

Speaker Representation Learning

Theories, Applications and Practice

Shuai Wang¹, Bing Han²

¹Shenzhen Research Institute of Big Data, ²Shanghai Jiao Tong University

December, 2023



深圳市大数据研究院
Shenzhen Research Institute of Big Data



上海交通大学
SHANGHAI JIAO TONG UNIVERSITY

Speaker Modeling: Background, Applications and Trends

Discriminative Speaker Representation Learning

Self-supervised based Speaker Representation Learning

Multi-modal Speaker Representation Learning

Efficiency and Robustness

Towards the Interpretability

Beyond Speaker Recognition

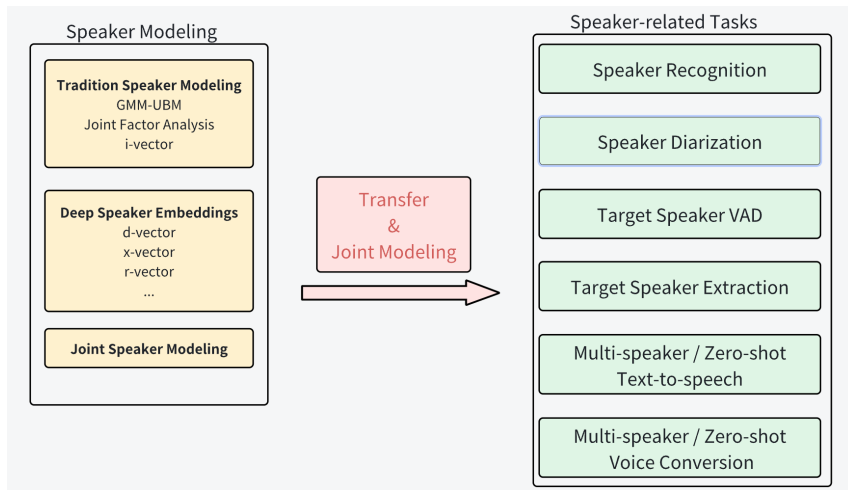
Practice: Speaker Representation Learning with Wespeaker

Speaker Modeling

Speaker modeling aims to capture information related to the identity of the speaker while neglecting other attributes.



Speaker modeling in different tasks



Applications of Speaker Modeling

Embedding Extraction for speaker verification

Speaker Verification: My voice is my password

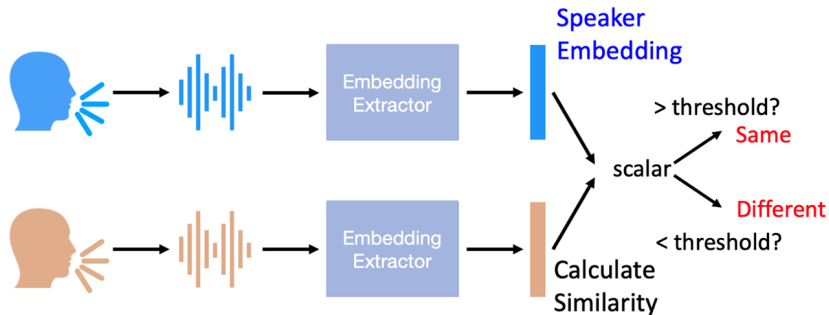


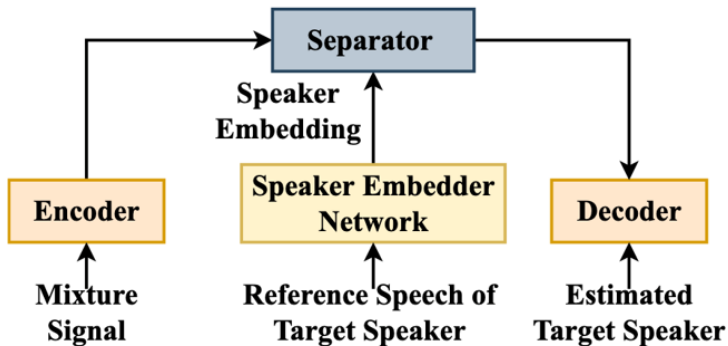
Figure adapted from Hung-yi Lee's DLHLP20 slides¹

¹<https://speech.ee.ntu.edu.tw/~tlkagk/courses/DLHLP20/>

Applications of Speaker Modeling

Reference/Cue modeling for Target speech extraction

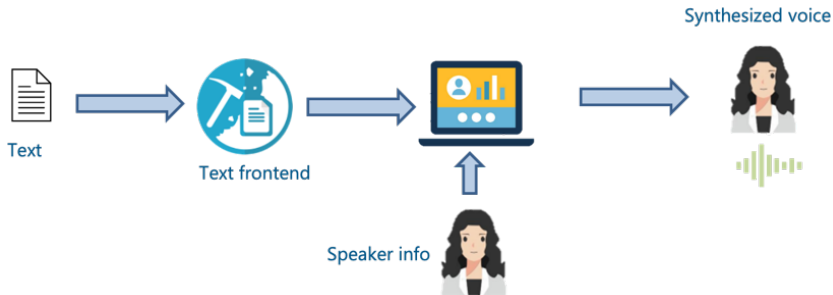
Target Speech Extraction: I only hear your voice



Applications of Speaker Modeling

Target voice identifier for Text-to-speech

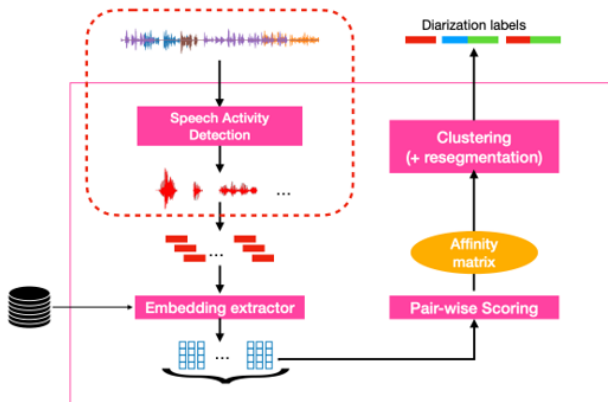
Speaker modeling for the target voice.



Applications of Speaker Modeling

Embedding clustering based Speaker diarization

Speaker diarization: Who Spoke When?



4

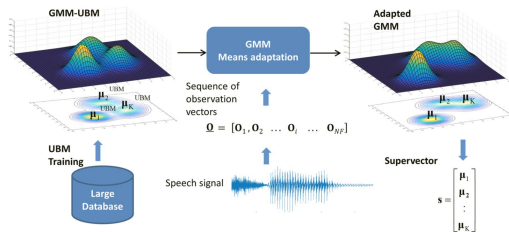


Figure: GMM-UMB^a model for speaker modeling

^aZheng, Zhang, and Xu, "Text-independent speaker identification using gmm-umb and frame level likelihood normalization".

Gaussian Mixture Model (GMM)^a

- ▶ $p(\mathbf{x}) = \sum_1^K c_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$ s.t. $\sum_1^K c_k = 1$
- ▶ Any distribution can be approximated by a weighted linear combination of several Gaussian distributions
- ▶ When using GMM to model the speaker's acoustic features, the number of Gaussians can be considered as the types of sounds produced.

Universal Background Model (UBM)

- ▶ Usually, a person's registered voice is limited (a few seconds), making it difficult to train a GMM with this data.
- ▶ UBM on a large-scale dataset can be trained first and then adapted to a specific speaker's data.

GMM-Supervector

- ▶ Concatenate the mean vector of each gaussian to represent the speaker

^aReynolds et al., "Gaussian mixture models."

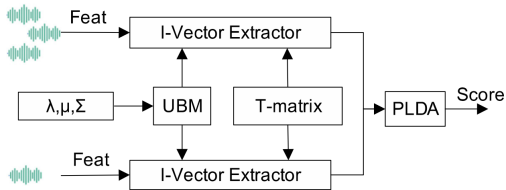


Figure: Block diagram of speaker recognition system based on i-vector

Drawbacks of GMM-Supervector

- ▶ Supervectors are extremely high-dimensional (often tens of thousands of dimensions), making them computationally challenging

Decompose Supervector to low-dimensional i-vector^a:

$$M(s) = m + Tw(s)$$

- ▶ $M(s)$: GMM-Supervector for speaker s
- ▶ m : speaker-independent supervector
- ▶ T : total variability matrix, capturing all sources of variability in the voice samples (both speaker-related and channel-related)
- ▶ $w(s)$: i-vector for speaker s

^aDehak et al., "Front-end factor analysis for speaker verification".

Speaker Modeling

d-vector

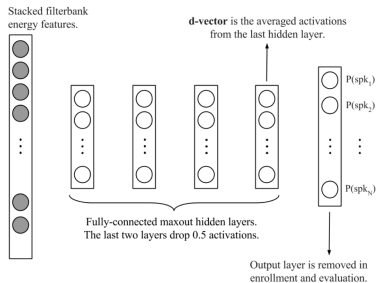


Figure: Architecture of d-vector

d-vector^a stands out as

- ▶ Very early attempts at applying deep neural networks to speaker information modeling
- ▶ Demonstrated a good complementarity with i-vector

^aVariani et al., "Deep neural networks for small footprint text-dependent speaker verification".

Table: Architecture of TDNN based speaker embedding extractor, T denotes the sequence length, N is the number of speakers

Layer	Layer context	Input \times output
frame1	$[t - 2, t + 2]$	200×512
frame2	$\{t - 2, t, t + 2\}$	1536×512
frame3	$\{t - 3, t, t + 3\}$	1536×512
frame4	$\{t\}$	512×512
frame5	$\{t\}$	512×1500
stats pooling	$[0, T]$	1500×3000
segment1	$\{0\}$	3000×512
segment2	$\{0\}$	512×512
projection	$\{0\}$	$512 \times N$

x-vector^a stands out as

- ▶ The first deep speaker embedding that beats traditional methods on well-recognized datasets (NIST SRE)
- ▶ The first work that introduces segment-level optimization
- ▶ Powerful variant ECAPA-TDNN^b

^aSnyder et al., "X-vectors: Robust dnn embeddings for speaker recognition".

^bDesplanques, Thienpondt, and Demuynck, "Ecapa-tdnn: Emphasized channel attention, propagation and aggregation in tdnn based speaker verification".

Table: Architecture of ResNet34 based speaker embedding extractor, T denotes the sequence length, N is the number of speakers

Layer name	Structure	Output
Input	–	$40 \times T \times 1$
Conv2D-1	3×3 , Stride 1	$40 \times T \times 32$
ResNetBlock-1	$\begin{bmatrix} 3 \times 3, 32 \\ 3 \times 3, 32 \end{bmatrix} \times 3$, Stride 1	$40 \times T \times 32$
ResNetBlock-2	$\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 4$, Stride 2	$20 \times \frac{T}{2} \times 64$
ResNetBlock-3	$\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 6$, Stride 2	$10 \times \frac{T}{4} \times 128$
ResNetBlock-4	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 3$, Stride 2	$5 \times \frac{T}{8} \times 256$
StatsPooling	–	10×256
Flatten	–	2560
Dense	–	256
Projection	–	N

r-vector^a stands out as

- ▶ The winner system of both tracks of VoxSRC 2019.
- ▶ Utilized as the embeddings in the winner systems of all 4 tracks in DIHARD 2019.

^aZeinali et al., “But system description to voxceleb speaker recognition challenge 2019”.

Speaker Modeling

Data: From well-labeled recordings to unlabeled in-the-wild massive data

Labeled data:

- ▶ Labeling the data is costly
- ▶ Automatically collected data using speaker information, like voxceleb, suffering from the privacy issue
 - ▶ The voxceleb dataset is no longer accessible from the official website

Unlabeled data:

- ▶ Easily to obtain
- ▶ Covers a wider range of real data
- ▶ No privacy issue

Training paradigm: Supervised to Unsupervised/Semi-supervised/Self-supervised

Speaker Representation Learning: Trends

Model: From shallow to deep

- ▶ GMM, i-vector can be regarded as single-layer MLP
- ▶ d-vector, j-vector, x-vector: less than 10 layers.
- ▶ ResNet based models (common setup: 34 layers, extended to 293 layers or even 500+ layers in challenges)

Speaker Representation Learning: Trends

Modality: From uni-modal to multi-modal

- ▶ Pure audio modality
- ▶ Audio-visual speaker embedding

Speaker Representation Learning: Trends

Paradigm: Training from scratch to leveraging pretrained models

- ▶ Training speaker discriminative models from scratch
- ▶ Leveraging large pretrained speech models such as WavLM
- ▶ Semi-supervised: Finetuning on the DINO models
- ▶ Semi-supervised: Iteratively clustering and supervised finetuning

Speaker Representation Learning: Trends

Task: From single task to cross-task

- ▶ Pretrained embeddings to be used in different tasks
- ▶ Explicit joint optimization with the specific task
- ▶ Implicit speaker modeling in related tasks

Speaker Modeling: Background, Applications and Trends
Discriminative Speaker Representation Learning
Self-supervised based Speaker Representation Learning
Multi-modal Speaker Representation Learning
Efficiency and Robustness
Towards the Interpretability
Beyond Speaker Recognition
Practice: Speaker Representation Learning with Wespeaker

Discriminative Loss Functions

Philosophy

Optimize the neural network towards the speaker classification direction.

