

Deep Speaker Representation Learning

Theory, Applications, and Practice

Shuai Wang



- ① Speaker Modeling: Background, Applications, and Trends
- ② Discriminative Speaker Representation Learning
- ③ Self-supervised Speaker Representation Learning
- ④ Multi-modal Speaker Representation Learning
- ⑤ Robustness and Interpretability
- ⑥ Speaker Modeling in Related Tasks
- ⑦ Practice: WeSpeaker and WeSep

Speaker Modeling



Definition

Speaker modeling aims to characterize and recognize an individual's unique traits by analyzing patterns embedded in speech signals.

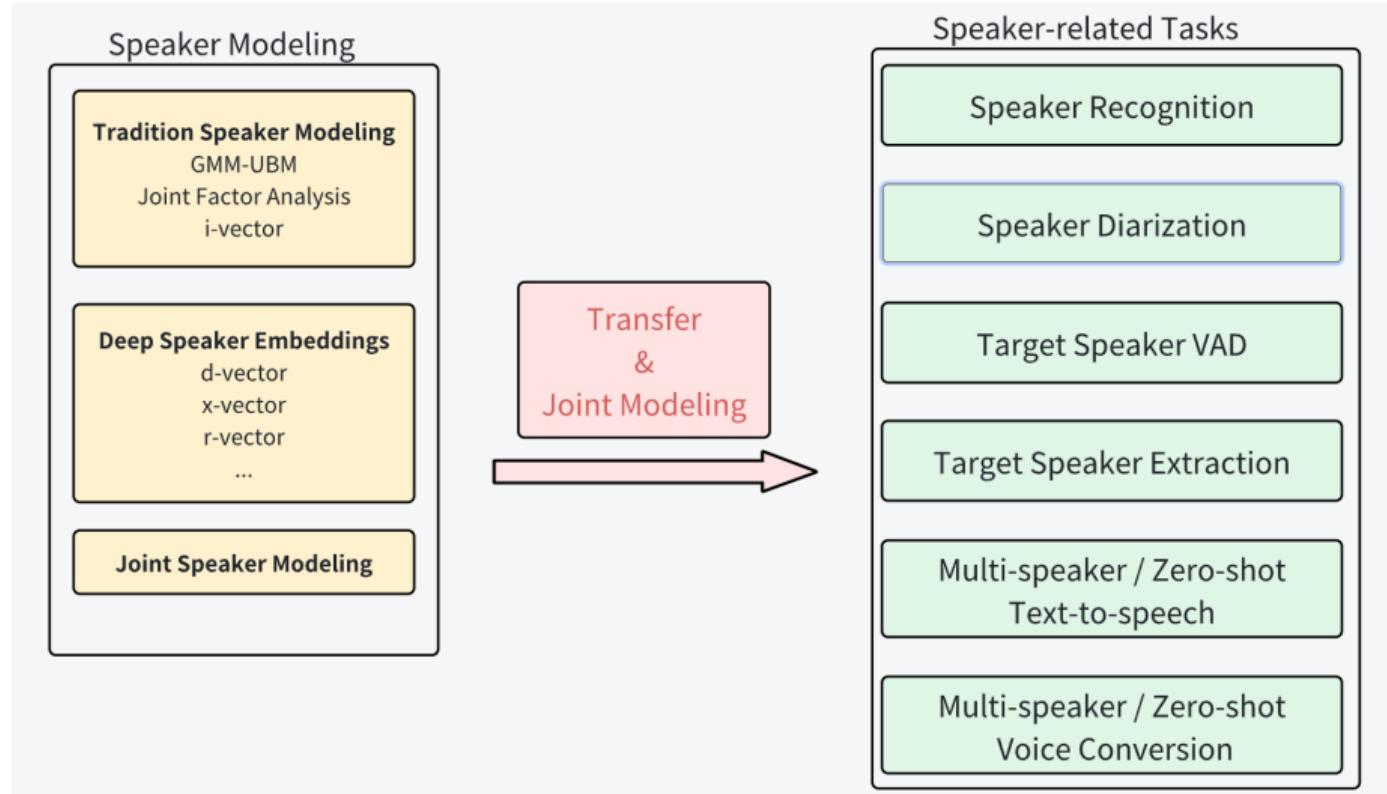
Main Applications:

- Speaker recognition
- Speaker diarization
- Voice cloning
- Speech synthesis
- Target speaker extraction

Real-world Impact:

- Biometric authentication
- Surveillance systems
- Personalized services
- Forensic analysis
- Privacy protection

Applications of Speaker Modeling



Applications of Speaker Modeling

Embedding Extraction for Speaker Verification

Speaker verification: voice as a password

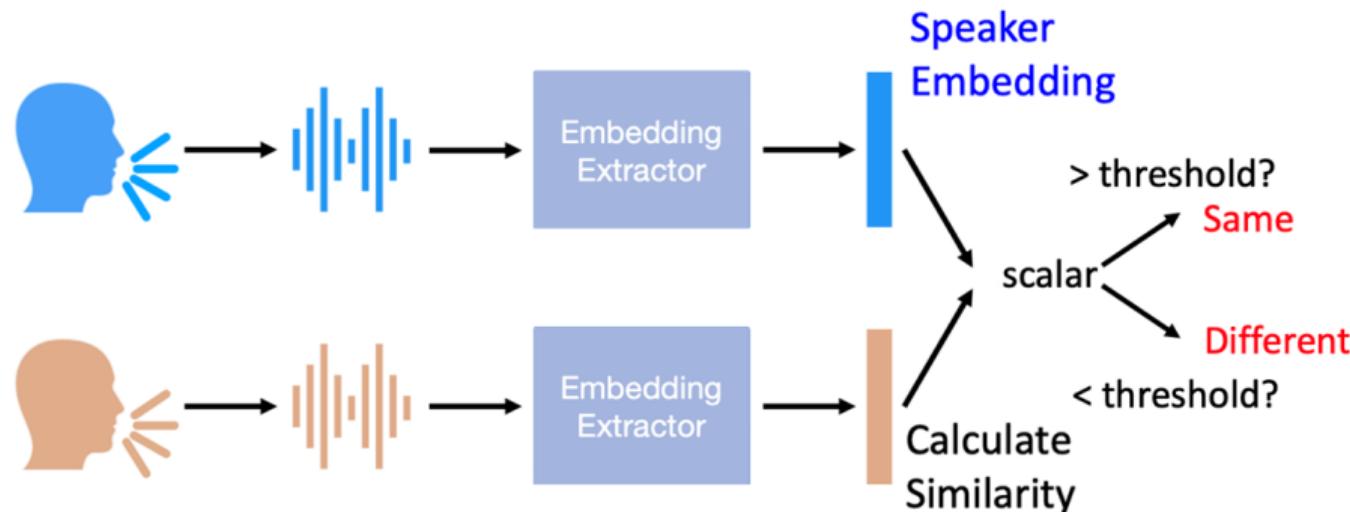
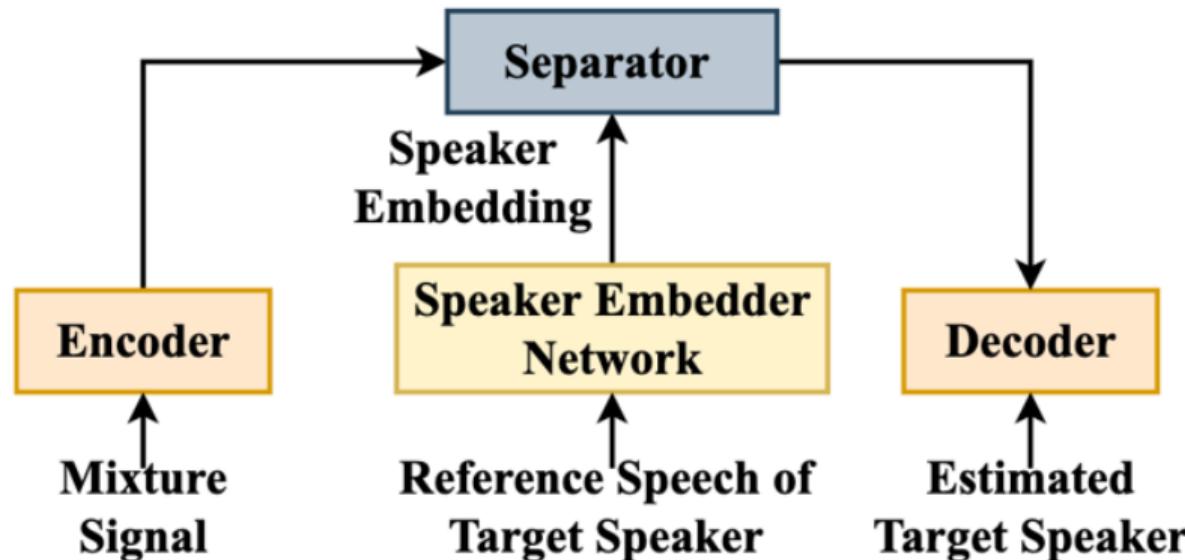


