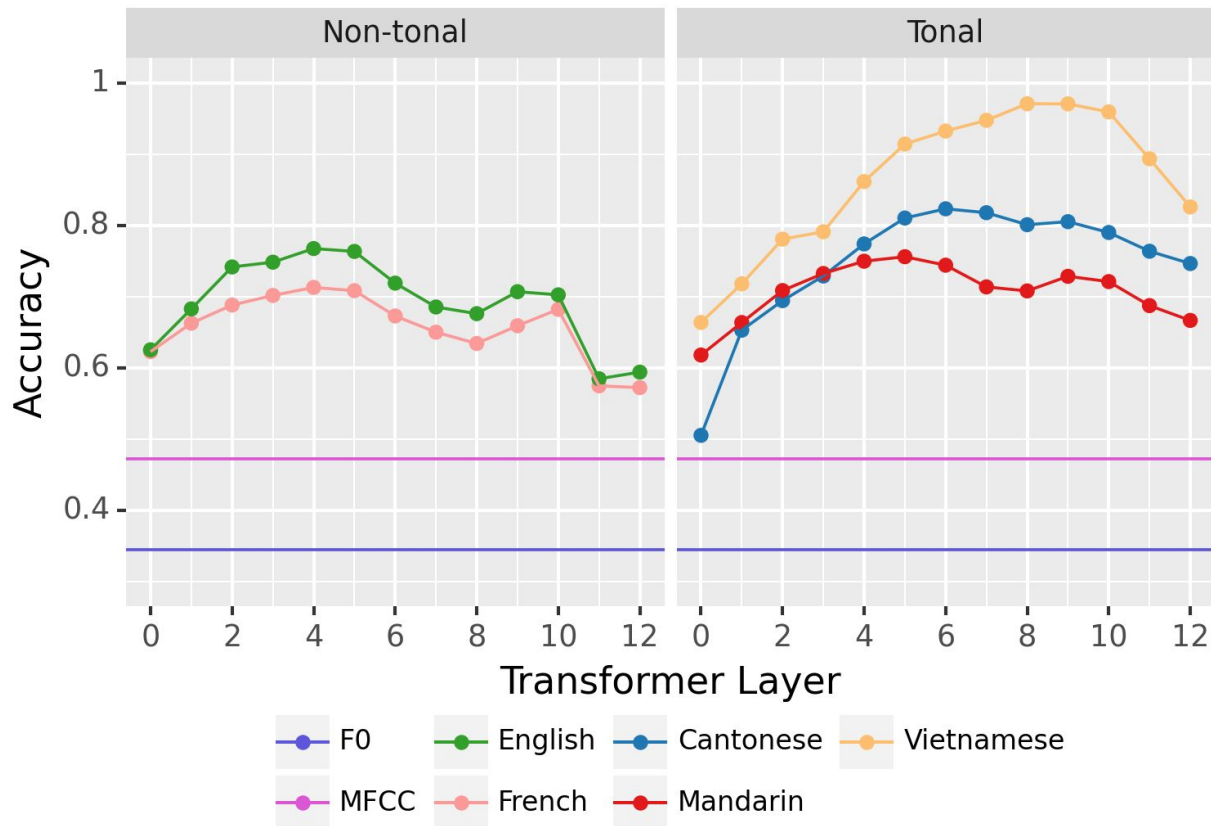


Image from [Shen et. al, 2024](#)

Languages can have different dialects

English	I don't know what to do
Jordanian Arabic	مش عارف شو اعمل
Palestinian Arabic	شو بدني اعمل
Emirati Arabic	معرف شو اسوي
Modern Arabic	لا اعلم ماذا افعل
Egyptian Arabic	مش عارف اعمل ايه
Tunisian Arabic	منعرفش
Algerian Arabic	ما على بالي
Kuwaiti Arabic	ما ادري شو اسوي

Image from [Bani-Hani et al.](#)
2017

Next week!



(a) XLSR-nl layer 15



(b) LD

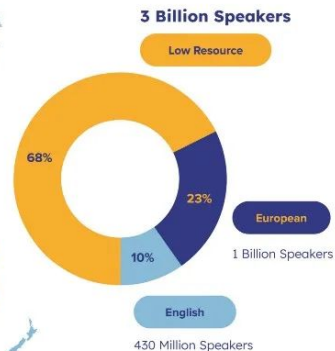
Image from [Bartelds & Wieling, 2022](#)