Discriminative speaker Representation learning

Comparison with ASR training

Compared with the classic phoneme classification based ASR system:

ASR: **Close-set problem**, the classes are the pre-defined phonemes/senones in the inference stage

Speaker: **Open-set problem**, usually we assume the speakers in the inference stage are not present in the training set

Discriminative Loss Functions

Classification based loss functions provide discriminative supervision signals, while margin-based ones²³ are superior to the standard softmax-based one⁴.

$$L_{\text{softmax}} = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{\mathbf{W}_{y_i}^T \mathbf{x}_i + \mathbf{b}_{y_i}}}{\sum_{j=1}^c e^{\mathbf{W}_j^T \mathbf{x}_i + \mathbf{b}_j}} \quad (1)$$

$$L_{\text{AAM-Softmax}} = -\frac{1}{N} \sum_{i=1}^N \log \frac{e^{s(\cos(\theta_{y_i, i} + m))}}{Z} \quad (2)$$

²Huang, Wang, and Yu, "Angular Softmax for Short-Duration Text-independent Speaker Verification."

³Cai, Chen, and Li, "Exploring the encoding layer and loss function in end-to-end speaker and language recognition system".

⁴Xiang et al., "Margin matters: Towards more discriminative deep neural network embeddings for speaker recognition".

Discriminative speaker Representation learning

The essence of learning

The margin-based loss only enlarges the inter-speaker distance.



Figure: The decision boundary change after adding the margin-based loss.

An extra center loss^{abc} can be applied to minimize the within-class variance.

$$\begin{aligned} L &= L_{\text{AAM-Softmax}} + L_C \\ &= L_{\text{AAM-Softmax}} + \frac{1}{2} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{c}_{y_i}\|^2 \end{aligned}$$

^aCai, Chen, and Li, "Exploring the encoding layer and loss function in end-to-end speaker and language recognition system".

^bLi et al., "Deep Discriminative Embeddings for Duration Robust Speaker Verification."

^cWang et al., "Discriminative neural embedding learning for short-duration text-independent speaker verification".

Speaker Modeling: Background, Applications and Trends
Discriminative Speaker Representation Learning
Self-supervised based Speaker Representation Learning
Multi-modal Speaker Representation Learning
Efficiency and Robustness
Towards the Interpretability
Beyond Speaker Recognition
Practice: Speaker Representation Learning with Wespeaker

Self-supervised based Speaker Representation Learning

- ▶ Leveraging large pre-text pretraining models.
 - ▶ Self-supervised Pretrained Speech Models
 - ▶ ASR Model Initialization
 - ▶ Efficient Finetuning
- ▶ Self-supervised Learning Approach
 - ▶ SimCLR/MoCo/DINO
 - ▶ Stage-wise Iterative Training

Self-supervised based Speaker Representation Learning

Finetuning Approach

► Self-Supervised Pretrained Speech Models

- Wav2Vec^a
- HuBERT^b
- WavLM^c
- UniSpeech^d

^aBaevski et al., “wav2vec 2.0: A framework for self-supervised learning of speech representations”.

^bHsu et al., “Hubert: Self-supervised speech representation learning by masked prediction of hidden units”.

^cChen et al., “Wavlm: Large-scale self-supervised pre-training for full stack speech processing”.

^dChen et al., “Unispeech-sat: Universal speech representation learning with speaker aware pre-training”.

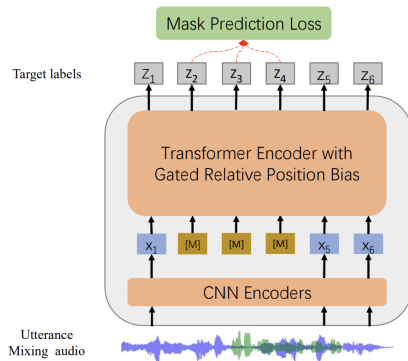


Figure: Model Architecture of WavLM

Self-supervised based Speaker Representation Learning

Finetuning Approach

Finetuning SSL Speech Models on Speaker Verification Task^a

- ▶ Replace Fbank with representation from pre-trained models.
- ▶ Learnable weighted sum

^aChen et al., "Large-scale self-supervised speech representation learning for automatic speaker verification".

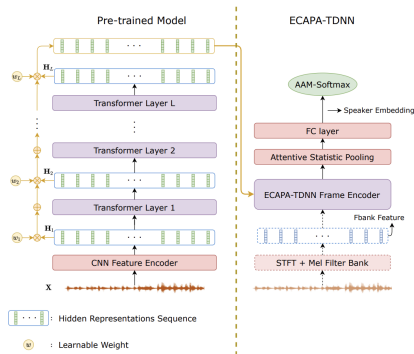


Figure: Leverage Representations from Pre-trained Model

Self-supervised based Speaker Representation Learning

Finetuning Approach

Efficient finetuning Self-supervised Model with adapters on Speaker Verification^a

- ▶ Frozen the large pretrained model
- ▶ Use adapters for efficient finetuning on speaker tasks.

^aPeng et al., "Parameter-efficient transfer learning of pre-trained Transformer models for speaker verification using adapters".

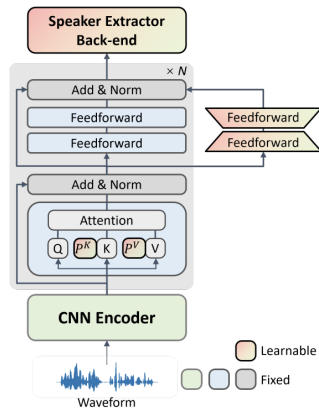


Figure: Schematic diagram of efficient finetuning

Self-supervised based Speaker Representation Learning

Finetuning Approach

Finetuning ASR Models on Speaker Verification Task^{ab}

- ▶ Pre-train model with ASR dataset.
- ▶ Initialize for speaker task training.

^aLiao et al., "Towards a unified conformer structure: from asr to asv task".

^bCai et al., "Pretraining Conformer with ASR for Speaker Verification".

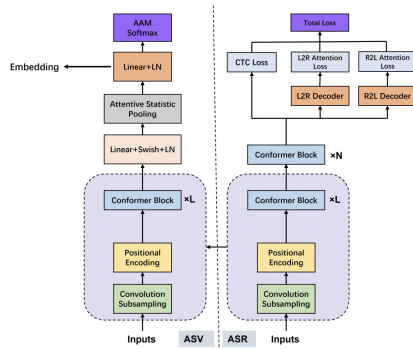


Figure: Schematic diagram of ASR transferring

Self-supervised based Speaker Representation Learning

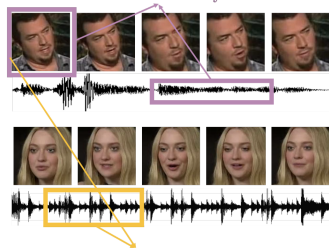
Self-supervised Learning Approach

Assumption of self-supervised learning on Speaker Verification Task.^a

- ▶ Segments from same utterances belong to same speaker.
- ▶ Segments from different utterances belong to different speakers.

^aHuh et al., "Augmentation adversarial training for self-supervised speaker recognition".

Within the same track = same identity but different content



Different tracks = different identity and different content

Figure: Schematic diagram of assumption.

Self-supervised based Speaker Representation Learning

Metric based Loss Functions

Metric learning based loss functions provide contrastive supervision signals, such as Triplet, Prototypical, GE2E⁵ and Angular Prototypical⁶.

$$L_{\text{Triplet}} = \frac{1}{N} \sum_{j=1}^N \max(0, \|\mathbf{x}_{j,0} - \mathbf{x}_{j,1}\|_2^2 - \|\mathbf{x}_{j,0} - \mathbf{x}_{k \neq j,1}\|_2^2 + m) \quad (3)$$

$$L_{\text{Prototypical}} = -\frac{1}{N} \sum_{j=1}^N \log \frac{e^{\mathbf{S}_{j,j}}}{\sum_{k=1}^N e^{\mathbf{S}_{j,k}}}, \text{ where } \mathbf{S}_{j,k} = \|\mathbf{x}_{j,M} - \mathbf{c}_k\|_2^2 \quad (4)$$

⁵Wan et al., "Generalized end-to-end loss for speaker verification".

⁶Chung et al., "In defence of metric learning for speaker recognition".

Self-supervised based Speaker Representation Learning

SimCLR

Based on SimCLR^a framework, adapt to speaker task^b

- ▶ Crop two segments from utterance and construct the positive and negative pairs.
- ▶ Use metric loss to attract the positive pairs and repel the negative pairs.

^aChen et al., "A simple framework for contrastive learning of visual representations".

^bZhang, Zou, and Wang, "Contrastive self-supervised learning for text-independent speaker verification".

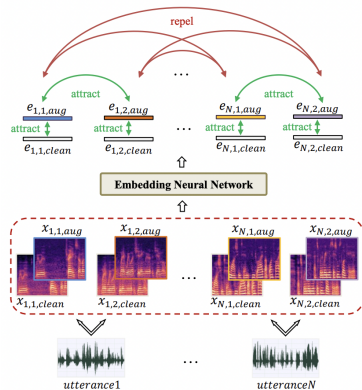


Figure: Schematic diagram of simclr on speaker task.

Self-supervised based Speaker Representation Learning

DINO

Based on DINO^a framework, adapt to speaker task^{bc}

- ▶ Crop several segments from one utterance and only construct the positive pairs.
- ▶ Use cross entropy loss to attract the positive pairs.

^aCaron et al., "Emerging properties in self-supervised vision transformers".

^bHan, Chen, and Qian, "Self-supervised speaker verification using dynamic loss-gate and label correction".

^cChen et al., "A comprehensive study on self-supervised distillation for speaker representation learning".

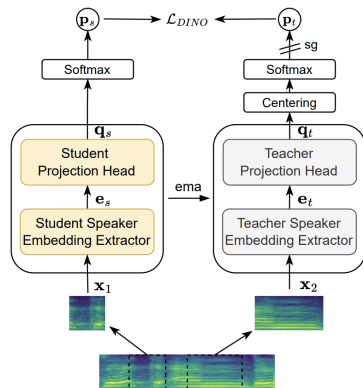


Figure: Schematic diagram of DINO on speaker task.

Self-supervised based Speaker Representation Learning

Stage-wise Iterative Training

Two stages based iterative framework^{abc}.

- ▶ I: Contrastive training
- ▶ II: Iterative clustering and representation learning.

^aCai, Wang, and Li, "An iterative framework for self-supervised deep speaker representation learning".

^bHan, Chen, and Qian, "Self-Supervised Learning with Cluster-Aware-DINO for High-Performance Robust Speaker Verification".

^cTao et al., "Self-supervised speaker recognition with loss-gated learning".

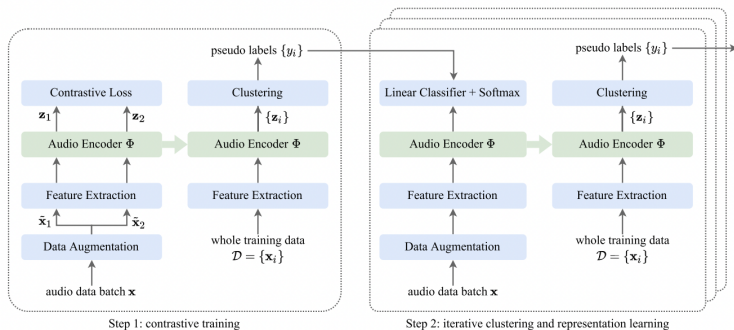


Figure: Schematic diagram of iterative framework for SSL speaker verification.