Figure adapted from Hung-Yi Lee's DLHLP20 slides¹

Applications of Speaker Modeling

Reference/Cue Modeling for Target Speech Extraction

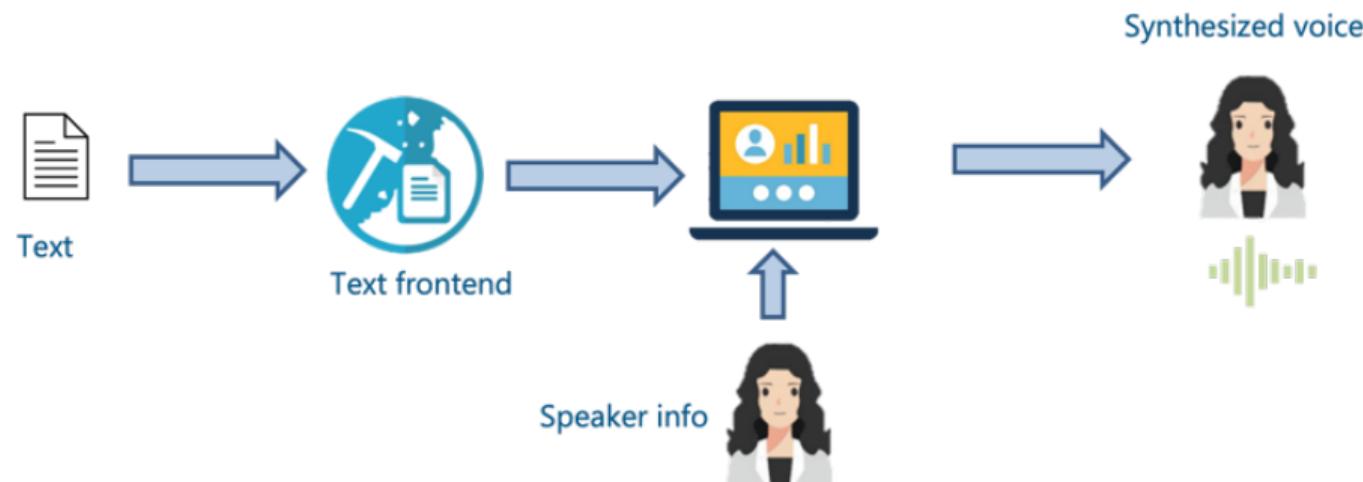
Target speech extraction: listen to the target person



Applications of Speaker Modeling

Target Speaker Identifier for TTS

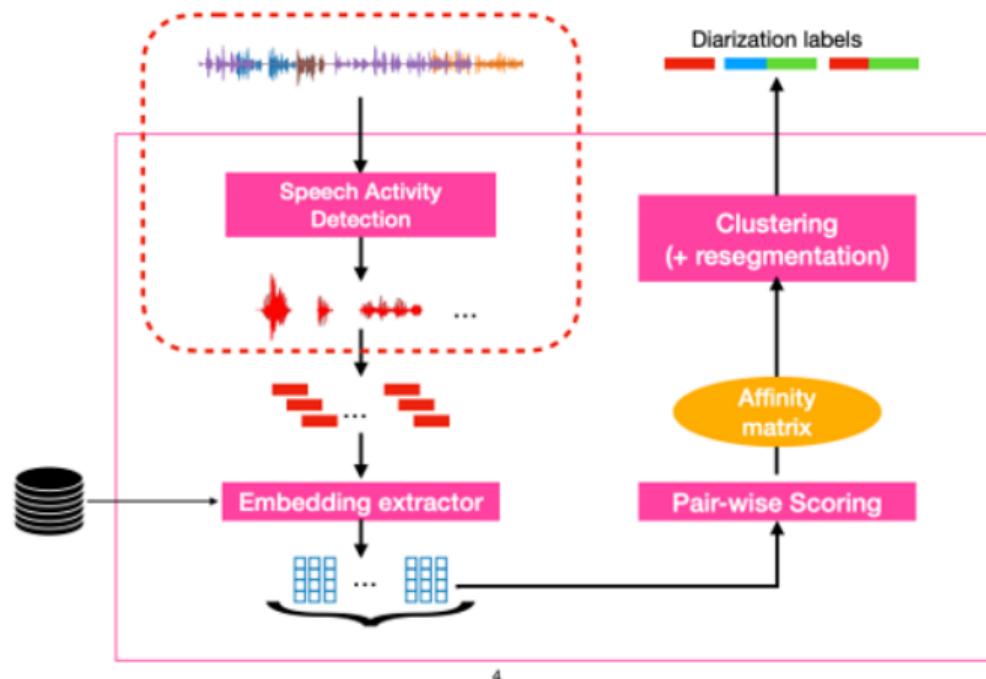
Speaker modeling for the target speaker in speech synthesis.



Applications of Speaker Modeling

Clustering-based Speaker Diarization

Speaker diarization: who spoke when?



1. From VQ to GMM (1980s–1990s)

- Vector quantization → Gaussian mixture models
- Key advance: modeling uncertainty with covariance matrices

2. From GMM-EM to GMM-UBM (2000s)

- Universal background model to improve generalization
- MAP adaptation replaces EM-only training for speaker models

3. From Supervectors to i-vectors (2010s)

- Dimensionality reduction and channel compensation
- Low-dimensional speaker representations

4. From Generative to Discriminative (2015–present)

- Deep neural networks for speaker embeddings
- End-to-end discriminative training

Methods and Representations

- Vector Quantization (VQ): discrete codebook representation; distance-based matching.
- Gaussian Mixture Models (GMM): continuous-density generative modeling, trained via EM.
- Likelihood-based scoring replaces heuristic distances; better fits acoustic feature distributions.

Limitations of VQ

- No explicit uncertainty modeling; fragile under channel/noise variations.
- Fixed, hand-crafted codebooks; limited capacity and adaptability.

Why GMM Dominated

- Soft assignment and covariance modeling capture within-speaker variability.
- Principled maximum-likelihood training; extendable to adaptation and compensation.

Key Innovations

- Universal Background Model (UBM): speaker-independent acoustic space shared by all speakers.
- MAP adaptation: derive speaker models from UBM using limited enrollment data (with relevance factor control).
- Likelihood-ratio scoring: $\log p(\mathbf{X} | \text{spk}) - \log p(\mathbf{X} | \text{UBM})$ for verification.

Practical Impact

- Robust under scarce enrollment data; improved cross-channel generalization.
- Established reproducible, scalable baselines for large-scale evaluations (telephone, microphone speech).

Speaker Modeling

GMM-UBM & GMM Supervector

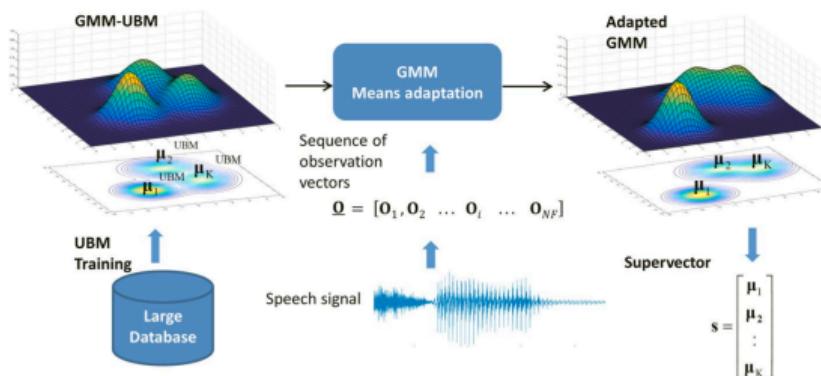


Figure: GMM-UBM^a for speaker modeling

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

Gaussian Mixture Model (GMM)^a

- $p(\mathbf{x}) = \sum_1^K c_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)$ s.t. $\sum_1^K c_k = 1$
- Any distribution can be approximated by a weighted linear combination of several Gaussians.
- When modeling speaker acoustics with GMMs, the number of Gaussians can be viewed as types of produced sounds.

Universal Background Model (UBM)

- Enrollment speech is often limited (a few seconds), making training a GMM difficult with such data.
- Train a UBM on large data and adapt it to specific speakers.

GMM Supervector

- Concatenate the mean vector of each Gaussian to represent a speaker.

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

GMM-based speaker modeling – test scoring

- In GMM-UBM systems, a likelihood ratio is typically used for scoring a test utterance.
Given an utterance Y , two hypotheses are:

$H_0 : Y$ comes from target speaker S

$H_1 : Y$ does not come from target speaker S

- The score Λ is determined by the log-likelihood ratio:

$$\Lambda = \frac{1}{T} \log \frac{p(Y | H_0)}{p(Y | H_1)} = \begin{cases} \geq \theta & \text{accept } H_0 \\ < \theta & \text{accept } H_1 \end{cases}$$

where T is the total number of frames of Y and θ is a preset threshold.

