# mHuBERT-147: A Compact Multilingual HuBERT Model

Marcely Zanon Boito, Vivek Iyer, Nikolaos Lagos, Laurent Besacier, Ioan Calapodescu

07/2024

Contact: marcely.zanon-boito@naverlabs.com

**NAVER LABS** 

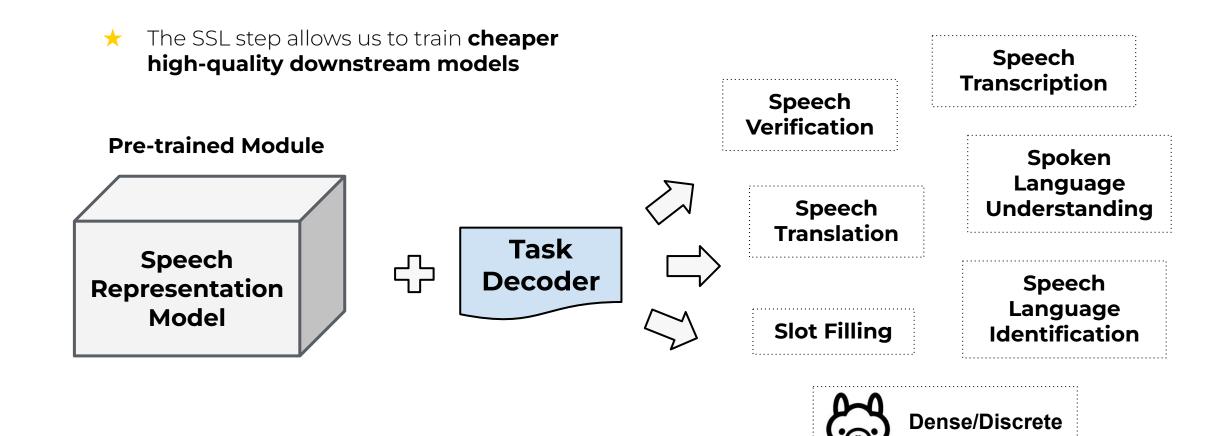


## Speech Representation Learning: learning contextualized high-dimensional speech embeddings for downstream tasks



**LLM Integration** 

#### Application: Any speech-to-something application



#### **Multilingual Speech Representation Models:**

- ★ One single back-bone for (multilingual) speech applications
- ★ Zero-shot unseen languages

Model	Training Approach	# Languages	# Datasets	# Hours
<b>XLSR-53</b> (Conneau et al. 2020)	wav2vec 2.0	53	3	56K
<b>XLS-R</b> (Babu et al. 2021)	wav2vec 2.0	128	5	436K
MMS (Pratap et al. 2023) [SOTA]	wav2vec 2.0	1,362	6	491K
<b>WavLabLM</b> (Chen et al. 2023)	WavLM	136	6	40K

#### **Multilingual Speech Representation Models:**

★ One single back-bone for (multilingual) speech applications

★ Zero-shot unseen languages

Model	Training Approach	# Languages	# Datasets	# Hours	
<b>XLSR-53</b> (Conneau et al. 2020)	wav2vec 2.0	53	3	56K	
<b>XLS-R</b> (Babu et al. 2021)	wav2vec 2.0	128	5	436K	
MMS (Pratap et al. 2023) [SOTA]	wav2vec 2.0	1,362	6	491K	
<b>WavLabLM</b> (Chen et al. 2023)	WavLM	136	6	40K	
mHuBERT-147	HuBERT	147	17	90K	

#### mHuBERT-147: A Compact Multilingual HuBERT Model

- A. **HIGH-QUALITY DATA:** Prioritizing open-license <u>smaller collections and dataset</u> <u>diversity over data quantity</u>
- B. Hubert Training: Most of the multilingual models follow the wav2vec 2.0 training approach
- C. COMPACT SIZE: Existing models are quite large (317M to 2B parameters), resulting in large downstream applications

#### DATA: 90,430 hours of diverse high-quality data

- → We gather 17 datasets across 147 languages
- → We filter popular large datasets **removing noise/music**
- → We downsample high-resource languages (>2,000 hours per source)