Languages can have little data available to train models

NLP Solutions by Language



Population Size of Languages

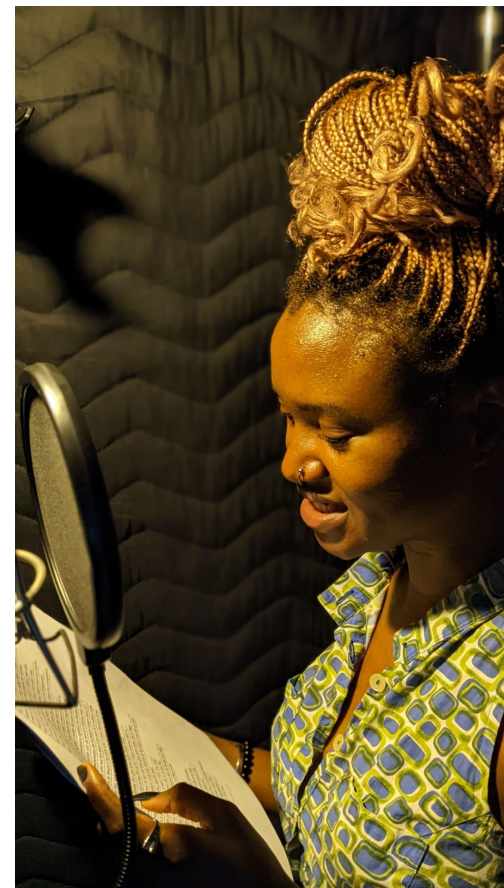


Making datasets

ÌròyìnSpeech: A multi-purpose Yorùbá Speech Corpus



Lecture 12:
Speech Recognition Beyond English



[Ògúnrèní et. al, 2024](#)

Making monolingual versions of large pretrained speech models

HUBERT-TR: REVIVING TURKISH AUTOMATIC SPEECH RECOGNITION WITH SELF-SUPERVISED SPEECH REPRESENTATION LEARNING

*Ali Safaya *, Engin Erzin*

KUIS AI Center
Computer Engineering Department
Koç University

Using Radio Archives for Low-Resource Speech Recognition: Towards an Intelligent Virtual Assistant for Illiterate Users

Moussa Doumbouya,¹ Lisa Einstein,^{1,2} Chris Piech²

¹ GNCode

² Stanford University

moussa@gncode.org, lisae@stanford.edu, piech@cs.stanford.edu

Different scripts leverage CTC - no need for huge language model



Languages can have codeswitching



What is code-switching?

The mixing of words, phrases and sentences from two distinct grammatical (sub) systems across sentence boundaries within the same speech event.

(Bokomba, 1988)

I'll tell you exactly when I have to leave, at ten o'clock. **Y son las nueve y cuarto.**

Off-the-shelf multilingual models don't work well in this scenario.

[Ògúnrimí et. al, 2023](#)

How can we improve the performance of large multilingual models on code-switched data?

[Ògúnrémí et. al, 2023](#)

Data: South African Soap Opera Clips



Four South-African
languages
code-switched
with English

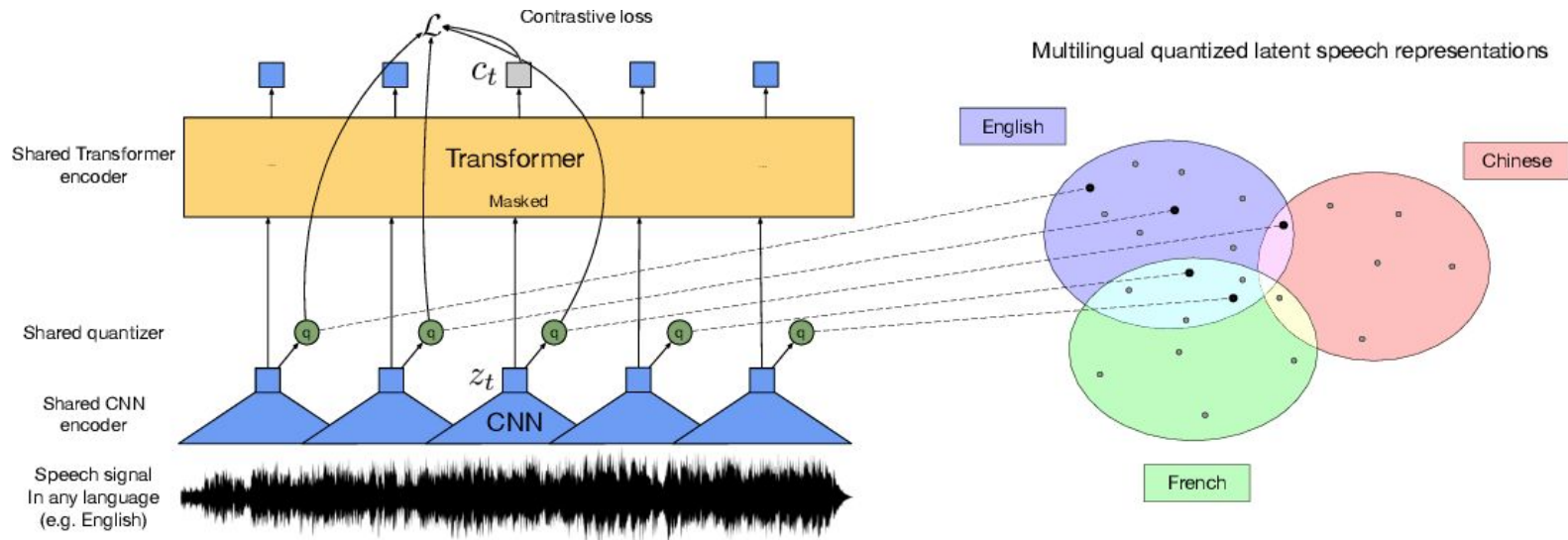
3 - 6 hours per
language

[Ògúnremí et. al](#), 2023

Data: South African Soap Opera Clips

Lang par	Train	Dev	Test	Total
Eng-Zul	4.81h	0.13h	0.51h	5.45h
Eng-Xho	2.68h	0.23h	0.23h	3.14h
Eng-Tsn	2.33h	0.23h	0.30h	2.86h
Eng-Sot	2.36h	0.21h	0.26h	2.83h