- Concretely, $p(Y | H_0)$ is the probability density of the features of Y on speaker S 's GMM, and $p(Y | H_1)$ is that on the impostor model. In GMM-UBM systems, the UBM serves as the impostor model.

From High- to Low-dimensional Space

- Supervector SVM: concatenate adapted GMM means into a very high-dimensional vector; powerful but channel-sensitive.
- Total-variability (T) model: joint factor analysis; utterance-level latent variable (i-vector) summarizing speaker and channel.
- i-vector extraction: posterior inference in $\mathcal{N}(\mathbf{m} + \mathbf{T}\mathbf{w}, \Sigma)$ to produce compact $\mathbf{w} \in \mathbb{R}^d$.

Back-ends and Compensation

- Length normalization, LDA/WCCN for inter-/intra-class variance control.
- PLDA or cosine scoring for calibrated verification.

Speaker Modeling

i-vector

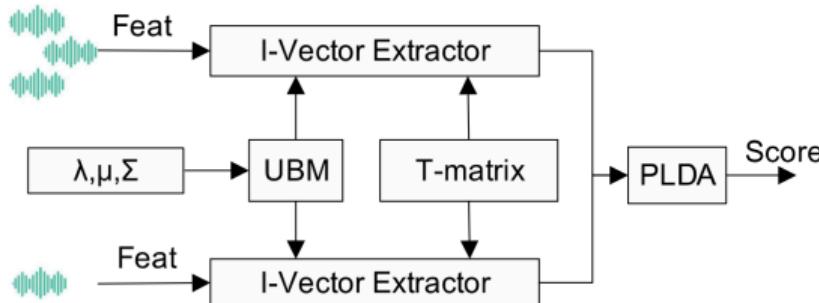


Figure: Block diagram of an i-vector based speaker recognition system

Drawbacks of GMM Supervector

- Supervectors are extremely high-dimensional (often tens of thousands), making computation challenging.

Factorize the supervector into a low-dimensional i-vector^a:

$$M(s) = m + T w(s)$$

- $M(s)$: GMM supervector of speaker s
- m : speaker-independent supervector
- T : total-variability matrix capturing all variability sources (speaker- and channel-related)
- $w(s)$: i-vector of speaker s

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

i-vector Extraction

- ① **UBM training:** train the universal background model with large data
- ② **T-matrix training:** learn the total-variability matrix \mathbf{T}
- ③ **Posterior inference:** compute $\mathbf{w} = E[\mathbf{w}|\mathbf{X}]$
- ④ **Length normalization:** $\mathbf{w}_{norm} = \frac{\mathbf{w}}{\|\mathbf{w}\|}$

Back-end Compensation

- **LDA:** linear discriminant analysis
 - Maximize between-class variance
 - Minimize within-class variance

Scoring

- **PLDA:** probabilistic linear discriminant analysis
 - Accounts for within-speaker variability
 - Provides a probabilistic interpretation

Neural Embeddings and Training

- x-vector/TDNN, ResNet/ECAPA encoders; attentive statistics pooling for temporal aggregation.
- Large-margin classification losses (AM-Softmax, AAM-Softmax) for discriminative speaker spaces.
- Data augmentation and domain adaptation for robustness (reverb, noise, codecs, channels).

Trends and Performance

- Self-supervised pretraining (e.g., wav2vec 2.0, HuBERT) as frontend or via end-to-end fine-tuning.
- Simple cosine/PLDA back-ends remain competitive with well-calibrated embeddings.
- Continuous improvements on open benchmarks (e.g., VoxCeleb) under constrained scoring (EER/minDCF).

Problem Statement

Speaker representation learning:

Given an utterance $\mathbf{O} = \{\mathbf{o}_1, \dots, \mathbf{o}_T\} \in \mathbb{R}^{T \times D}$, learn a mapping function \mathcal{F} to extract a speaker representation:

$$\mathbf{v} = \mathcal{F}(\mathbf{O}) \in \mathbb{R}^d$$

where

- \mathbf{o}_t : frame-level acoustic feature at time t
- T : number of frames
- D : feature dimension
- \mathbf{v} : fixed-length speaker embedding
- d : embedding dimension

Objective: embeddings from the same speaker should be close, and those from different speakers should be far apart.

Speaker Encoder Architectures

Frame-level vs. Segment-level Optimization

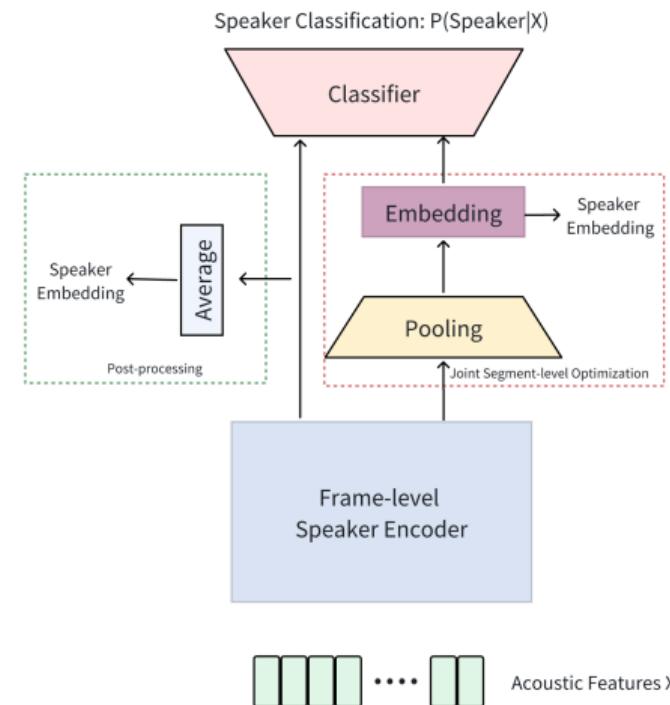
Frame-level (d-vector)

- Train with frame-level labels
- Aggregate after training
- Simple but limited performance

$$\mathbf{v} = \frac{1}{T} \sum_{t=1}^T \mathbf{f}_t \quad (1)$$

Segment-level (x-vector)

- Train with utterance-level labels
- Integrate aggregation in-network
- Better performance



Speaker Encoder Architectures

1D vs. 2D Convolution

1D Convolution

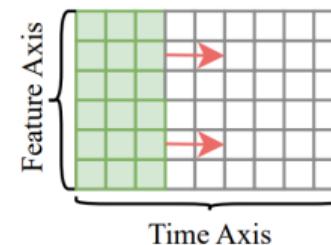
- Applied along time
- Lower computational cost
- Potentially large receptive field
- Simple architecture
- Limited frequency modeling

2D Convolution

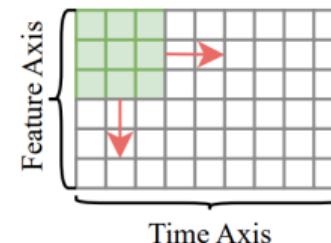
- Applied along time and frequency
- Better time-frequency modeling
- Higher computational cost
- Better performance potential



a) 1D Convolution for raw wav input



b) 1D Convolution for spectrogram input



c) 2D Convolution for spectrogram input

Speaker Encoder Architectures

d-vector

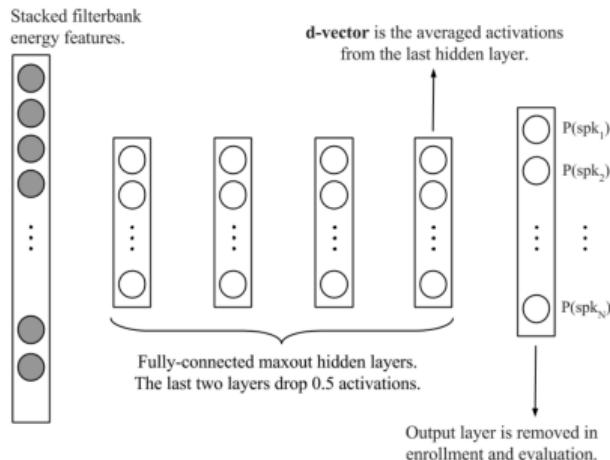


Figure: Architecture of d-vector

Labels for d-vector^a:

- Early attempt to apply DNNs to speaker information modeling
- Demonstrated good complementarity to i-vector
- Frame-level embeddings averaged to utterance vector; trained with CE on speaker IDs
- Typical encoders: TDNN/LSTM/CNN with temporal mean pooling (no attentive statistics)
- Pros: simple training, suitable for short utterances, low latency
- Limitations: phonetic content leakage; weaker than x-vector for segment aggregation