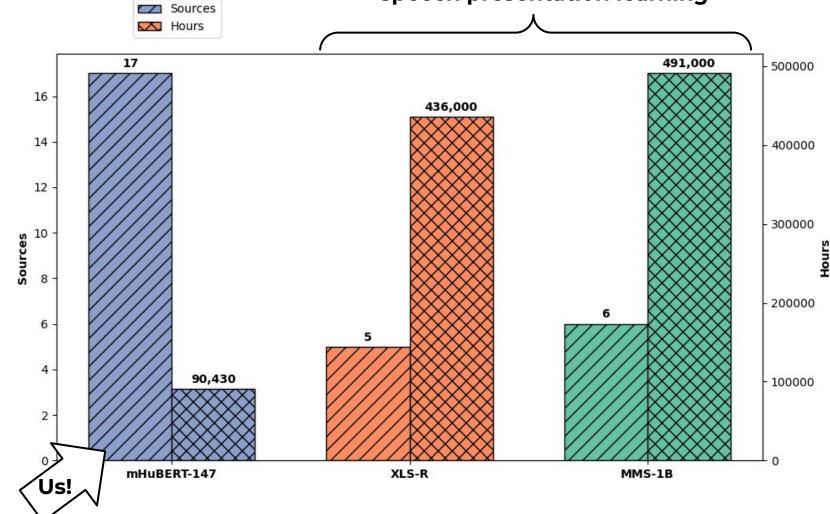
Dataset	<b>Full Names and References</b>	# Languages	# Hours (filtered)	License
Aishells	Aishell [4] and AISHELL-3 [41]	1	212	Apache License 2.0
B-TTS	BibleTTS [20]	6	358	CC BY-SA 4.0
Clovacall	ClovaCall [17]	1	38	MIT
CV	Common Voice version 11.0 [1]	98	14,943	CC BY-SA 3.0
	High quality TTS data for Javanese, Khmer,			
G-TTS	Nepali, Sundanese, and Bengali Languages [42]	9	27	CC BY-SA 4.0
	High quality TTS data for four South			
	African languages [46]			
<b>IISc-MILE</b>	IISc-MILE Tamil and Kannada ASR Corpus [26,27]	2	406	CC BY 2.0
JVS	Japanese versatile speech [44]	1	26	CC BY-SA 4.0
Kokoro	Kokoro Speech Dataset [19]	1	60	CC0
kosp2e	Korean Speech to English Translation Corpus [9]	1	191	CC0
MLS	Multilingual LibriSpeech [32]	8	50,687	CC BY 4.0
MS	MediaSpeech [21]	1	10	CC BY 4.0
Samrómur	Samrómur Unverified 22.07 [43]	1	2,088	CC BY 4.0
TH-data	THCHS-30 [48] and THUYG-20 [37,38]	2	46	Apache License 2.0
VL	VoxLingua107 [45]	107	5,844	CC BY 4.0
VP	VoxPopuli [47]	23	15,494	CC0

#### **DATA:** How much is 90K hours of speech?

SOTA models for multilingual speech presentation learning

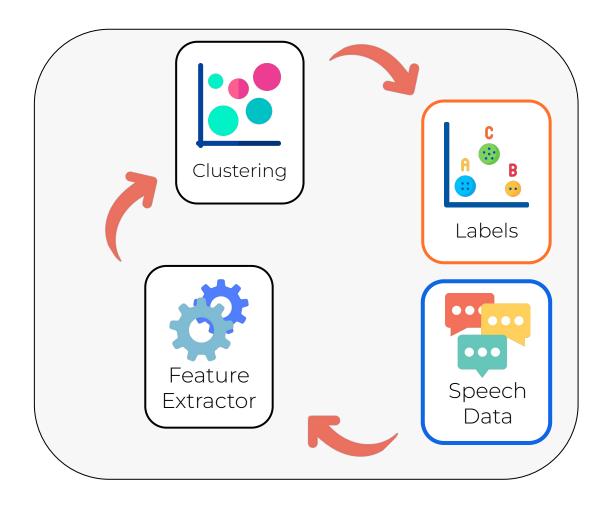
**LESS** training data compared to SOTA approaches

But **MORE** dataset diversity, **BETTER** language ratio



TRAINING APPROACH: Hidden Units BERT (Hsu et al. 2021)