Model: wav2vec 2.0 XLSR



[Ògúnremí et. al, 2023](#)

Does incorporating language information help?

what if etholwa amaphoyisa kuqala

<eng> what if </eng> <zul> etholwa amaphoyisa kuqala
</zul>

TAGS

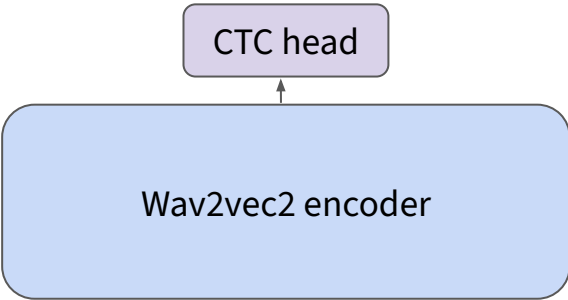
WHAT IF etholwa amaphoyisa kuqala

CASING

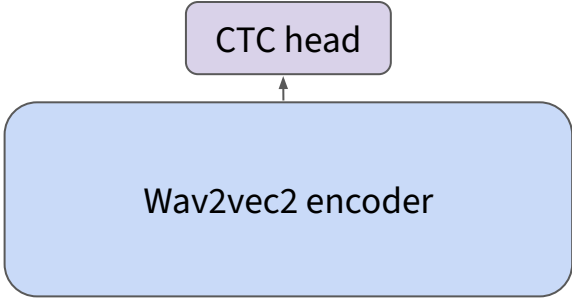
[Ògúnrèní et. al, 2023](#)

We fine-tune (with a CTC head) first on the language pair along with additional data, then on the language pair itself

Step 1



Step 2



language +
pair

A: monolingual data
B: the rest of the soap opera
corpus

language
pair

We find that finetuning with utterances in the same domain (soap opera data) but different, neighbouring languages improve performance over finetuning a single language pair.

[Ògúnremí et. al, 2023](#)

Language is varied

wá

wà

napenda utambuzi wa hotuba

Λατρεύω την αναγνώριση ομιλίας

나는 음성 인식을 좋아해요

मलाई बोली पहिचान मन पर्छ

ฉันชอบการรู้จำคำพูด

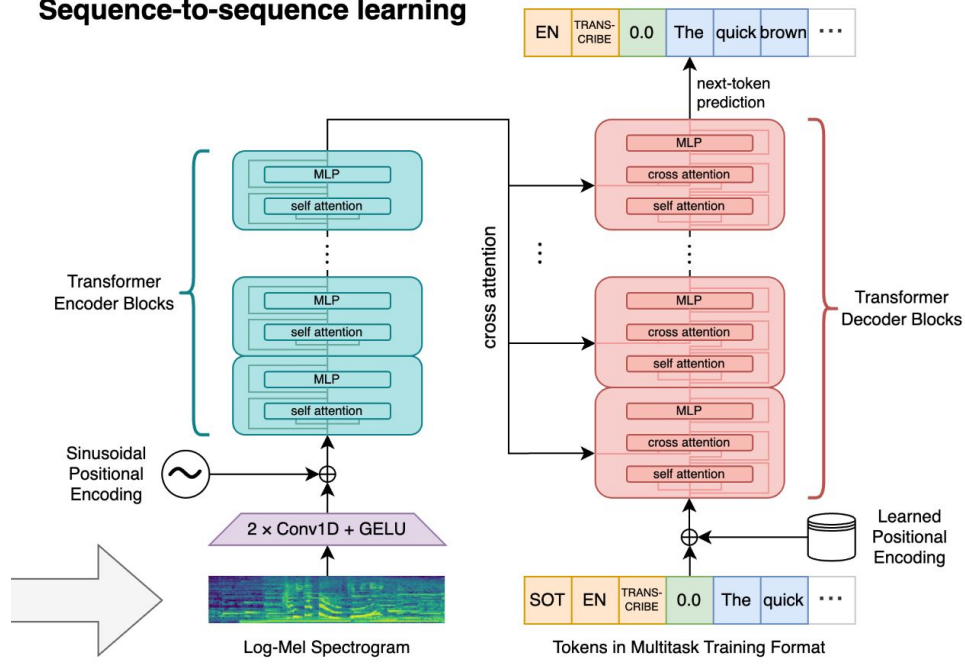
நான் பேச்சு அங்கீகாரத்தை
விரும்புகிறேன்

أنا أحب التعرف على الكلام



Surprisingly, all you need to do is chuck a bunch of data into a model.

Sequence-to-sequence learning



Thank You