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

Speaker Encoder Architectures

x-vector

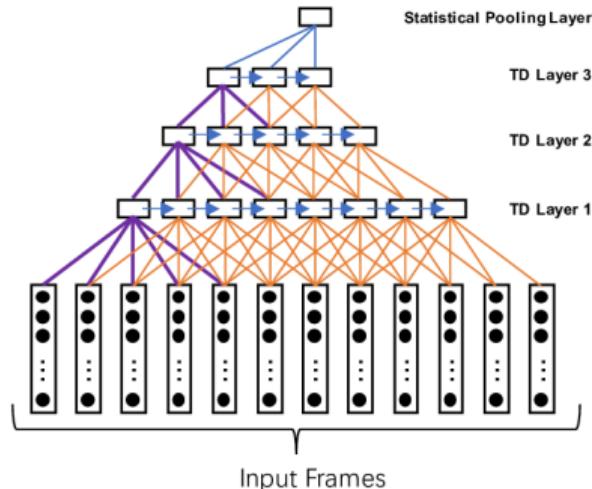


Figure: TDNN architecture used by x-vector

Labels for x-vector^a:

- First deep speaker embedding surpassing traditional methods on standard datasets (NIIST SRE)
- First to introduce segment-level optimization
- Strong variant: ECAPA-TDNN^b

Temporal pooling:

$$= \frac{1}{T} \sum_{t=1}^T \mathbf{h}_t \quad (\mathbf{h}_t \in \mathbb{R}^D)$$

Statistics pooling:

$$= \sqrt{\frac{1}{T} \sum_{t=1}^T (\mathbf{h}_t - \bar{\mathbf{h}})^{\odot 2}}$$
$$\mathbf{v} = [\quad ; \quad] \in \mathbb{R}^{2D}$$

^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".

Speaker Modeling

Data: From Labeled Studio to Unlabeled In-the-wild



Labeled data:

- Expensive to annotate
- Automatically collected speaker-labeled data such as VoxCeleb has privacy concerns
 - The VoxCeleb dataset is no longer accessible from the official website

Unlabeled data:

- Easy to obtain
- Covers a broader range of real-world conditions
- Less privacy concern

Training paradigms: from supervised to unsupervised/semi-supervised/self-supervised

Trends in Speaker Representation Learning

Models: From Shallow to Deep



- GMM and i-vector can be seen as one-layer MLPs
- d-vector, j-vector, x-vector: fewer than 10 layers
- ResNet-based models (common: 34 layers; up to 293 or 500+ in challenges)

Trends in Speaker Representation Learning

Modality: From Single-modality to Multi-modality



- Pure audio modality
- Audio-visual speaker embeddings

Trends in Speaker Representation Learning



Paradigm: From Training from Scratch to Leveraging Pretrained Models

- Train speaker-discriminative models from scratch
- Leverage large pretrained speech models such as WavLM
- Semi-supervised: iterative clustering and supervised fine-tuning

Trends in Speaker Representation Learning



Tasks: From Single-task to Cross-task

- Pretrained embeddings used across tasks
- Explicit joint optimization with task-specific objectives
- Implicit speaker modeling in related tasks

Technical Challenges:

- Channel mismatch
- Language dependency
- Short utterance handling
- Computational efficiency
- Model interpretability

Practical Challenges:

- Data privacy concerns
- Real-time processing
- Scalability issues
- Cross-domain generalization
- Evaluation standardization

- **Large-scale Pre-training:** Foundation models for speaker modeling
- **Few-shot Learning:** Rapid adaptation with limited data
- **Multimodal Integration:** Audio-visual speaker modeling
- **Edge Computing:** Efficient deployment on mobile devices
- **Federated Learning:** Privacy-preserving distributed training
- **Explainable AI:** Interpretable speaker modeling systems

- 1 Speaker Modeling: Background, Applications, and Trends
- 2 Discriminative Speaker Representation Learning
- 3 Self-supervised Speaker Representation Learning
- 4 Multi-modal Speaker Representation Learning
- 5 Robustness and Interpretability
- 6 Speaker Modeling in Related Tasks
- 7 Practice: WeSpeaker and WeSep

Problem Setup

Given an utterance $\mathbf{O} = \{\mathbf{o}_t\}_{t=1}^T \in \mathbb{R}^{T \times D}$, learn \mathcal{F} to produce a fixed-length embedding $\mathbf{v} = \mathcal{F}(\mathbf{O}) \in \mathbb{R}^d$.

- Embeddings of the same speaker should be close; different speakers should be far apart.
- Use classification-driven objectives aligned with cosine/angle scoring at inference.
- Introduce explicit margins to strengthen open-set generalization.

Supervised Loss Functions

Speaker Verification vs. ASR (Motivation)



ASR (classical acoustic modeling)

- Closed-set at inference: phone/state labels are fixed.
- Objective: maximize classification accuracy over a known label set.

Speaker verification

- Open-set at inference: speakers are unseen during training.
- Requires compact clusters and clear margins under cosine/angles.

Supervised Loss Functions

Softmax Baseline for Speaker Classification



Model

Embedding $\mathbf{v} = \mathcal{F}(\mathbf{O}) \in \mathbb{R}^d$, classifier $W = [\mathbf{w}_1, \dots, \mathbf{w}_C] \in \mathbb{R}^{d \times C}$ with bias \mathbf{b} .

- Logits: $s_j = \mathbf{w}_j^\top \mathbf{v} + b_j$.

$$\mathcal{L}_{\text{CE}} = -\log \frac{e^{s_y}}{\sum_{j=1}^C e^{s_j}}.$$

Supervised Loss Functions

From Classification to Embeddings: Core Challenge



Ideal embeddings should be:

- **Intra-class compact:** same-class samples cluster tightly
- **Inter-class separable:** different classes are far apart

Limitations of Softmax

- **Objective misalignment:** only enforces correct classification, not compactness/separation explicitly
- **Open-set issue:** limited generalization to unseen identities; insufficient safety margin

Solution

Introduce **margin losses** to directly optimize embedding quality.

Supervised Loss Functions

Normalized Softmax on Unit Hypersphere



Normalization

Enforce $\|\mathbf{v}\| = 1$, $\|\mathbf{w}_j\| = 1$, and set $b_j = 0$. Then

$$s_j = s \cos \theta_j,$$

where θ_j is the angle between \mathbf{v} and \mathbf{w}_j , and $s > 0$ is a scale (inverse temperature).

$$\mathcal{L}_{\text{Norm-CE}} = -\log \frac{e^{s \cos \theta_y}}{\sum_{j=1}^C e^{s \cos \theta_j}}.$$

Geometric view

Classification equals selecting the smallest angle on the unit hypersphere, consistent with cosine scoring.

Supervised Loss Functions

SphereFace: Multiplicative Angular Margin

Form

$$\mathcal{L} = -\log \frac{e^{s \cos(m\theta_y)}}{e^{s \cos(m\theta_y)} + \sum_{j \neq y} e^{s \cos \theta_j}}, \quad m \geq 1.$$

- Decision boundary: from $\theta_1 = \theta_2$ to $m\theta_1 = \theta_2$.
- Pros: pioneering angular-margin formulation.

Caveat

Highly nonlinear; training may be unstable, often requires annealing/special schedules.

Supervised Loss Functions

CosFace (AM-Softmax): Additive Cosine Margin



Form

$$\mathcal{L} = -\log \frac{e^{s(\cos \theta_y - m)}}{e^{s(\cos \theta_y - m)} + \sum_{j \neq y} e^{s \cos \theta_j}}, \quad m > 0.$$

- Interpretation: enforce $\cos \theta_y \geq \cos \theta_j + m$ (constant margin in cosine space).
- Pros: simple and stable optimization; no annealing required.

Note

Margin is constant in cosine but not in angle; angular gap varies with θ .

Supervised Loss Functions

ArcFace (AAM-Softmax): Additive Angular Margin

Form

$$\mathcal{L} = -\log \frac{e^{s \cos(\theta_y + m)}}{e^{s \cos(\theta_y + m)} + \sum_{j \neq y} e^{s \cos \theta_j}}, \quad m > 0.$$

- Interpretation: enforce $\theta_y + m \leq \theta_j$ (constant angular gap).
- Pros: clear geometry; strong open-set performance.

Implementation hints

Ensure numerical stability (clip $\cos \theta$ to $[-1, 1]$; avoid unstable inverse-trig operations).

Scope of Numerical Stabilization Tricks

Why Emphasize ArcFace?

Key Conclusion

Numerical stabilization (clipping $\cos \theta$, avoiding \arccos) is not unique to ArcFace, but only ArcFace must address it due to angle-dependent computation; other losses (SphereFace/CosFace) are less sensitive or naturally avoid the instability.