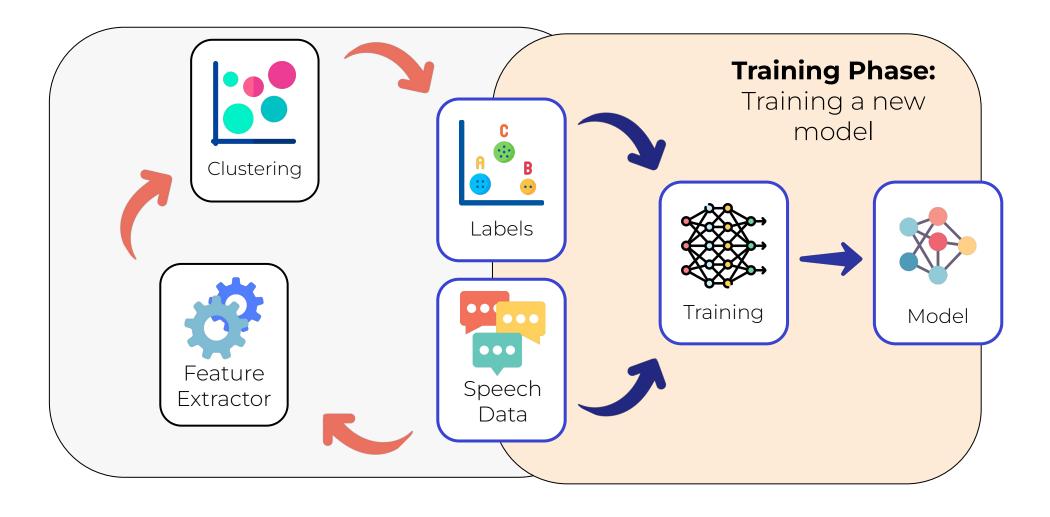
#### TRAINING APPROACH: Hidden Units BERT (Hsu et al. 2021)