Speaker Modeling: Background, Applications and Trends
Discriminative Speaker Representation Learning
Self-supervised based Speaker Representation Learning
Multi-modal Speaker Representation Learning
Efficiency and Robustness
Towards the Interpretability
Beyond Speaker Recognition
Practice: Speaker Representation Learning with Wespeaker

The complementation between audio and visual modality

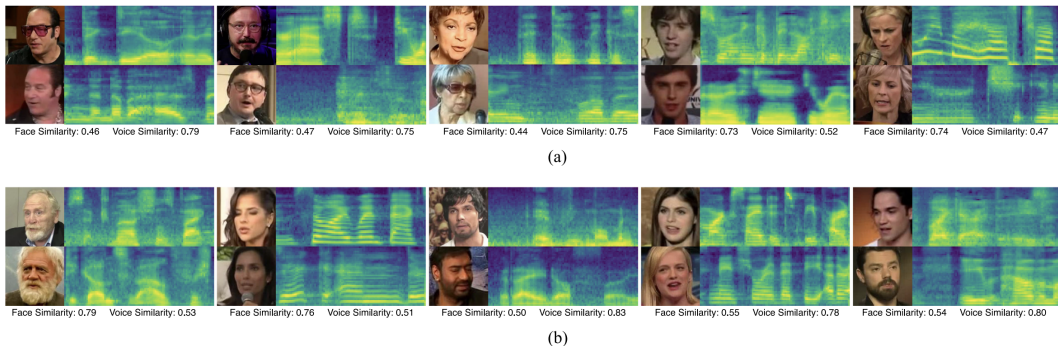
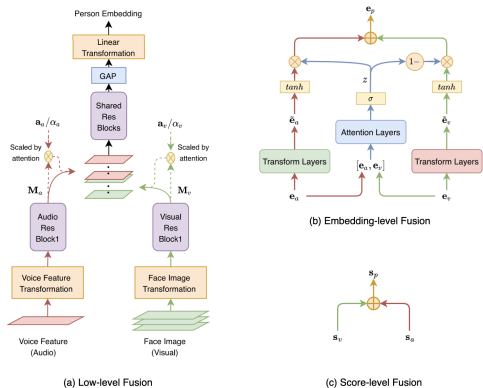


Figure: The speaker similarity based on the audio or visual information⁷

- ▶ The upper part shows the speaker's similarity to the same person
- ▶ The bottom part shows the speaker's similarity between different persons

⁷Qian, Chen, and Wang, "Audio-visual deep neural network for robust person verification".

Audio-Visual Information Fusion



- ▶ Embedding-level fusion performs better than low-level fusion
- ▶ The attention mechanism in embedding-level fusion makes it more noise-robust than score-level fusion

Figure: Audio-visual information fusion at different levels^a

^aQian, Chen, and Wang, "Audio-visual deep neural network for robust person verification".

Multi-Modal Knowledge Distillation

From audio-visual system to single-modality system

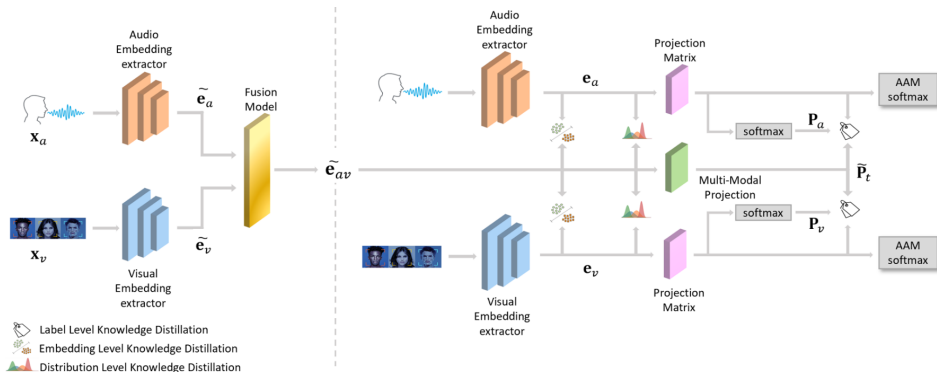


Figure: Knowledge distillation from audio-visual system to single-modality system⁸

⁸Zhang, Chen, and Qian, "Knowledge Distillation from Multi-Modality to Single-Modality for Person Verification".

Multi-Modal Knowledge Distillation

From visual system to audio system

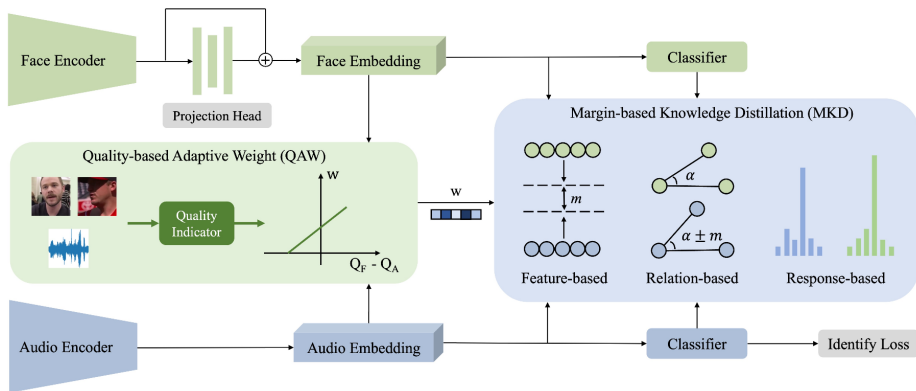


Figure: Knowledge distillation from visual system to audio system⁹

⁹Jin et al., "Cross-modal distillation for speaker recognition".

Speaker Modeling: Background, Applications and Trends
Discriminative Speaker Representation Learning
Self-supervised based Speaker Representation Learning
Multi-modal Speaker Representation Learning
Efficiency and Robustness
Towards the Interpretability
Beyond Speaker Recognition
Practice: Speaker Representation Learning with Wespeaker

In speaker representation learning, we mainly optimize the efficiency of the model from two perspectives: computational efficiency and memory efficiency

- ▶ Computation Efficiency
 - ▶ Knowledge Distillation
 - ▶ Network Quantization
 - ▶ Efficient Architecture Design
- ▶ Memory Efficiency
 - ▶ Reversible Neural Networks

Knowledge Distillation on Speaker Verification Task

- ▶ Knowledge distillation from teacher model to student model^a
- ▶ Self-knowledge distillation via feature enhancement^b
- ▶ Knowledge distillation from multi-modality to single-modality^c

^aWang et al., "Knowledge Distillation for Small Foot-print Deep Speaker Embedding".

^bLiu et al., "Self-Knowledge Distillation via Feature Enhancement for Speaker Verification".

^cZhang, Chen, and Qian, "Knowledge Distillation from Multi-Modality to Single-Modality for Person Verification".

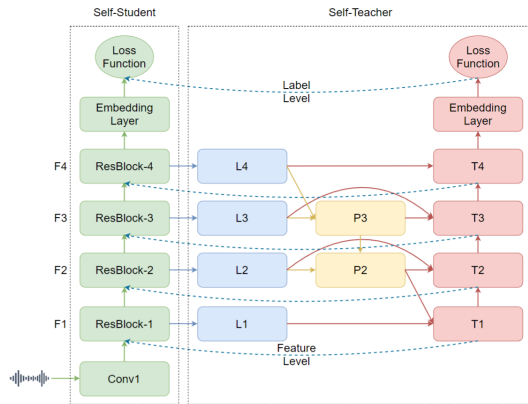


Figure: Schematic diagram of self-knowledge distillation via feature enhancement

Model Computation Efficiency

Quantization

Quantization achieves model compression by reducing the parameter precision

- ▶ Binary Neural Network^a
- ▶ Linear and PoT(Power of Two) quantization^b
- ▶ K-Means based quantization^c
- ▶ Static and adaptive quantizer for binary quantization^d

^aZhu, Qin, and Li, "Binary Neural Network for Speaker Verification".

^bLiu et al., "Self-Knowledge Distillation via Feature Enhancement for Speaker Verification".

^cWang et al., "Adaptive Neural Network Quantization For Lightweight Speaker Verification".

^dLiu, Wang, and Qian, "Extremely Low Bit Quantization for Mobile Speaker Verification Systems Under 1MB Memory".

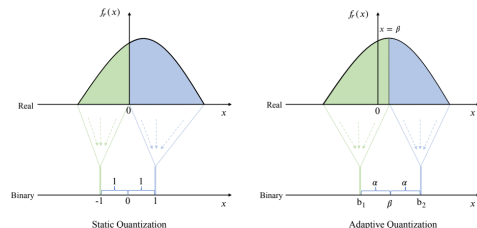


Figure: The overview of static and adaptive binary quantization

Efficient Architecture Design on Speaker Verification Task

- ▶ Depth-First Neural Architecture with Attentive Feature Fusion^a
- ▶ CS-CTCSCONV1D^b(Channel Split Time-Channel-Time Separable 1-dimensional Convolution)
- ▶ Asymmetric Enroll-Verify Structure(ECAPA-TDNNLite^c)

^aLiu, Chen, and Qian, "Depth-First Neural Architecture With Attentive Feature Fusion for Efficient Speaker Verification".

^bCai et al., "CS-CTCSCONV1D: Small footprint speaker verification with channel split time-channel-time separable 1-dimensional convolution".

^cLi et al., "Towards Lightweight Applications: Asymmetric Enroll-Verify Structure for Speaker Verification".

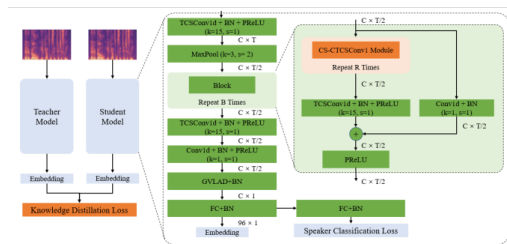


Figure: Schematic of CS-CTCSCONV1D

Model Computation Efficiency

Asymmetric Enroll-Verify Structure

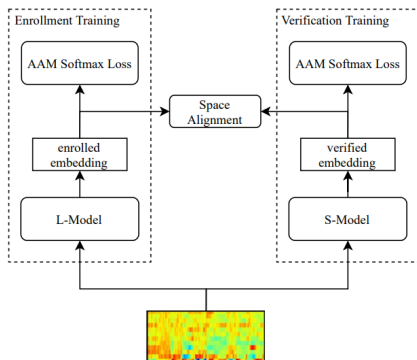


Figure: The training process of the asymmetric structure. Frame-wise input features are fed into the large-scale model and the small-scale model, respectively

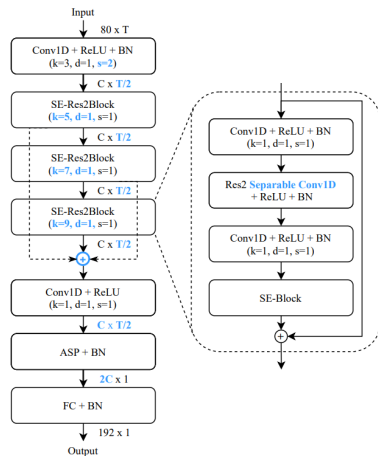


Figure: Schematic of ECAPA-TDNNLite

Model Computation Efficiency

Performance of Computational Efficient Models

Table: The experiment results of compressed/quantized ResNet34 and other full-precision compact architectures.

Model	Size (MB)	Bit-width (bit)	Vox1-O EER(%)
KMQAT-ResNet34 ¹⁰	3.45	4	0.957
PoT-ResNet34 ¹¹	3.45	4	1.09
TWN-ResNet34 ¹² (our impl.)	1.80	2	1.473
b-vector(adaptive) ¹³	0.97	1	1.72
ResNet34(binary) ¹⁴	0.66	1	5.355
CS-CTCConv1d	0.96	32	2.62
ECAPA-TDNNLite	1.2	32	3.07

¹⁰Wang et al., "Adaptive Neural Network Quantization For Lightweight Speaker Verification".

¹¹Li et al., "Model Compression for DNN-based Speaker Verification Using Weight Quantization".

¹²Li, Zhang, and Liu, "Ternary weight networks".

¹³Liu, Wang, and Qian, "Extremely Low Bit Quantization for Mobile Speaker Verification Systems Under 1MB Memory".

¹⁴Zhu, Qin, and Li, "Binary Neural Network for Speaker Verification".

Model Memory Efficiency

Training Memory Efficiency

Reversible Neural Networks^a (RevNets) alleviate the need to store activations in memory during back-propagation. Consequently, RevNets require nearly constant memory costs as the network depth increases.