Sensitivity to Numerical Anomalies

Loss	Require \arccos to recover θ ?	Sensitive to out-of-range $\cos \theta$?
CosFace	No (use $\cos \theta$ directly)	Low
SphereFace	No (use multiple-angle identities)	Medium
ArcFace	Yes (must recover θ)	High

- **ArcFace pain point:** \arccos is mandatory; out-of-range $\cos \theta$ leads to NaNs.
- **Others:** \arccos is optional; anomalies affect accuracy but not training stability.

Supervised Loss Functions

Three Margins at a Glance



Method	Margin domain	Decision boundary (1 vs 2)
SphereFace	Angle (multiplicative)	$\cos(m\theta_1) = \cos(\theta_2)$
CosFace	Cosine (additive)	$\cos(\theta_1) - m = \cos(\theta_2)$
ArcFace	Angle (additive)	$\cos(\theta_1 + m) = \cos(\theta_2)$

Constancy

CosFace: constant in cosine.

ArcFace: constant in angle.

SphereFace: depends on θ .

Practice

Prefer ArcFace or CosFace for stability and performance.

Optional: Add Center Loss

Form

$$\mathcal{L}_{\text{center}} = \frac{1}{2} \sum_i \|\mathbf{v}_i - \mathbf{c}_{y_i}\|_2^2.$$

- Complements margin softmax by explicitly reducing within-speaker variance.
- Combine as $\mathcal{L} = \mathcal{L}_{\text{ArcFace/CosFace}} + \lambda \mathcal{L}_{\text{center}}$ with small λ .

Supervised Loss Functions

More Complex: Adaptive Angular Margin

Form

$$\mathcal{L} = -\log \frac{e^{s \cos(\theta_y + m_y)}}{e^{s \cos(\theta_y + m_y)} + \sum_{j \neq y} e^{s \cos(\theta_j - m_j)}}$$

where $m_y, m_j > 0$ are class-adaptive margins.

- Idea: assign different margins to different classes to adapt to data properties
- Adjustment strategies:
 - Larger margin for classes with fewer samples
 - Larger margin for classes with higher intra-class dispersion
 - Larger margin for classes with higher fraction of hard examples
- Pros: handles class imbalance; more precise optimization for difficult classes

Implementation

Requires additional class statistics; margin policies should be designed carefully.

- ① Speaker Modeling: Background, Applications, and Trends
- ② Discriminative Speaker Representation Learning
- ③ Self-supervised Speaker Representation Learning
- ④ Multi-modal Speaker Representation Learning
- ⑤ Robustness and Interpretability
- ⑥ Speaker Modeling in Related Tasks
- ⑦ Practice: WeSpeaker and WeSep

- Leverage large pretrained models
 - Self-supervised pretrained speech models
 - ASR model initialization
 - Efficient fine-tuning
- Self-supervised learning methods
 - SimCLR/MoCo/DINO
 - Stage-wise iterative training

Self-supervised Speaker Representation Learning

Fine-tuning Methods

- 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".

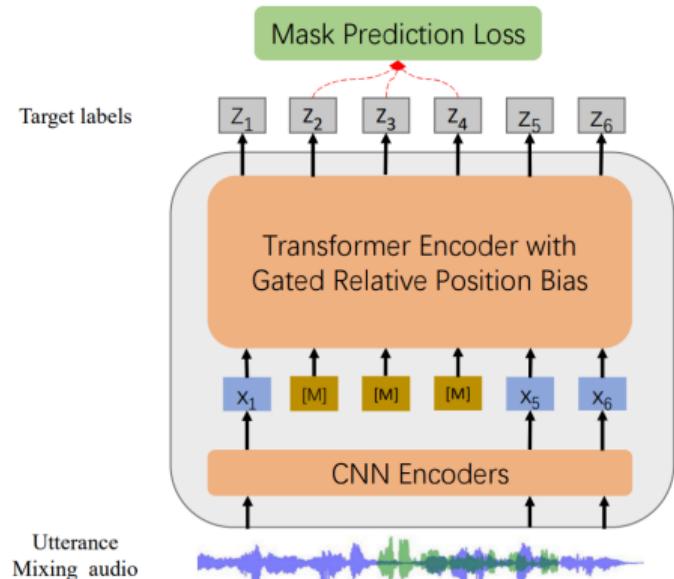
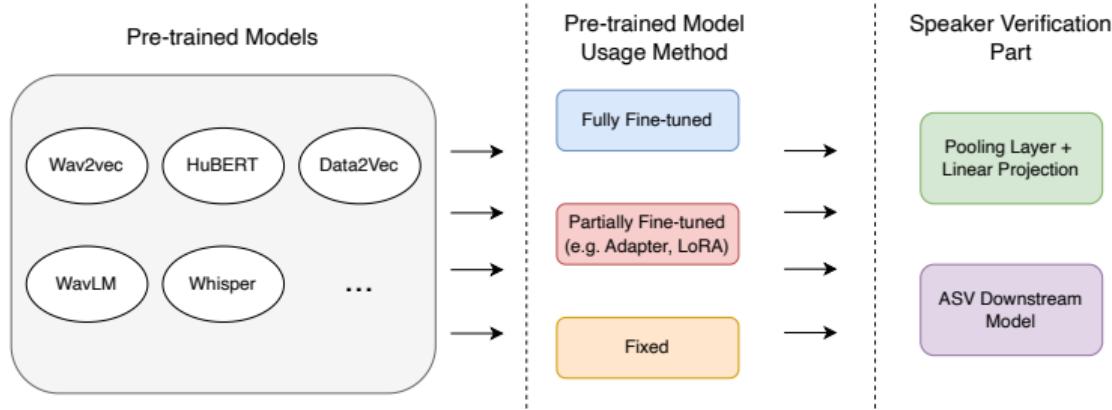


Figure: WavLM architecture

Leveraging Pretrained Models

Integration Strategies of Large Pretrained Speech Models



Integration Strategies:

- 1. Feature Extraction:** use pretrained features as input
- 2. Fine-tuning:** adapt pretrained models to speaker tasks
- 3. Multi-task Learning:** jointly optimize multiple tasks

Leveraging Pretrained Models

Fine-tuning Methods

Fine-tune **SSL speech models** on speaker verification^a

- Replace Fbank with pretrained representations
- Learnable weighted aggregation across layers