### Labeling Phase:

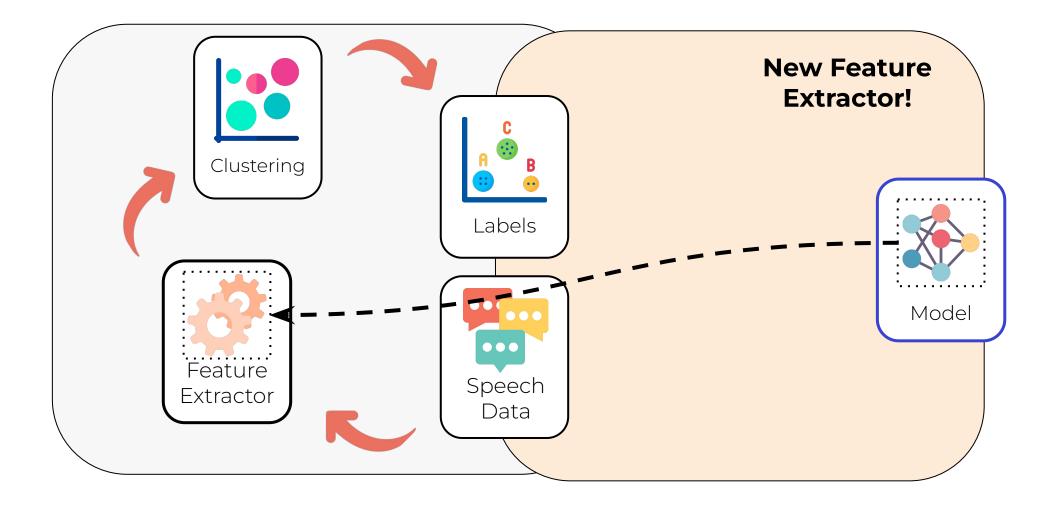
Creating fake discrete labels

#### TRAINING APPROACH: Hidden Units BERT (Hsu et al. 2021)



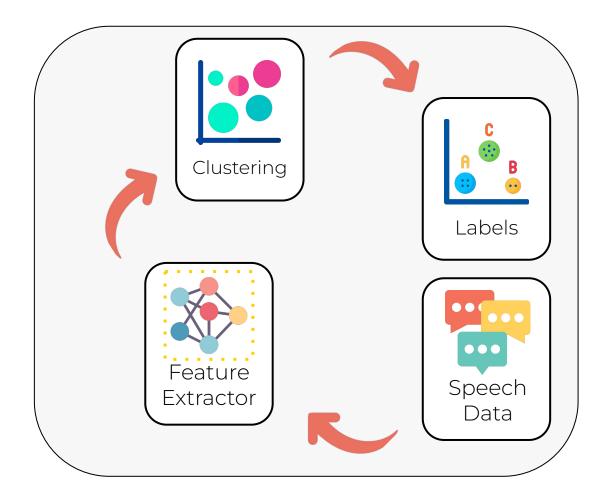
12

#### TRAINING APPROACH: Hidden Units BERT (Hsu et al. 2021)



13

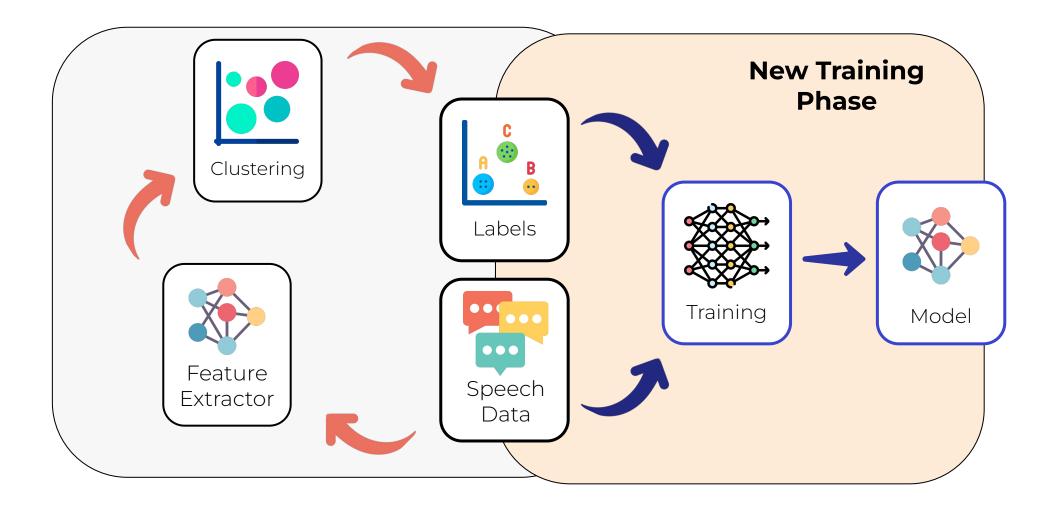
#### TRAINING APPROACH: Hidden Units BERT (Hsu et al. 2021)



### New Labeling Begins!

14

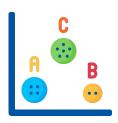
#### TRAINING APPROACH: Hidden Units BERT (Hsu et al. 2021)



#### Scaling the Approach for Multilingual Settings

#### A. Clustering/Labeling step are extremely expensive and slow

- We make training set smaller, but higher-quality
- We make clustering faster by using faiss + graph search for labeling (5.2x times faster)



#### **B.** Multilingual data distribution

 We propose a multilingual batching approach: two-step language, data source up-sampling



#### C. Training iterations from longer!

 We find that existing speech representation models tend to be quite undertrained (accounting for some grokking)



#### The mHuBERT-147 models

- ★ HIGH QUALITY DATA: +90K hours in 147 languages
- ★ **COMPACT SIZE:** 95M parameters, 1000 discrete units
- TRAINED FOR LONGER: 3 iterations, each for 2M updates (~20 days at 32xA100-80GB) at NAVER Cloud platform<sup>1</sup>

## Evaluation



mHuBERT-147 Evaluation: ML-SUPERB NAVER LABS

## **EVALUATION:** Multilingual Speech Representation Benchmark (ML-SUPERB)

- A. Two setups: 10min and 1h per language
- B. Two language settings: normal (123 languages) and few-shot (20 languages)
- C. Four tasks:
  - monolingual ASR
  - multilingual ASR
  - Language Identification (LID)
  - Multilingual ASR + LID

Final ranking SUPERB scores are computed considering SOTA and baseline scores

#### **ML-SUPERB:** Leaderboard Summary



SSL backbone	# Parameters	SUPERB Score 10min (1)	SUPERB Score 1h (1)	
MMS-1B	965M	983.5	948.1	
mHuBERT-147	95M	949.8	950.2	
MMS-300M	317M	824.9	844.3	
NWHC1 (MMS-300M variant)	317M	774.4	876.9	
NWHC2 (MMS-300M variant)	317M	759.9	873.3	
XLS-R-300M	317M	730.8	850.5	
WavLabLM-large-MS	317M	707.5	740.9	

Table: SUPERB scores (10min/1h).