- ▶ Partially reversible networks
- ▶ Fully reversible networks

^aLiu and Qian, “Reversible Neural Networks for Memory-Efficient Speaker Verification”.

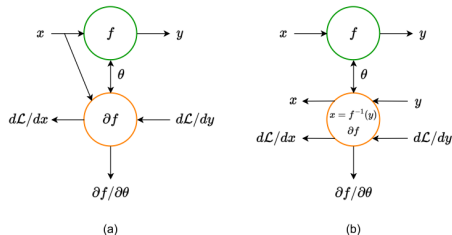


Figure: Comparison between non-reversible operator (a) and reversible operator (b)

Model Memory Efficiency

GPU Memory Usage vs. Parameter Number

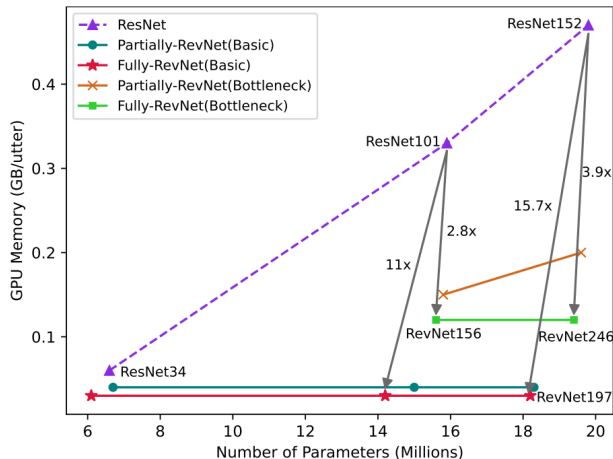


Figure: GPU Memory Usage vs. Parameter Number

Other Work on Model Efficiency

- ▶ Thin-ResNet¹⁵
- ▶ Fast-ResNet¹⁶
- ▶ ADMM¹⁷
- ▶ Small Footprint Text-Independent Speaker Verification¹⁸

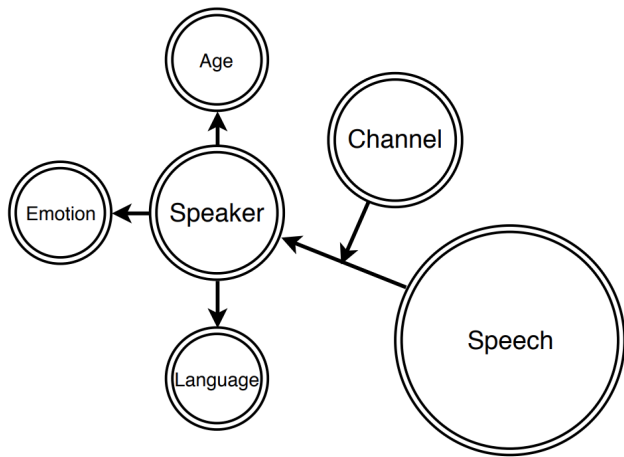
¹⁵Cai, Chen, and Li, "Exploring the encoding layer and loss function in end-to-end speaker and language recognition system".

¹⁶Chung et al., "In Defence of Metric Learning for Speaker Recognition".

¹⁷Xu et al., "Mixed Precision Low-Bit Quantization of Neural Network Language Models for Speech Recognition".

¹⁸Balian et al., "Small footprint text-independent speaker verification for embedded systems".

Robustness in Speaker Representation Learning



The recording environment also introduces variability in modeling speaker identity, influenced by factors like the recording device and microphone distance. To enhance model robustness across different devices, various domain adaptation methods are applied in speaker recognition, including

- ▶ Discrepancy-based alignment
- ▶ Adversarial learning
- ▶ Domain-specific adapter

Model Robustness to Device

Discrepancy-based alignment

Discrepancy-based alignment aims to minimize domain discrepancy in a latent feature space and facilitate learning domain-invariant representations. To achieve this goal, choosing a proper divergence measure is at the core of these methods. Widely used measures include MMD¹⁹, correlation alignment (CORAL)²⁰, etc.

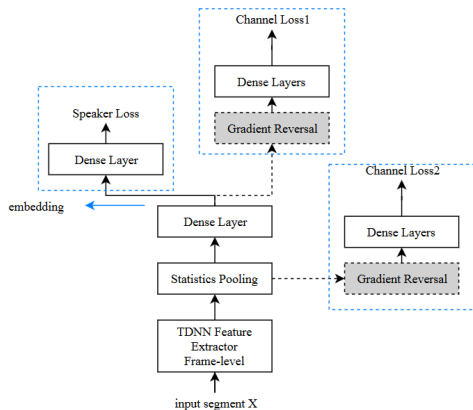
$$\mathcal{L}_{\text{mmd}} \triangleq \sup_{\phi \in \Phi} (\mathbb{E}_S [\phi(S)] - \mathbb{E}_T [\phi(T)]) \quad (5)$$

¹⁹Li, Han, and Song, "CDMA: Cross-Domain Distance Metric Adaptation for Speaker Verification".

²⁰Li, Zhang, and Chen, "The coral++ algorithm for unsupervised domain adaptation of speaker recognition" 

Model Robustness to Device

Adversarial learning



Adversarial learning employs a domain classifier to eliminate discriminative domain information from features. Min-max optimization in domain-adversarial training minimizes the domain gap and enforces domain-invariant feature extraction^a.

^aChen et al., "Channel invariant speaker embedding learning with joint multi-task and adversarial training".

Figure: Structure of channel-level adversarial learning^a

Model Robustness to Device

Domain-specific adapter

Instead of directly aligning domains with discrepancy measures, incorporating additional modules like domain-specific adapters helps capture and mitigate domain variances, resulting in domain-invariant embeddings.

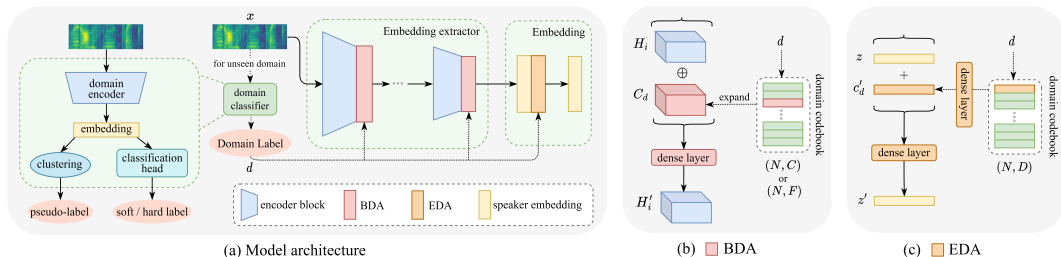


Figure: Framework with domain-specific adapters²¹

²¹Huang et al., "Enhancing Cross-Domain Speaker Verification through Multi-Level Domain Adapters".

Model Robustness to language mismatch

Language mismatch between datasets



Observation: In real-world scenarios, speaker verification systems may degrade when training on one language and test it on another.

Model Robustness to language mismatch

Language mismatch between enroll / test

Over 40% of the world's population is bilingual, this mismatch happens when the languages used are different for enrollment and test.

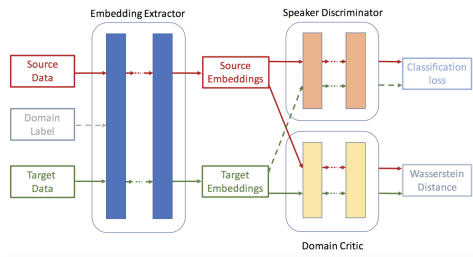


Figure: Structure of language-mismatch adversarial learning

Adversarial learning employs a language classifier to eliminate discriminative language information from features. Min-max optimization in domain-adversarial training minimizes the language gap and enforces language-invariant feature extraction^{ab}.

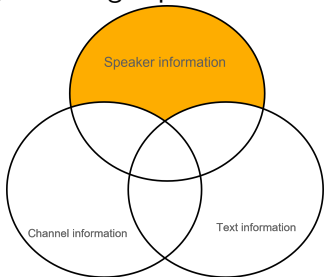
^aRohdin et al., "Speaker verification using end-to-end adversarial language adaptation".

^bXia, Huang, and Hansen, "Cross-lingual text-independent speaker verification using unsupervised adversarial discriminative domain adaptation".

Robustness

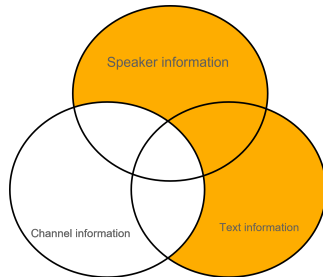
Robustness to text mismatch

Besides the speaker information, the text or content is the most crucial information conveyed through speech.



For text-independent speaker tasks, we only need **speaker information**

Enroll: **Hey Siri**; Test: **whatever to say**



For text-dependent speaker tasks, we also need **content information**

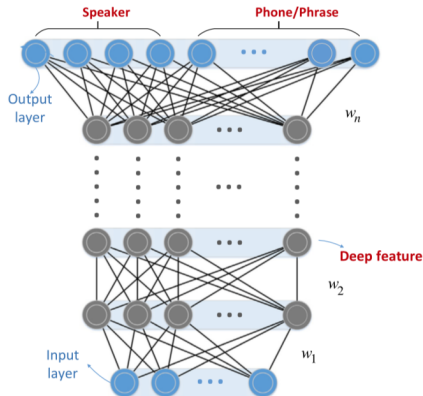
Enroll: **Hey Siri**; Test: **Hey Siri**

The representation of **Content Information**

- ▶ Phoneme index
- ▶ Phoneme posteriors predicted by ASR
- ▶ Hidden layer outputs from the ASR Model
- ▶ Phrase number (Fix-phrase datasets)
- ▶ Normalized phoneme distribution

Text Robustness

multi-task learning in the d-vector framework²²



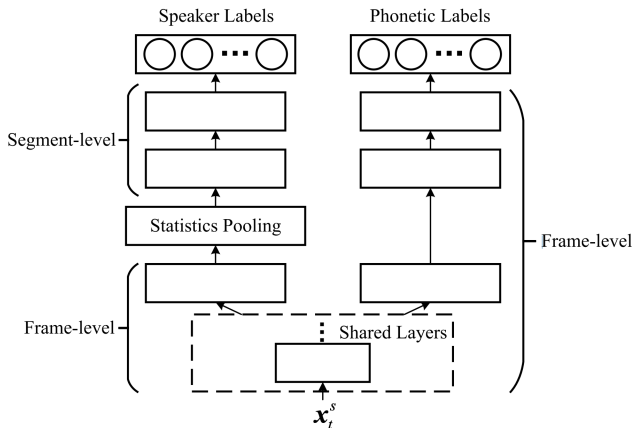
- ▶ Text-dependent task
- ▶ Multi-task at the **frame-level**
- ▶ Performance improved

Explicitly modeling phonetic information helps the text-dependent speaker verification task, which is intuitive.

²²Liu, Yuan, et al. "Deep feature for text-dependent speaker verification." *Speech Communication* 73 (2015): 1-13.

Text Robustness

multi-task learning in the x-vector framework²³

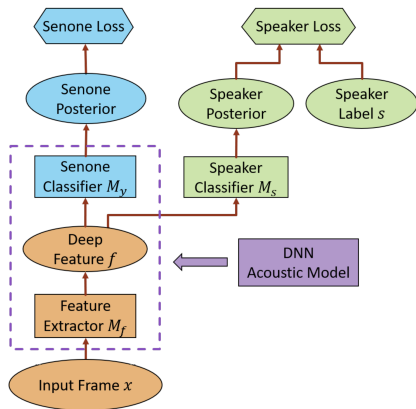


- ▶ Text-independent task
- ▶ Multi-task at the frame-level
- ▶ Performance improved!

²³Liu, Yi, et al. "Speaker Embedding Extraction with Phonetic Information." Proc. Interspeech 2018 (2018): 2247-2251.