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

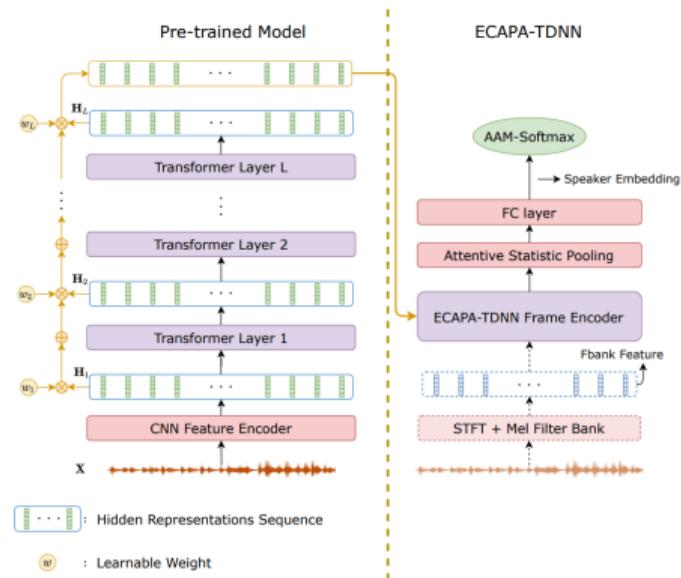


Figure: Using pretrained representations

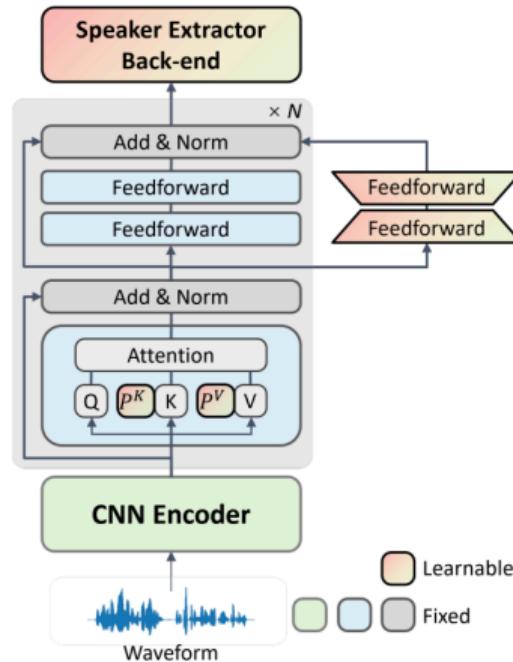
Leveraging Pretrained Models

Fine-tuning Methods

Efficient fine-tuning for SSL models with adapters on speaker verification^a

- Freeze large pretrained model
- Use adapters for parameter-efficient tuning

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



Leveraging Pretrained Models

Fine-tuning Methods

Adopt large SSL model w2v-BERT 2.0 for SOTA performance²

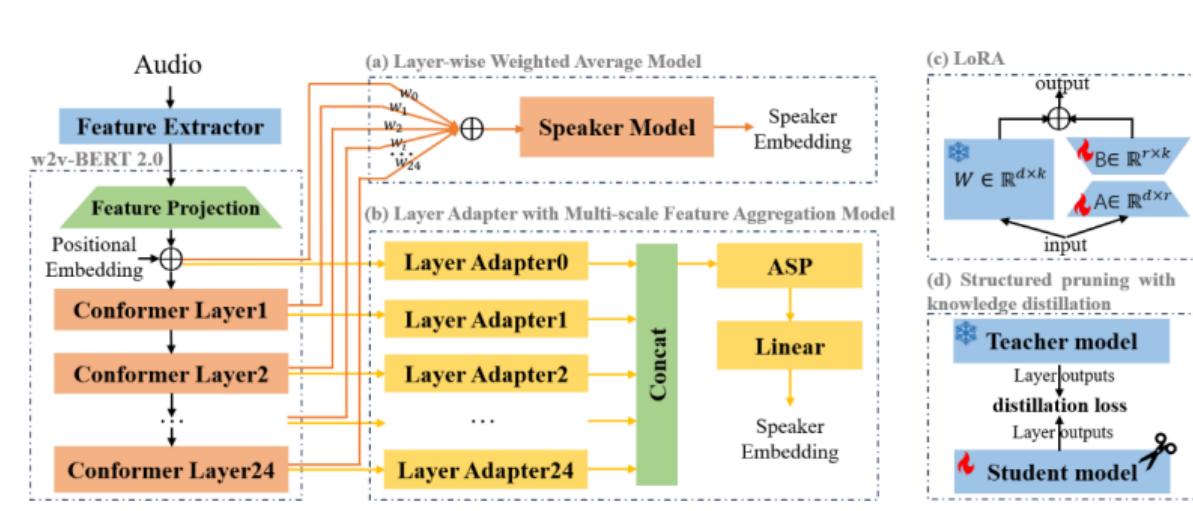


Figure: Enhancing SV with w2v-BERT 2.0 and distillation-guided structured pruning

²Li, Cheng, and Li, "Enhancing Speaker Verification with w2v-BERT 2.0 and Knowledge Distillation guided Structured Pruning".

Leveraging Pretrained Models

Fine-tuning Methods

Fine-tune **ASR models** on speaker verification^{ab}

- Pretrain with ASR datasets
- Initialize for speaker tasks

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

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

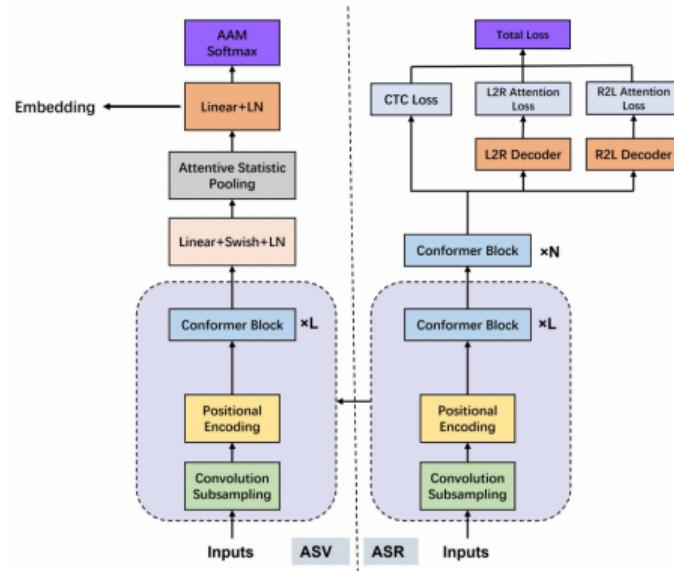


Figure: ASR transfer illustration

Self-supervised Speaker Representation Learning



Metric-based Loss Functions

Metric learning losses provide contrastive supervision, e.g., Triplet, Prototypical, GE2E³, and Angular Prototypical⁴.

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

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

³Wan et al., "Generalized end-to-end loss for speaker verification".

⁴Chung et al., "In defence of metric learning for speaker recognition".

Self-supervised Speaker Representation Learning

Assumption for SSL on SV

Assumption for SSL on speaker verification^a

- Segments from the same utterance belong to the 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: Illustration of the assumption

Self-supervised Speaker Representation Learning

SimCLR

Based on SimCLR^a, adapted to speaker tasks^b

- Crop two segments from an utterance to construct positive/negative pairs
- Use metric loss to attract positives and repel negatives

^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".

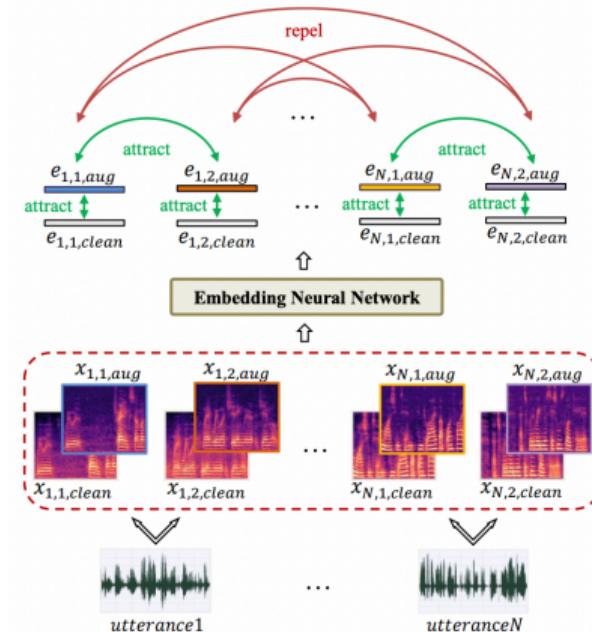


Figure: SimCLR for speaker task

Based on DINO^a, adapted to speaker tasks^{bc}

- Crop several segments from one utterance and build only positive pairs
- Use cross-entropy to attract 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".

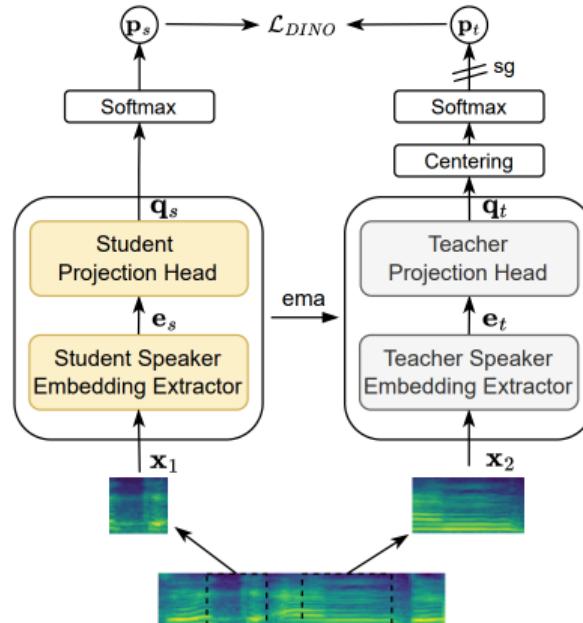


Figure: DINO for speaker task

Self-supervised Speaker Representation Learning

Stage-wise Iterative Training

Two-stage 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".