#### **ML-SUPERB:** Leaderboard Summary

★ mHuBERT-147 is competitive while being much more compact!

SSL backbone	# Parameters	SUPERB Score 10min (1)	SUPERB Score 1h (↑)		
MMS-1B	965M	983.5 ①	948.1		
mHuBERT-147	95M	949.8	950.2		
MMS-300M	317M	824.9	844.3		
NWHC1 (MMS-300M variant)	317M	774.4	876.9		
NWHC2 (MMS-300M variant)	317M	759.9	873.3		
XLS-R-300M	317M	730.8	850.5		
WavLabLM-large-MS	317M	707.5	740.9		

Table: SUPERB scores (10min/1h).



mHuBERT-147 Evaluation: ML-SUPERB NAVER LABS

#### **ML-SUPERB:** Detailed Scores

SSI backbone	# Monolingual ASR CER (+)		Multilingual ASR (normal/few-shot) CER (↓)		LID ACC (1)		Multilingual ASR+LID (normal/few-shot) ACC - CER/CER (↓/↓)		
	Params	10min	1h	10min	1h	10min	1h	10min	1h
MMS-1B	965M	33.3	25.7	21.3/30.2	18.1/30.8	84.8	86.1	73.3 - <b>26.0/25.4</b>	74.8 - 25.5/ <b>24.8</b>
mHuBERT-147	95M	34.2	26.3	23.6/33.2	22.0/32.9	85.3	91.0	<b>81.4</b> - 26.2/34.9	<b>90.0</b> - <b>22.1</b> /33.5

**Table:** Detailed ML-SUPERB scores for the two best models.

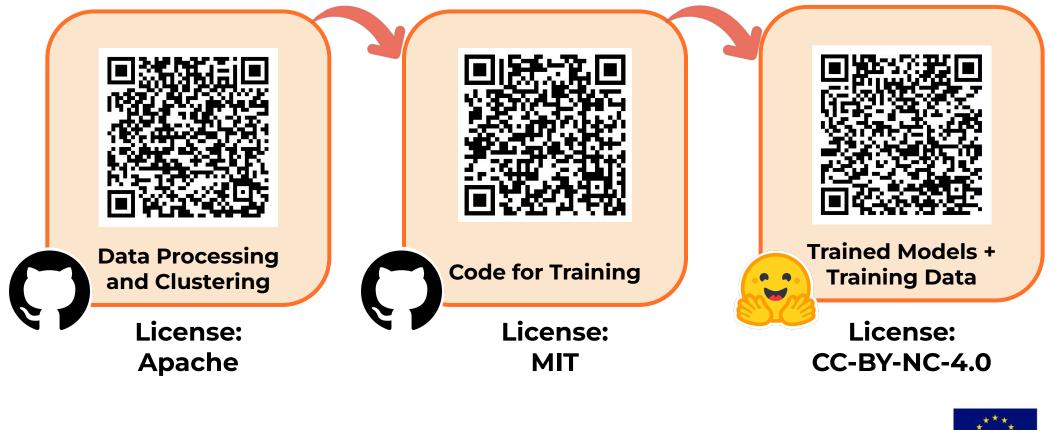


- Omitted from the table:
  - Competitive to MMS-300M
     We beat it in all tasks but 3: few-shot CER LID+ASR (10min/1h) and monolingual ASR (10min)
  - mHuBERT-147 > XLS-R and WavLabLM in all tasks

#### **Summarizing:**

- ★ We proposed the first multilingual HuBERT model, covering 147 languages
- ★ We approached data collection differently
- We proposed a two-level language-source up-sampling for training
- We produced a compact yet powerful foundation model that should be more exploitable for compact applications

#### mHuBERT-147 is an output of the EU project UTTER<sup>1</sup>





# mHuBERT-147: A Compact Multilingual HuBERT Model

Marcely Zanon Boito, Vivek Iyer, Nikolaos Lagos, Laurent Besacier, Ioan Calapodescu

07/2024

Contact: marcely.zanon-boito@naverlabs.com

**NAVER LABS** 