Text Robustness

Speaker invariant training for ASR ²⁴

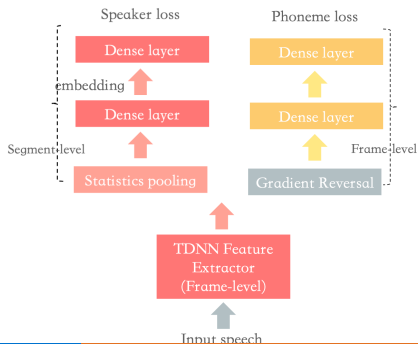
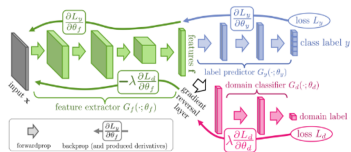


- ▶ Acoustic modelling
- ▶ Adversarial training suppressing the speaker effect
- ▶ Performance improved

²⁴Meng, Zhong, et al. "Speaker-invariant training via adversarial learning." 2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). IEEE, 2018.

Text Robustness

Frame-level multi-task/adversarial training



$$\mathcal{L}_s = \text{CE}(M_s(M_f(\mathbf{X})), \mathbf{y}^s)$$

$$\mathcal{L}_p = \frac{1}{N} \sum_{i=1}^N \text{CE}(M_p(M_f(\mathbf{x}_i)), \mathbf{y}_i^p)$$

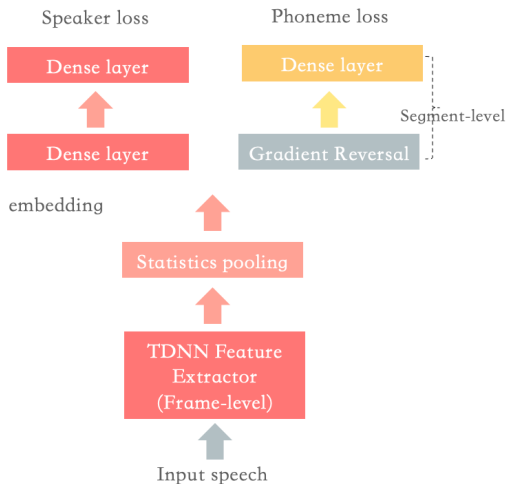
$$\mathcal{L}_{total} = \mathcal{L}_s + \mathcal{L}_p$$

Systems	voxceleb1_O	voxceleb1_E	voxceleb1_H
x-vector baseline	2.361	2.470	4.260
FRM-MT	2.165	2.198	3.911
FRM-ADV	3.143	3.214	5.419

Frame-level multi-task / adversarial training

Text Robustness

Segment-level multi-task/adversarial training



$$\mathcal{L}_s = \text{CE}(M_s(M_f(\mathbf{X})), \mathbf{y}^s)$$

$$\mathcal{L}_p = \text{CE}(M_p(M_f(\mathbf{x}_i)), \mathbf{y}^p)$$

$$\mathcal{L}_{total} = \mathcal{L}_s + \mathcal{L}_p$$

For a given segment \mathbf{x} with N frames, segment-level phoneme label \mathbf{y}^p is

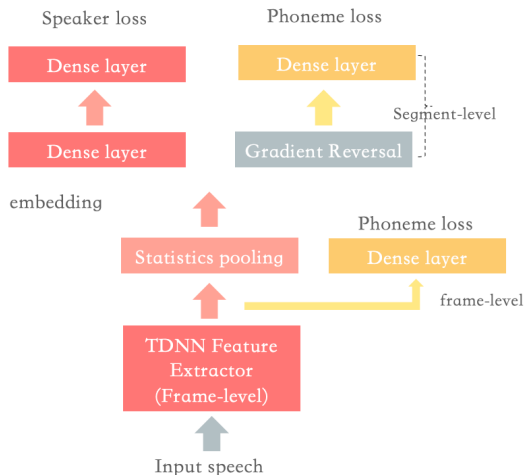
$$\mathbf{y}^p = \{y_1, y_2, \dots, y_C\}$$

$$y_c = \frac{N_c}{N}$$

where C is the size of the chosen phoneme set. N_c denotes the number of occurrences of the c -th phoneme in \mathbf{x}

Text Robustness

Frame-level multitask + segment-level adversarial learning



Systems	voxceleb1_O	voxceleb1_E	voxceleb1_H
x-vector baseline	2.361	2.470	4.260
SEG-MT	2.175	2.330	4.059
SEG-ADV	2.154	2.198	3.923

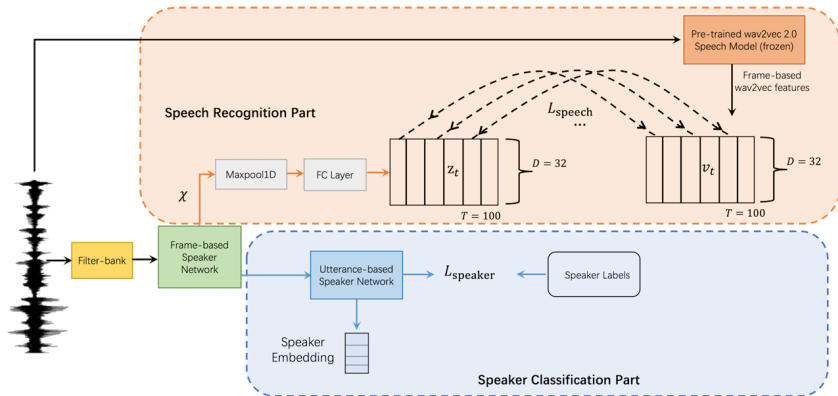
Segment-level Multitask/Adversarial training

Systems	voxceleb1_O	voxceleb1_E	voxceleb1_H
x-vector baseline	2.361	2.470	4.260
FRM-MT	2.165	2.198	3.911
SEG-ADV	2.154	2.198	3.923
COMBINE	2.013	2.030	3.819

Frame-level multi-task + segment-level adv training

Text Robustness

Multi-task training with high-level content representation²⁵



²⁵Jin, Tu, and Mak, "Phonetic-aware speaker embedding for far-field speaker verification".

Extract phonetic bottleneck (PBN) from a pretrained ASR model and combine it with the filterbanks²⁶

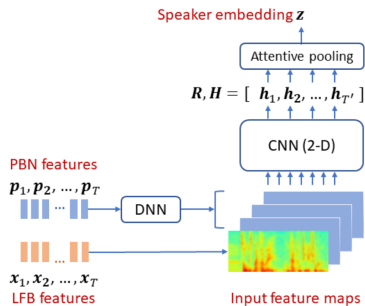


Fig. 1: Implicit phonetic attention by combining LFB and PBN features at the input layer (LFB: log filter bank; PBN: phonetic bottleneck).

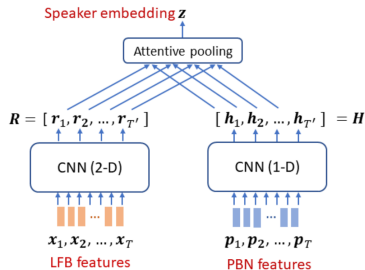


Fig. 2: Explicit phonetic attention by routing LFB and PBN features through separate networks (LFB: log filter bank; PBN: phonetic bottleneck).

²⁶Zhou T, Zhao Y, Li J, et al. CNN with phonetic attention for text-independent speaker verification, *ASRU*

Text Robustness

Phoneme-aware speaker embedding learning

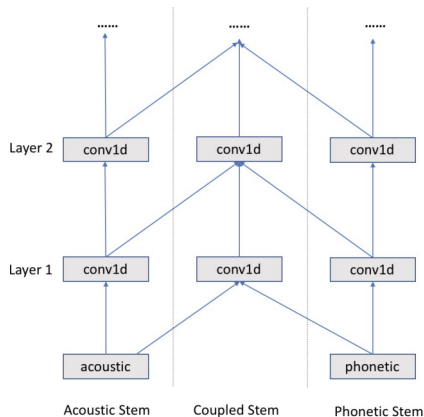


Table 1: Network configurations of PacNet

Layer 7	Linear	In=1024 Out=1000		
Layer 6	Pooling	In=1024 Out=1024		
Layer 5	Conv1d	In=2048 Out=1024		
Layer 4	Conv1d kernel=5	Out=512 In=512	Out=1024 In=2048	Out=512 In=512
Layer 3	Conv1d kernel=5	Out=512 In=512	Out=1024 In=2048	Out=512 In=512
Layer 2	Conv1d kernel=5	Out=512 In=512	Out=1024 In=2048	Out=512 In=512
Layer 1	Conv1d kernel=5	Out=512 In=40	Out=1024 In=140	Out=512 In=100
Stem		Acoustic	Coupled	Phonetic

► Triplet loss instead of softmax loss

27

²⁷Zheng, Lei, and Suo, "Phonetically-Aware Coupled Network For Short Duration Text-Independent Speaker Verification."



The general idea: **Decomposition and Reconstruction**. The application is far more than speaker modeling

Applications

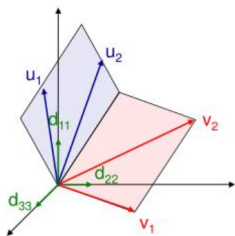
- ▶ Speaker representation learning
- ▶ Voice conversion
- ▶ Speech synthesis/ Voice Cloning

Speech representation disentanglement

Back to old times

Joint Factor Analysis for speaker representation learning²⁸

$$\mathbf{M} = \mathbf{M}^{\text{UBM}} + \mathbf{V}\mathbf{y} + \mathbf{D}\mathbf{z} + \mathbf{U}\mathbf{x}$$

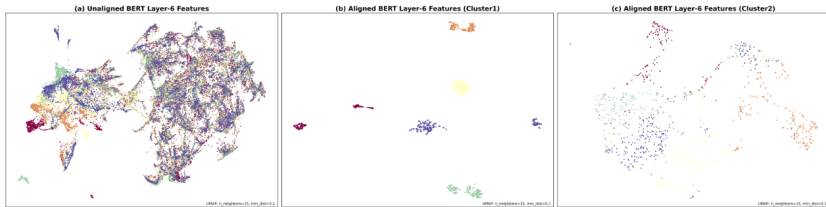
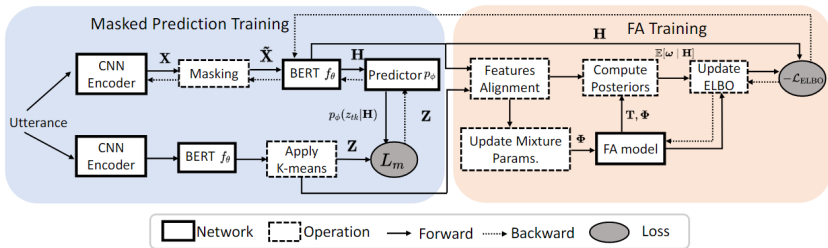


- ▶ Gaussian priors assumed for factors \mathbf{y} , \mathbf{z} , \mathbf{x}
- ▶ \mathbf{M}^{UBM} , \mathbf{V} , \mathbf{D} , \mathbf{U} are estimated using EM algorithm
- ▶ \mathbf{V} captures main speaker variability (Eigen voices)
- ▶ \mathbf{D} captures channel variability
- ▶ \mathbf{U} captures residual variability

²⁸Kenny, Patrick, et al. "Joint factor analysis versus eigenchannels in speaker recognition." TASLP 2007

Speech representation disentanglement

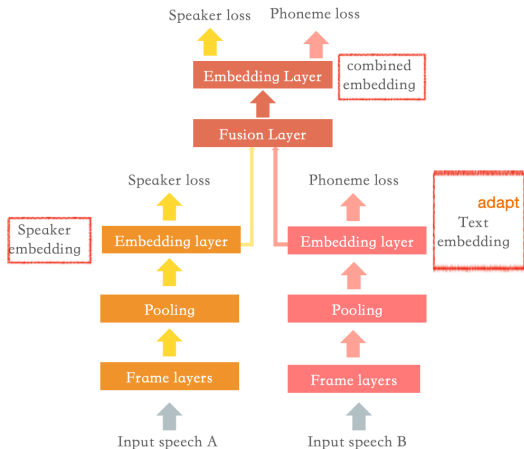
Neural Factor Analysis: Neglect the phoneme variations by additional alignment ²⁹



²⁹Lin W W, He C H, .et, Self-supervised Neural Factor Analysis for Disentangling Utterance-level Speech Representations, ICML 2023

Speech representation disentanglement

Decouple and Reorganization of Phonetic Information³⁰

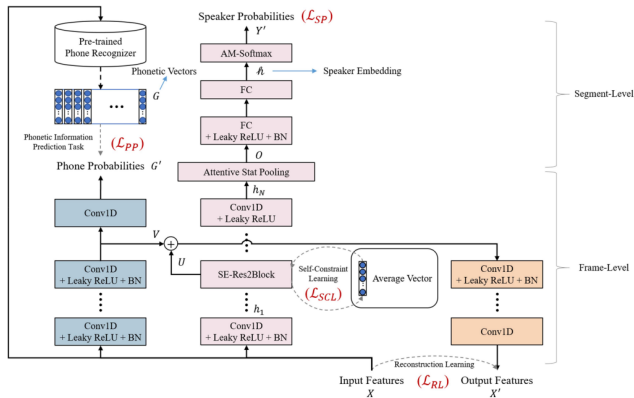


- ▶ Segment-level reconstruction
- ▶ Decoupling the speaker and text information
- ▶ For the text-independent task, we neglect the text information
- ▶ For the text-dependent task, we use the combined embedding
- ▶ **Text-adaptive task:** Modify the text information in the embedding while keeping the speaker identity. (Change the enrollment keyword)

³⁰Yang Y*, Wang S*, Gong X, et al. Text adaptation for speaker verification with speaker-text factorized embeddings. ICASSP 2020

Speech representation disentanglement

Decouple and Reorganization of Phonetic Information³¹



- ▶ Frame-level Reconstruction
- ▶ Center frame-level speaker representations towards its mean
- ▶ Coarse-grained phoneme categories (Vowel, semi-vowel, affricate, ...)