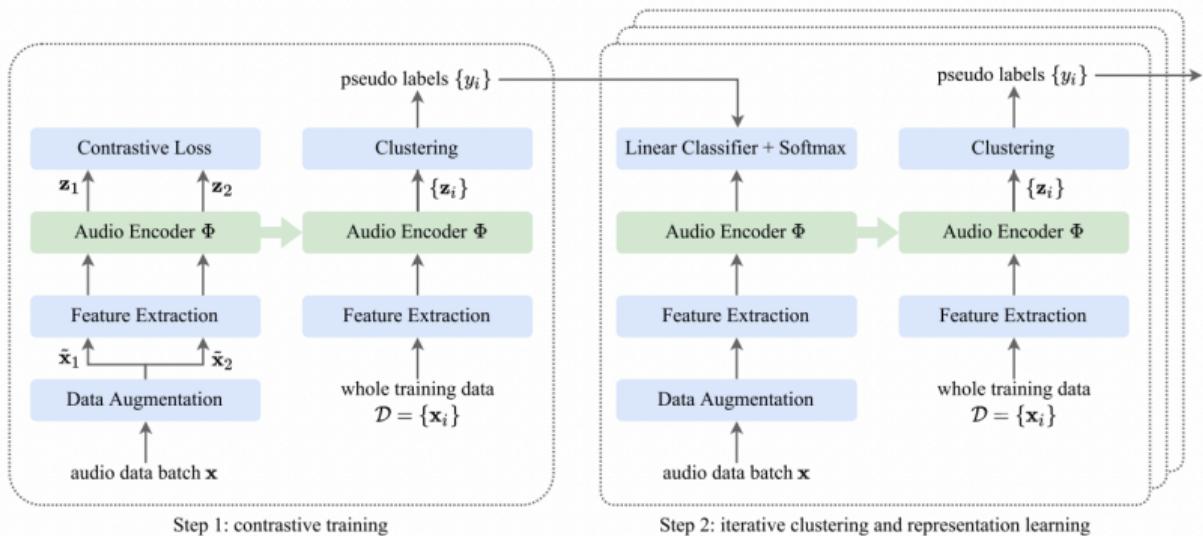
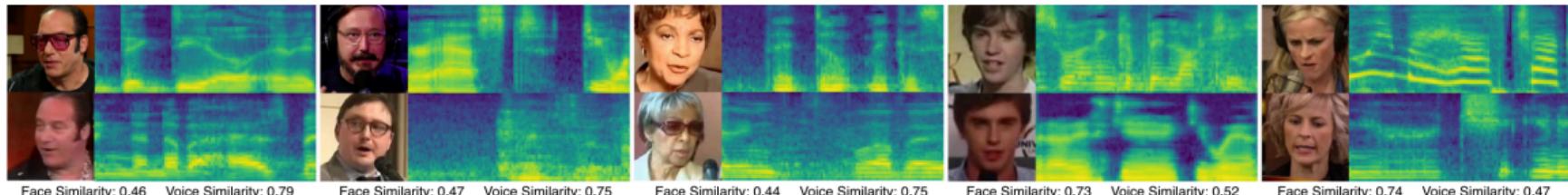


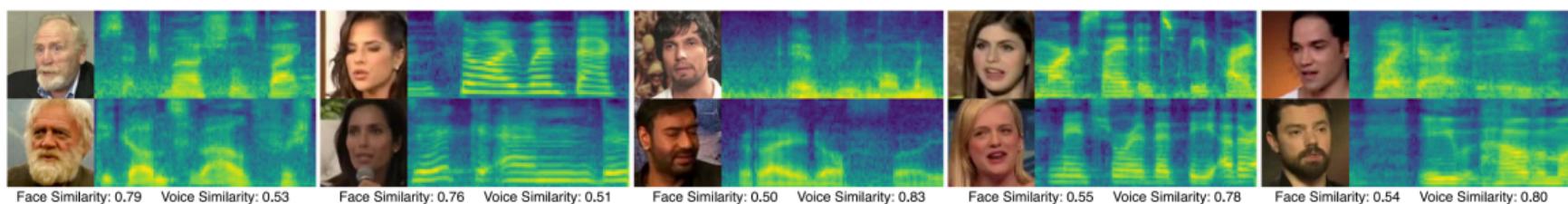
Figure: Iterative framework for SSL speaker verification

- ① Speaker Modeling: Background, Applications, and Trends
- ② Discriminative Speaker Representation Learning
- ③ Self-supervised Speaker Representation Learning
- ④ Multi-modal Speaker Representation Learning
- ⑤ Robustness and Interpretability
- ⑥ Speaker Modeling in Related Tasks
- ⑦ Practice: WeSpeaker and WeSep

Complementarity between Audio and Visual Modalities



(a)

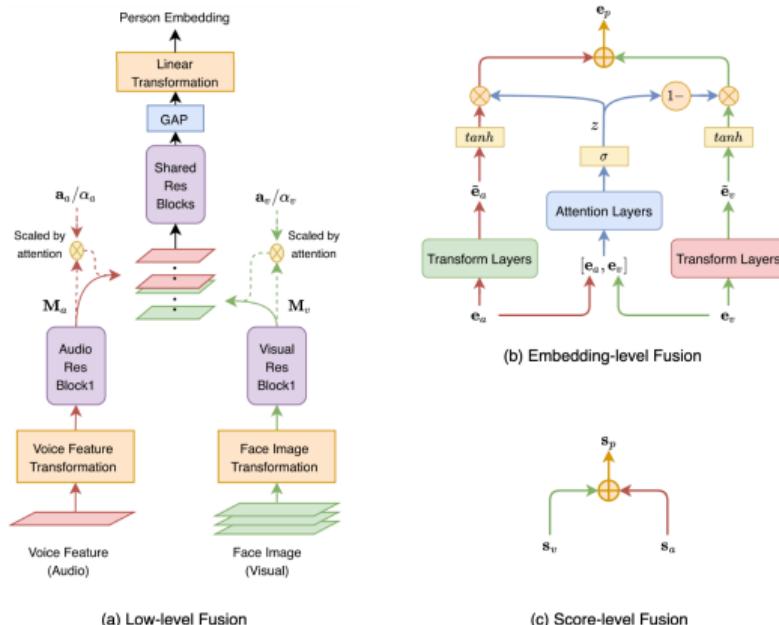


(b)

Figure: Speaker similarity from audio or visual information⁵

- Upper part shows similarity to the same person
- Lower part shows similarity between different persons

⁵Qian, Chen, and Wang, "Audio-visual deep neural network for robust person verification".



- Embedding-level fusion outperforms low-level fusion
- Attention in embedding-level fusion is more noise-robust than score-level fusion

Figure: Different levels of audio-visual fusion^a

Multi-modal Knowledge Distillation

From Audio-Visual Systems to Single-modality Systems

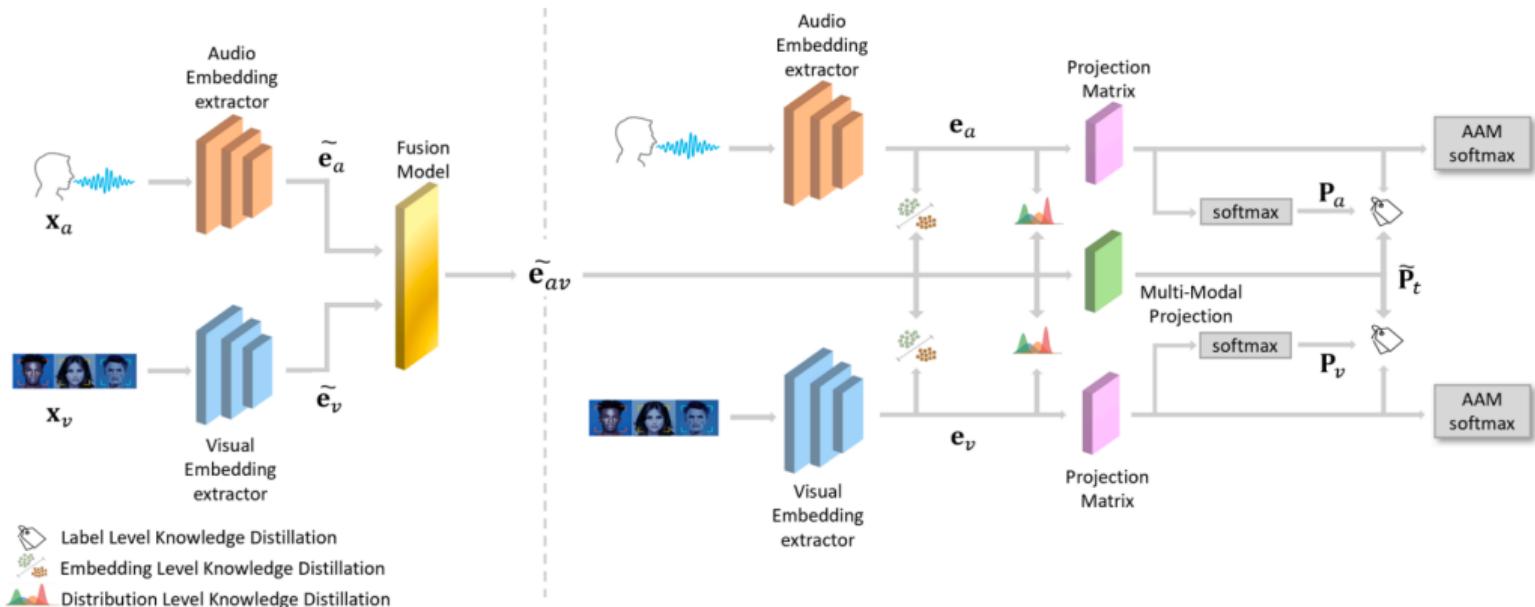


Figure: Knowledge distillation from audio-visual to single-modality systems⁶

⁶Zhang, Chen, and Qian, "Knowledge Distillation from Multi-Modality to Single-Modality for Person Verification".