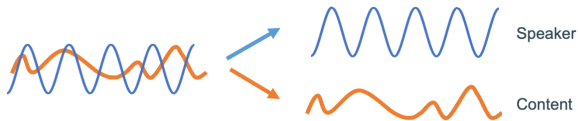
Fig. 3. The architecture of the proposed DROP-TDNN x-vector system. DROP-TDNN consists of three training procedures, including phonetic information prediction, reconstruction and speaker recognition.

³¹Hong Q B, Wu C H, Wang H M. Decomposition and Reorganization of Phonetic Information for Speaker Embedding Learning. TASLP 2023

Speech representation disentanglement

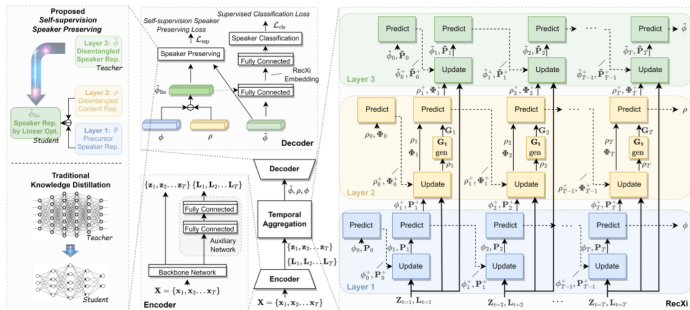
RecXi with multiple Gaussian Inference³²

Disentangle **Static** and **Dynamic** components in Speech



by three **Gaussian** inference layers

and a novel **speaker preserving self-supervision**



³²Liu T C, Disentangling Voice and Content with Self-Supervision for Speaker Recognition

Speech representation disentanglement

Codec Approach: Speech Tokenizer



Ensure the first layer representations contain content-related information, the subsequent residual layers will naturally fill in the gaps with remaining details—specifically, modeling the paralinguistic information.³³

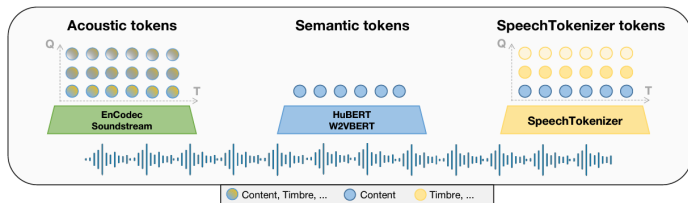


Figure 1: Illustration of information composition of different discrete speech representations. Speech tokens are represented as colored circles and different colors represent different information.

³³Zhang X, Zhang D, Li S, et al. SpeechTokenizer: Unified Speech Tokenizer for Speech Large Language Models[J]. arXiv preprint arXiv:2308.16692, 2023.

Speech representation disentanglement

Codec Approach: Speech Tokenizer

Semantic Distillation to enable the disentanglement

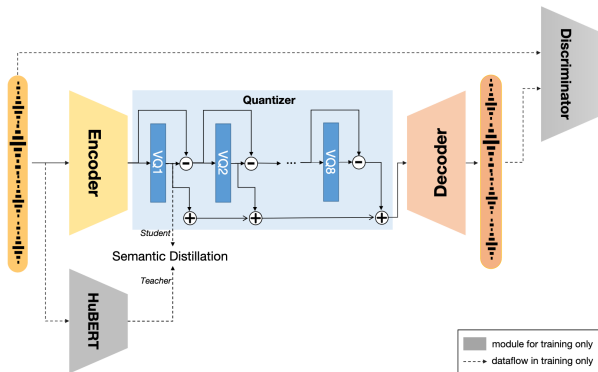


Figure 2: Illustration of SpeechTokenizer framework.

- ▶ **Continuous distillation** Output of hubert's 9-th layer/Average across all layers

$$\mathcal{L}_{\text{distill}} = \frac{1}{T} \sum_{t=1}^T \log \sigma(\cos(Aq_1^t, s^t))$$

- ▶ **Discrete distillation** pseudo-label prediction

$$\mathcal{L}_{\text{distill}} = -\frac{1}{T} \sum_{t=1}^T u^t \log(\text{Softmax}(Aq_1^t))$$



Assume hubert is a perfect semantic encoder

Speaker Modeling: Background, Applications and Trends
Discriminative Speaker Representation Learning
Self-supervised based Speaker Representation Learning
Multi-modal Speaker Representation Learning
Efficiency and Robustness
Towards the Interpretability
Beyond Speaker Recognition
Practice: Speaker Representation Learning with Wespeaker

Explore model capacity via probing tasks

Analyze the information encoded



Assumption: If a certain attribute is encoded in the speaker representation, the accuracy of a classifier predicting this property depends on how well it's embedded.^{34, 35, 36, 37}

- ▶ Speaker-related attributes: identity, gender, and speaking rate.
- ▶ Text-related factors: spoken terms, word order, and utterance length.
- ▶ Channel-related elements include the handset ID and noise type.

³⁴Wang, Qian, and Yu, "What does the speaker embedding encode?"

³⁵Belinkov and Glass, "Analyzing hidden representations in end-to-end automatic speech recognition systems".

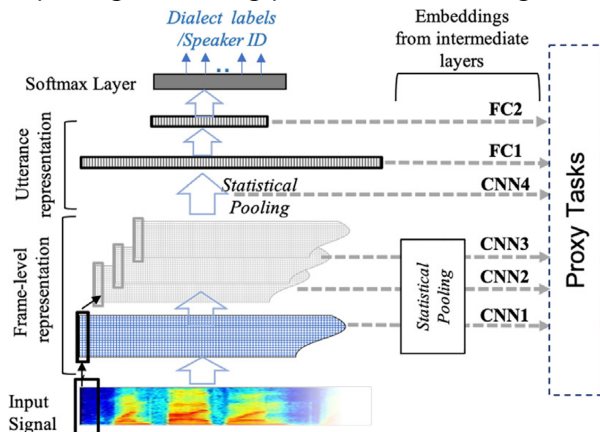
³⁶Raj et al., "Probing the information encoded in x-vectors".

³⁷Zhao et al., "Probing Deep Speaker Embeddings for Speaker-related Tasks".

Explore model capacity via probing tasks

Analyze the information encoded

Illustration of the paradigm: Probing pretrained embeddings with proxy tasks³⁸

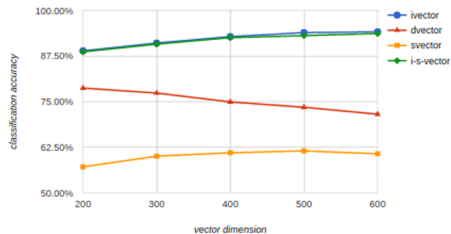


³⁸Chowdhury, Durrani, and Ali, "What do end-to-end speech models learn about speaker, language and channel information? a layer-wise and neuron-level analysis".

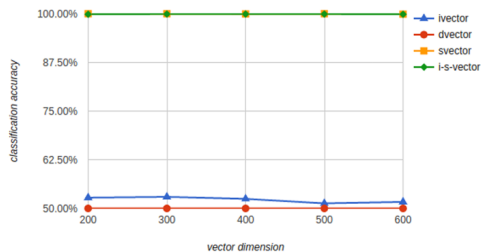
Explore model capacity via probing tasks

Analyze the information encoded

Examples of the speaker identity task and word order task³⁹



Speaker identity task



Word order task

³⁹Wang, Qian, and Yu, "What does the speaker embedding encode?"

Measuring the importance through visualization

In the context of speaker modeling, f is the speaker classifier, c represents the class, θ represents the trainable model parameters.

$$y^c = f_c(x; \theta)$$

For the k -th activation map A^k (e.g. k represents the k -th channel), each entry w_{ij}^{kc} is defined as

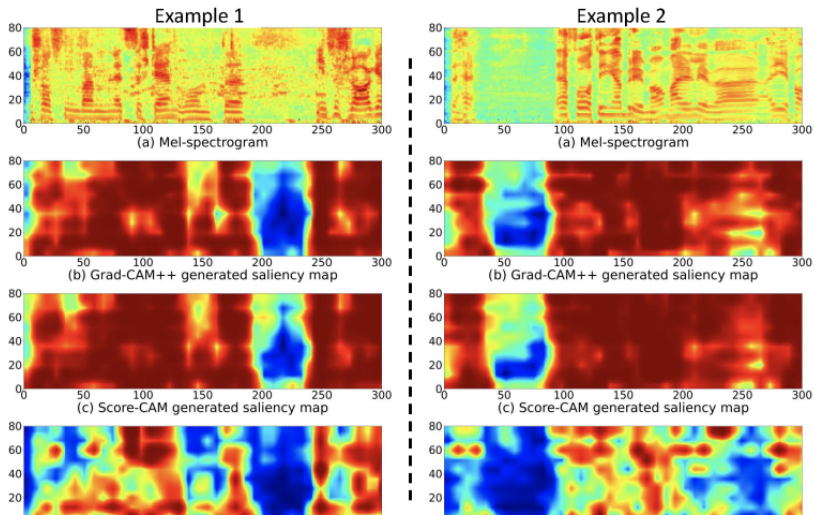
$$w_{ij}^{kc} = \text{ReLU} \left(\frac{\partial y^c}{\partial A_{ij}^k} \right)$$

Saliency map is defined as the linear combination

$$S_{ij}^c = \text{ReLU} \left(\sum_k w_{ij}^{kc} \cdot A_{ij}^k \right)$$

Measuring the importance through visualization

Visualization in Speaker Recognition^{40, 41}



Speaker Modeling: Background, Applications and Trends
Discriminative Speaker Representation Learning
Self-supervised based Speaker Representation Learning
Multi-modal Speaker Representation Learning
Efficiency and Robustness
Towards the Interpretability
Beyond Speaker Recognition
Practice: Speaker Representation Learning with Wespeaker

New paradigm to model speakers

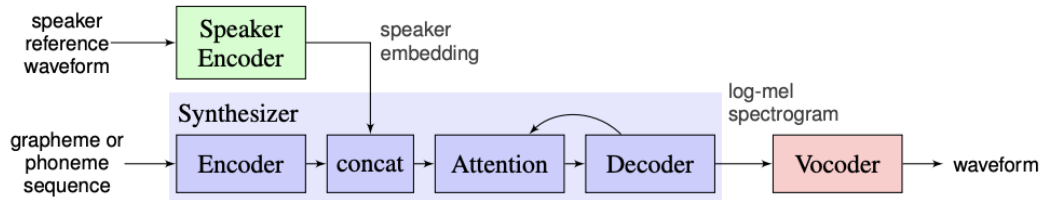
Different tasks, different approaches



1. Pretrained speaker embeddings as additional inputs
2. Joint training to learn task-specific embeddings
3. Implicit speaker modeling

New paradigm to model speakers

Example: Explicit speaker modeling for Zero-shot TTS^{42, 43, 44}



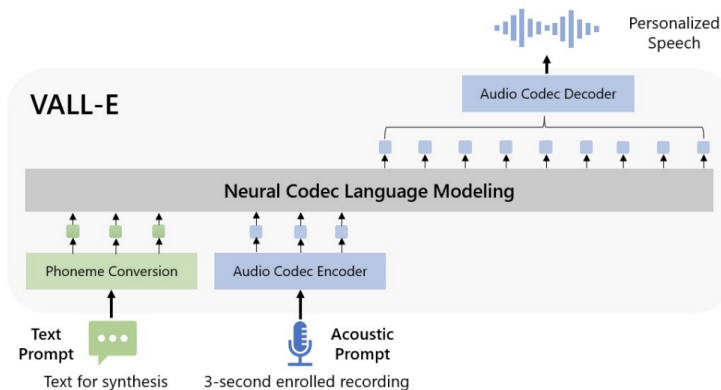
⁴²Jia et al., "Transfer learning from speaker verification to multispeaker text-to-speech synthesis".

⁴³Casanova et al., "Yourtts: Towards zero-shot multi-speaker tts and zero-shot voice conversion for everyone".

⁴⁴Wu et al., "Adaspeech 4: Adaptive text to speech in zero-shot scenarios".

New paradigm to model speakers

Example: Implicit speaker modeling for Zero-shot TTS^{45, 46, 47}



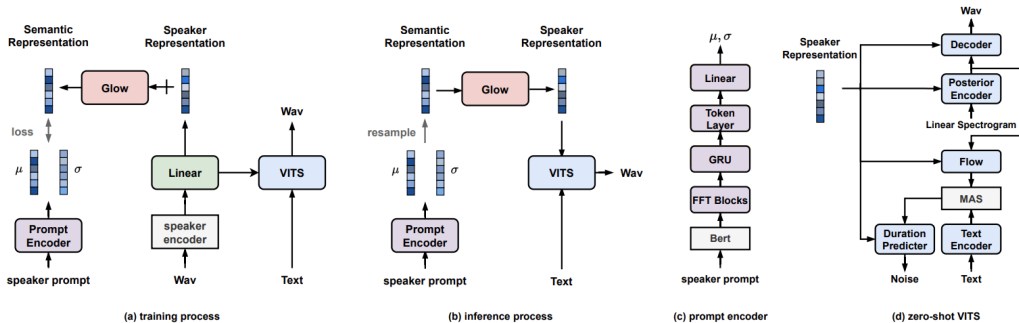
⁴⁵Wang et al., "Neural codec language models are zero-shot text to speech synthesizers".

⁴⁶Du et al., "UniCATS: A Unified Context-Aware Text-to-Speech Framework with Contextual VQ-Diffusion and Vocoding".

⁴⁷Le et al., "Voicebox: Text-guided multilingual universal speech generation at scale".

New paradigm to model speakers

Example: towards controlability and new voice generation^{48, 49, 50, 51}



⁴⁸Zhang et al., "PromptSpeaker: Speaker Generation Based on Text Descriptions".

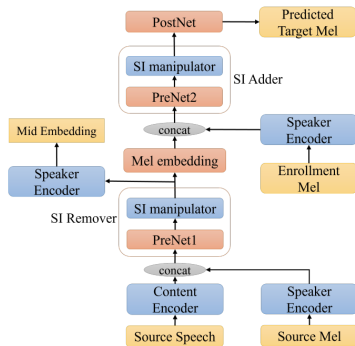
⁴⁹Stanton et al., "Speaker generation".

⁵⁰Shimizu et al., "PromptTTS++: Controlling Speaker Identity in Prompt-Based Text-to-Speech Using Natural Language Descriptions".

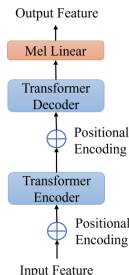
⁵¹Bilinski et al., "Creating new voices using normalizing flows".

New paradigm to model speakers

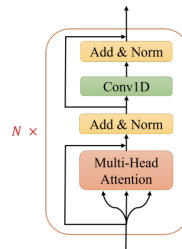
Example: Explicit speaker modeling for zero-shot voice conversion⁵²⁵³⁵⁴



(a) Proposed System



(b) SI Manipulator



(c) Encoder and decoder architecture

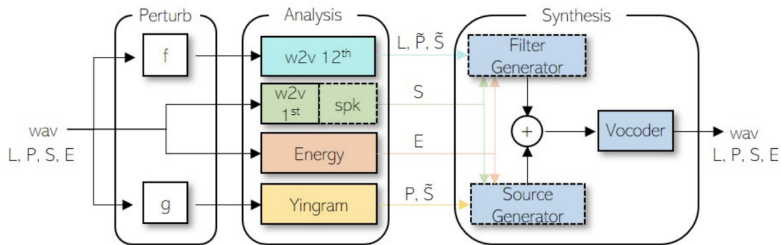
⁵²Zhang et al., "SIG-VC: A Speaker Information Guided Zero-Shot Voice Conversion System for Both Human Beings and Machines".

⁵³Chen and Duan, "ControlVC: Zero-Shot Voice Conversion with Time-Varying Controls on Pitch and Rhythm".

⁵⁴Hussain et al., "ACE-VC: Adaptive and Controllable Voice Conversion Using Explicitly Disentangled Self-Supervised Speech Representations".

New paradigm to model speakers

Example: Implicit speaker modeling for zero-shot voice conversion^{55, 56, 57}



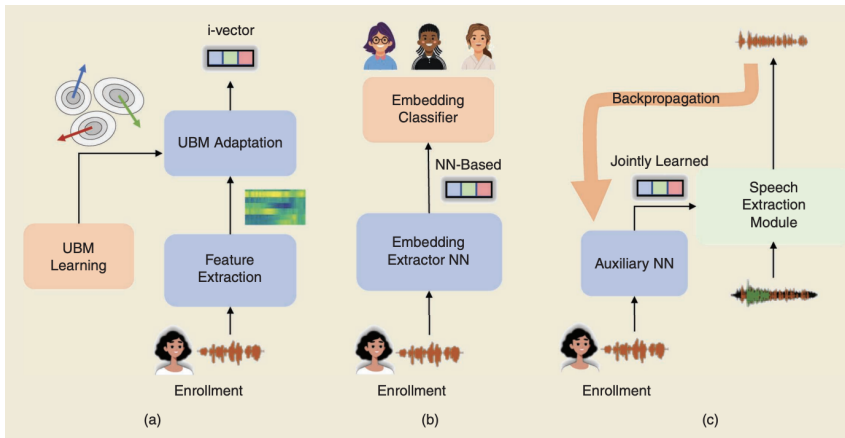
⁵⁵Choi et al., "Neural analysis and synthesis: Reconstructing speech from self-supervised representations".

⁵⁶Wu and Lee, "One-shot voice conversion by vector quantization".

⁵⁷Wu, Chen, and Lee, "Vqvc+: One-shot voice conversion by vector quantization and u-net architecture".

New paradigm to model speakers

Example: Explicit speaker modeling for target speaker extraction⁵⁸⁵⁹⁶⁰



⁵⁸Zmolikova et al., "Neural Target Speech Extraction: An overview".

⁵⁹Delcroix et al., "Single channel target speaker extraction and recognition with speaker beam".

⁶⁰Ge et al., "Spex+: A complete time domain speaker extraction network".

New paradigm to model speakers

Example: Implicit speaker modeling for target speaker extraction^{61, 62}

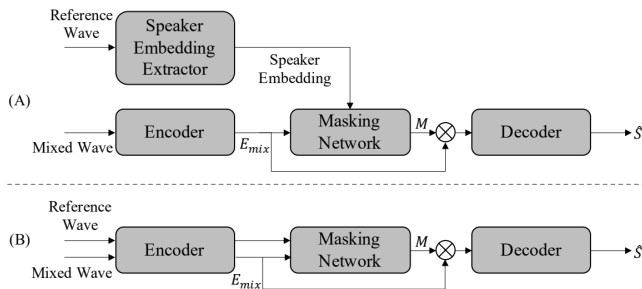


Figure 1: (A) is the diagram of a typical time-domain target speaker extraction method. (B) is the diagram of our proposed method. \otimes is an operation for element-wise product.

⁶¹Zeng et al., "SEF-Net: Speaker Embedding Free Target Speaker Extraction Network".

⁶²Yang et al., "Target Speaker Extraction with Ultra-Short Reference Speech by VE-VE Framework"

Speaker Modeling: Background, Applications and Trends
Discriminative Speaker Representation Learning
Self-supervised based Speaker Representation Learning
Multi-modal Speaker Representation Learning
Efficiency and Robustness
Towards the Interpretability
Beyond Speaker Recognition
Practice: Speaker Representation Learning with Wespeaker

Wespeaker is a speaker embedding learning toolkit designed for **research** and **production** purposes, it characterized by

- ▶ Lightweight codebase
- ▶ SOTA performance
- ▶ Discriminative and SSL based paradigms
- ▶ Runtime/Deployment support
- ▶ Adopted by research groups from both companies and academic institutions:

- ▶ Tencent
- ▶ Meituan
- ▶ China Telecom
- ▶ NVIDIA

- ▶ Tsinghua University
- ▶ The Chinese University of Hong Kong (Shenzhen)
- ▶ Shanghai Jiao Tong University
- ▶ University of Science and Technology of China
- ▶ National University of Singapore
- ▶ Institute for Infocomm Research (I2R)
- ▶ Brno University of Technology.

- ▶ Data Preparation
 - ▶ Data Downloading
 - ▶ Formating
 - ▶ Transformation
- ▶ Model Training
 - ▶ On-the-fly Data Augmentation
 - ▶ Model Selection
 - ▶ Large-margin Fine-tuning
- ▶ Backend Scoring
 - ▶ As-norm
 - ▶ PLDA / A-PLDA
 - ▶ Score Calibration (coming soon)

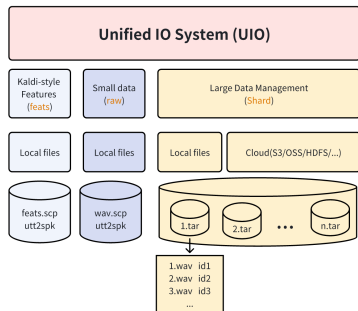


Figure: Unified I/O system

Unified I/O system

- ▶ Also adopted in wenet ASR toolkit
- ▶ Inspired by webdataset and tfrecord

Idea

- ▶ Raw: load wav and label files from disk (small data)
- ▶ Shard:
 - ▶ Pack a set of small files into a bigger shard
 - ▶ Read and decompress the shard files on-the-fly
- ▶ Feat: Compatible with kaldi-style feature files
- ▶ Effectively loading large-scale datasets

Step 1: Download and prepare metadata

```
if [ ${stage} -le 1 ] && [ ${stop_stage} -ge 1 ]; then
  echo "Prepare datasets ..."
  ./local/prepare_data.sh --stage 2 --stop_stage 4 --data ${data}
fi
```

Step 2: Covert train and test data

```
if [ ${stage} -le 2 ] && [ ${stop_stage} -ge 2 ]; then
  echo "Covert train and test data to ${data_type}..."
  for dset in vox2_dev vox1; do
    if [ ${data_type} == "shard" ]; then
      python tools/make_shard_list.py --num_utts_per_shard 1000 \
        --num_threads 16 \
        --prefix shards \
        --shuffle \
        ${data}/${dset}/wav.scp ${data}/${dset}/utt2spk \
        ${data}/${dset}/shards ${data}/${dset}/shard.list
    else
      python tools/make_raw_list.py ${data}/${dset}/wav.scp \
        ${data}/${dset}/utt2spk ${data}/${dset}/raw.list
    fi
  done
fi
```

Step3: Start training

```

if [ ${stage} -le 3 ] && [ ${stop_stage} -ge 3
]; then
echo "Start training ..."
num_gpus=$(echo $gpus | awk -F ' ' '{print NF
}')
torchrun --standalone --nnodes=1 --
nproc_per_node=$num_gpus \
wespeaker/bin/train.py --config $config \
--exp_dir ${exp_dir} \
--gpus $gpus \
--num_avg ${num_avg} \
--data_type "${data_type}" \
--train_data ${data}/vox2_dev/${data_type
}.List \
--train_label ${data}/vox2_dev/utt2spk \
--reverb_data ${data}/rirs/lmdb \
--noise_data ${data}/musan/lmdb \
${checkpoint:+--checkpoint $checkpoint}
fi

```

Dataset Config:

```

dataset_args:
speed_perturb: True
num_frms: 200
aug_prob: 0.6
# prob to add reverb
& noise aug per
sample
fbank_args:
num_mel_bins: 80
frame_shift: 10
frame_length: 25
dither: 1.0
spec_aug: False
spec_aug_args:
num_t_mask: 1
num_f_mask: 1
max_t: 10
max_f: 8
prob: 0.6

```

Data Augmentation:

```

# add noise
dataset = Processor(
dataset, processor
.add_reverb_noise,
reverb_data,
noise_data,
resample_rate,
aug_prob)
# speed perturb
dataset = Processor(
dataset, processor
.speed_perturb,
len(spk2id_dict))
# specaug
dataset = Processor(
dataset, processor
.spec_aug, **
configs['
spec_aug_args'])

```

Model Arch:

- ▶ ResNet Series
- ▶ TDNN
- ▶ ECAPA-TDNN
- ▶ RepVGG
- ▶ CAM++

Pooling Methods:

- ▶ TSTP
- ▶ ASTP
- ▶ MQMHASTP

Loss Function:

- ▶ add_margin
- ▶ arc_margin
- ▶ sphere
- ▶ sphereface2
- ▶ intertopk
- ▶ subcenter

Model Config:

```
model: ResNet34
      # ECAPA, CAMPPlus,
      # REPVGG,
      # ResNet152
model_args:
  feat_dim: 80
  embed_dim: 256
  pooling_func: "TSTP"
              # TSTP, ASTP,
              # MQMHASTP
  two_emb_layer: False
  projection_args:
    project_type: "
      arc_margin"
    # add_margin,
    # arc_margin,
    # sphere,
    # sphereface2,
    # softmax,
    # aam_intertopk
  scale: 32.0
```

Back-end Support:

- ▶ Cosine
- ▶ PLDA
- ▶ Adapt-PLDA

Others:

- ▶ Score normalization
- ▶ QMF based Calibration

Scoring:

```
if [ ${stage} -le 5 ] && [ ${stop_stage} -ge 5
]; then
echo "Score ..."
local /score.sh \
  --stage 1 --stop-stage 2 \
  --data ${data} \
  --exp_dir $exp_dir \
  --trials "$trials"
fi

if [ ${stage} -le 6 ] && [ ${stop_stage} -ge 6
]; then
echo "Score norm ..."
local /score_norm.sh \
  --stage 1 --stop-stage 3 \
  --score_norm_method $score_norm_method \
  --cohort_set_vox2_dev \
  --top_n $top_n \
  --data ${data} \
  --exp_dir $exp_dir \
  --trials "$trials"
fi
```

Onnx Inference Demo

To use the pretrained model in `pytorch` format, please directly refer to the `run.sh` in corresponding recipe.

As for extracting speaker embeddings from the `onnx` model, the following is a toy example.

```
# Download the pretrained model in onnx format and save it as onnx_path
# wav_path is the path to your wave file (16k)
python wespeaker/bin/infer_onnx.py --onnx_path $onnx_path --wav_path $wav_path
```

You can easily adapt `infer_onnx.py` to your application, a speaker diarization example can be found in [the voxconverse recipe](#)

Model List

Datasets	Languages	Checkpoint (pt)	Runtime Model (onnx)
VoxCeleb	EN	ResNet34 / ResNet34_LM	ResNet34 / ResNet34_LM
VoxCeleb	EN	ResNet152_LM	ResNet152_LM
VoxCeleb	EN	ResNet221_LM	ResNet221_LM
VoxCeleb	EN	ResNet293_LM	ResNet293_LM
VoxCeleb	EN	CAM++ / CAM++_LM	CAM++ / CAM++_LM
CNCeleb	CN	ResNet34 / ResNet34_LM	ResNet34 / ResNet34_LM

Figure: Pretrained Model List

Export Jit:

```
if [ ${stage} -le 7 ] && [ ${stop_stage} -ge 7
]; then
echo "Export the best model ..."
python wespeaker/bin/export_jit.py \
  --config $exp_dir/config.yaml \
  --checkpoint $exp_dir/models/avg_model.pt \
  --output_file $exp_dir/models/final.zip
fi
```

Export Onnx:

```
exp=exp # Change it to your experiment dir
onnx_dir=onnx
python wespeaker/bin/export_onnx.py \
  --config $exp/config.yaml \
  --checkpoint $exp/avg_model.pt \
  --output_model $onnx_dir/final.onnx
```

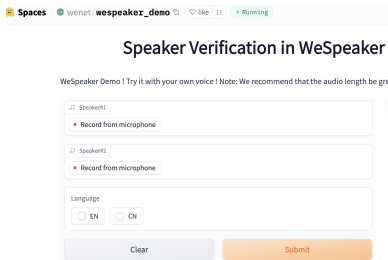


Figure: Wespeaker Demo Page

Command-line usage:

```
wespeaker --task embedding --audio_file audio.wav --  
output_file embedding.txt -g 0  
wespeaker --task embedding_kaldi --wav_scp wav.scp --  
output_file /path/to/embedding -g 0  
wespeaker --task similarity --audio_file audio.wav --  
audio_file2 audio2.wav -g 0
```

Python programming usage:

```
import wespeaker  
  
model = wespeaker.load_model('chinese')  
# set_gpu to enable the cuda inference, number < 0 means  
# using CPU  
model.set_gpu(0)  
embedding = model.extract_embedding('audio.wav')  
utt_names, embeddings = model.extract_embedding_list('wav.scp')  
similarity = model.compute_similarity('audio1.wav', 'audio2.wav')  
diar_result = model.diarize('audio.wav')
```


Table: Supervised results achieved using different architectures on the VoxCeleb dataset, “dev” of part 2 is used as the training set

Literature/Toolkits	Architecture	voxceleb1_O		voxceleb1_E		voxceleb1_H	
		EER(%)	minDCF	EER(%)	minDCF	EER(%)	minDCF
IDLab VoxSRC 2020 ⁶³	ECAPA-TDNN	0.870	0.107	1.120	0.132	2.120	0.210
BUT VoxSRC 2019 ⁶⁴	ResNet34	1.310	0.154	1.380	0.163	2.500	0.233
AsvSubtools ¹	ECAPA-TDNN	0.856	-	-	-	-	-
	Conformer	0.792	-	-	-	-	-
SpeechBrain ²	TDNN	3.23	-	-	-	-	-
	ECAPA-TDNN	0.90	-	-	-	-	-
	ECAPA-TDNN *	1.30	-	1.98	-	3.62	-
Nemo ³	TDNN	1.96	-	-	-	-	-
	ECAPA-TDNN	0.92	-	-	-	-	-
	titanet_large	0.66	-	-	-	-	-
Wespeaker	TDNN	1.590	0.166	1.641	0.170	2.726	0.248
	ECAPA-TDNN	0.728	0.099	0.929	0.100	1.721	0.169
	CAM++	0.654	0.087	0.805	0.092	1.576	0.164
	RepVGG	0.750	0.083	0.846	0.090	1.495	0.141
	ResNet34	0.723	0.069	0.867	0.097	1.532	0.146
	ResNet50	0.803	0.061	0.887	0.092	1.519	0.136
	ResNet101	0.542	0.052	0.758	0.079	1.398	0.128
	ResNet152	0.495	0.033	0.685	0.069	1.205	0.105
	ResNet221	0.505	0.045	0.676	0.067	1.213	0.111
ResNet293	0.447	0.043	0.657	0.066	1.183	0.111	

⁶³Desplanques, Thienpondt, and Demuyck, “Ecapa-tdnn: Emphasized channel attention, propagation and aggregation in tdnn based speaker verification”.

⁶⁴Zeinali et al., “But system description to voxceleb speaker recognition challenge 2019”.

Table: Results on the CNCeleb evaluation set

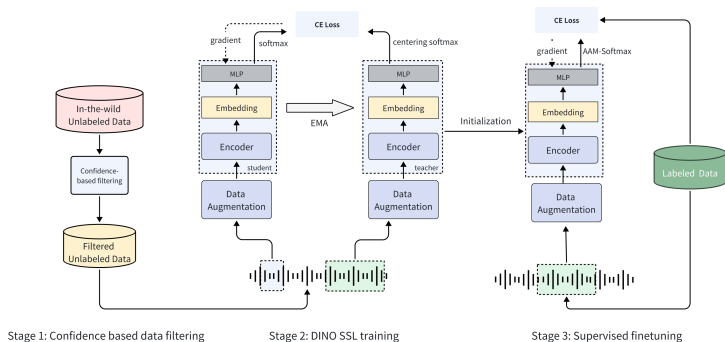
Toolkits	Architecture	EER(%)	minDCF
ASVSubtools	ResNet34	9.141	0.463
Wespeaker	TDNN	8.960	0.446
	ECAPA-TDNN	7.395	0.372
	CAM++	7.052	0.368
	ResNet34	6.492	0.354
	ResNet221	5.655	0.330

Table: Performance (EER%) of SSL-based systems on the VoxCeleb evaluation set

Toolkits	Paradigm	Architecture	VoxCeleb1_O	VoxCeleb1_E	VoxCeleb1_H
3Dspeaker	RDINO	ECAPA-TDNN (C1024)	3.16	-	-
	SimCLR	ECAPA-TDNN	8.523	9.417	14.907
	MoCo	ECAPA-TDNN	8.709	9.287	14.756
wespeaker	DINO	ResNet34	3.170	3.324	5.821
	DINO	ECAPA-TDNN (C512)	3.016	3.093	5.538
	DINO	ECAPA-TDNN (C1024)	2.627	2.665	4.644

Wespeaker

A comprehensive example of using Wespeaker^{65, 66}



Step 1: Clustering-based speaker diarization system, filter out low-quality segments.

Step 2: Train a DINO system on the filtered data

Step 3: Finetuning the pretrained DINO system in a supervised setup.

⁶⁵Wang et al., "Leveraging In-the-Wild Data for Effective Self-Supervised Pretraining in Speaker Recognition".

⁶⁶Yu et al., "AutoPrep: An Automatic Preprocessing Framework for In-the-Wild Speech Data".

Table: Comparison of performance on CNCeleb-Eval with other pretrain-finetune methods.

System	Pretraining Configurations			Finetuning Configurations		EER(%)	MinDCF
	Data	Model	Role	Data	Model		
⁶⁷	VoxCeleb2	ECAPA-TDNN	Init	CNCeleb1	ECAPA-TDNN	10.65	-
⁶⁸	VoxCeleb2	ECAPA-TDNN	Init	CNCeleb1	ECAPA-TDNN	8.710	0.422
⁶⁹	VoxCeleb2	ECAPA-TDNN	Init	CNCeleb1	ECAPA-TDNN	10.03	0.539
⁷⁰	CNCeleb1	HuBERT (94.6M)	Frontend	CNCeleb1	HuBERT + ECAPA-TDNN	10.86	-
⁷	CNCeleb-Train	HuBERT (94.6M)	Frontend	CNCeleb-Train	HuBERT + ECAPA-TDNN	8.890	-
⁷	CNCeleb-Train	Conformer (172.2M)	Frontend	CNCeleb-Train	Conformer + MHFA	7.730	0.406
⁷¹ *	Mix 94k hr	WavLM (94.7M)	Frontend	VoxCeleb2 + CNCeleb-Train	WavLM+MAM+MHFA	6.890	0.378
⁷² **	WenetSpeech	Conformer (18.8M)	Init	CNCeleb-Train	Conformer	7.420	0.443
Ours	WenetSpeech	ECAPA-TDNN	Init	CNCeleb1	ECAPA-TDNN	7.373	0.383
Ours	+ filtering	ECAPA-TDNN	Init	CNCeleb1	ECAPA-TDNN	7.339	0.377
Ours	WenetSpeech	ECAPA-TDNN	Init	CNCeleb-Train	ECAPA-TDNN	6.738	0.338
Ours	+ filtering	ECAPA-TDNN	Init	CNCeleb-Train	ECAPA-TDNN	6.474	0.331

⁶⁷Heo et al., "Self-supervised curriculum learning for speaker verification".⁶⁸Kang et al., "Augmentation adversarial training for self-supervised speaker representation learning".⁶⁹Han et al., "Improving dino-based self-supervised speaker verification with progressive cluster-aware training".⁷⁰Peng et al., "Improving speaker verification with self-pretrained transformer models".⁷¹Peng et al., "Parameter-efficient transfer learning of pre-trained Transformer models for speaker verification using adapters".⁷²Liao et al., "Towards a unified conformer structure: from asr to asv task".

Conclusions and Takeaway

- ▶ Speaker modeling is not all about speaker recognition.
- ▶ Speaker modeling is more than embedding learning.
- ▶ Customize the speaker modeling approach for the specific task.
- ▶ Try wespeaker! <https://github.com/wenet-e2e/wespeaker>

Email: wsstriving@gmail.com