Multi-modal Knowledge Distillation

From Visual Systems to Audio Systems

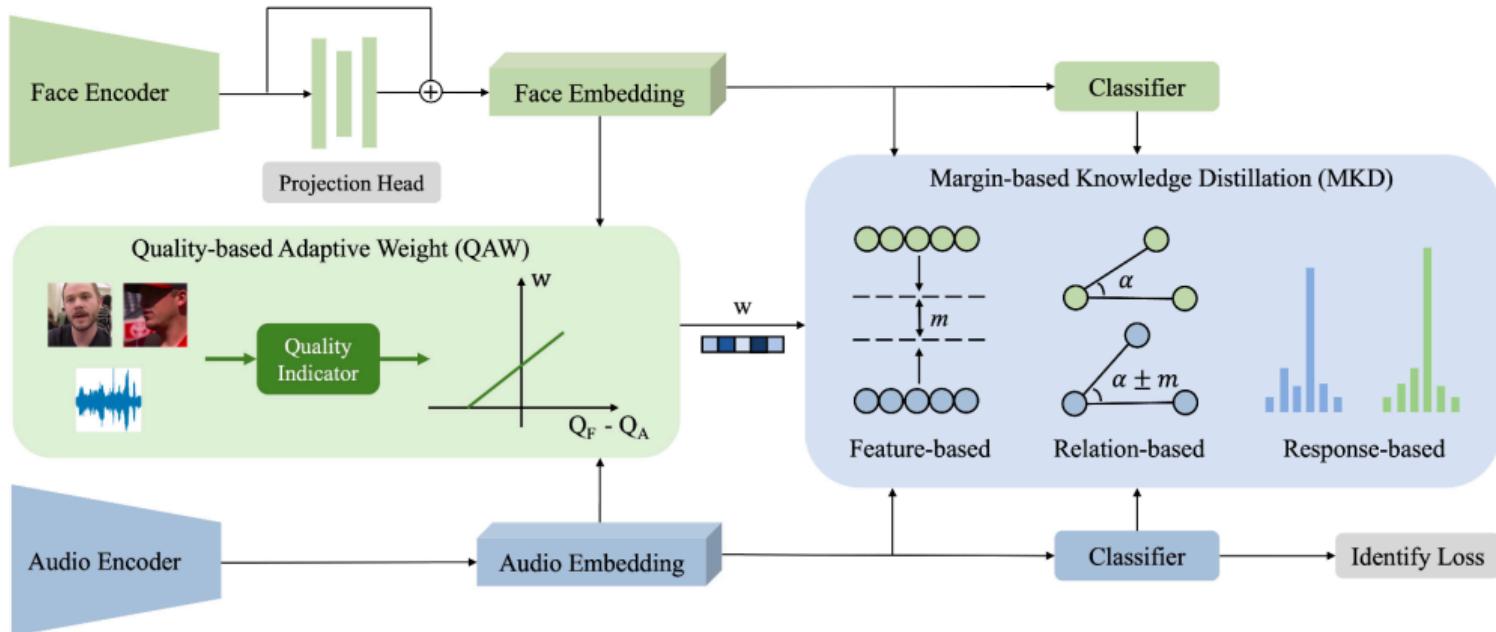
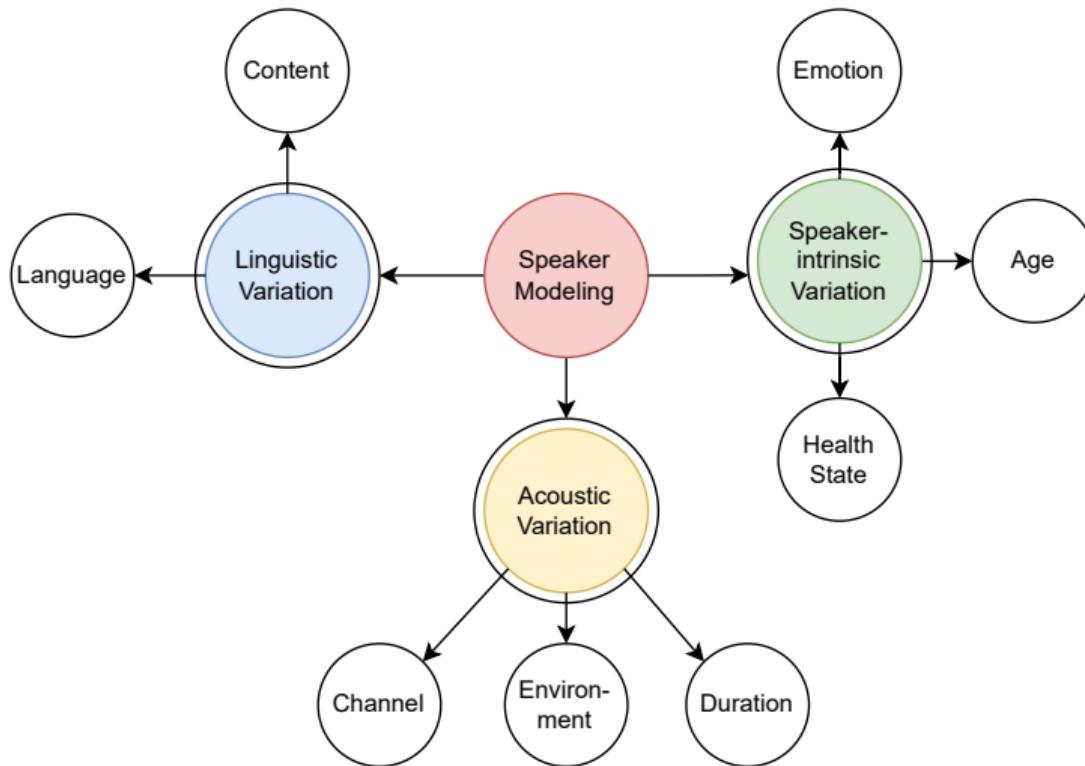


Figure: Knowledge distillation from visual to audio systems⁷

- ① Speaker Modeling: Background, Applications, and Trends
- ② Discriminative Speaker Representation Learning
- ③ Self-supervised Speaker Representation Learning
- ④ Multi-modal Speaker Representation Learning
- ⑤ Robustness and Interpretability
- ⑥ Speaker Modeling in Related Tasks
- ⑦ Practice: WeSpeaker and WeSep



Recording environments introduce variability to speaker identity modeling due to device and microphone distance, etc. Various domain adaptation methods are adopted in speaker recognition to enhance robustness across devices, including

- Discrepancy-based alignment
- Adversarial learning
- Domain-specific adapters

Discrepancy-based alignment aims to minimize domain differences in the latent space and promote domain-invariant representations. Proper divergence measures are central to these methods, e.g., MMD⁸, CORAL⁹, etc.

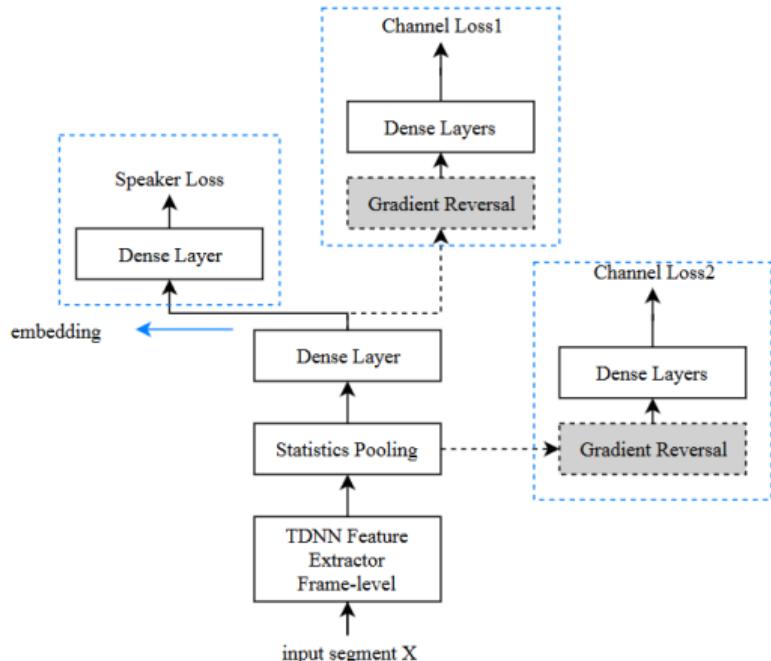
$$\mathcal{L}_{\text{mmd}} \triangleq \sup_{\phi \in \Phi} (\mathbb{E}_S [\phi(S)] - \mathbb{E}_T [\phi(T)]) \quad (2)$$

⁸Li, Han, and Song, "CDMA: Cross-Domain Distance Metric Adaptation for Speaker Verification".

⁹Li, Zhang, and Chen, "The coral++ algorithm for unsupervised domain adaptation of speaker recognition".

Model Robustness to Devices

Adversarial Learning



Adversarial learning uses a domain classifier to remove discriminative domain information from features. The min-max optimization in domain-adversarial training reduces domain gaps and enforces domain-invariant feature extraction^a.

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

Figure: Channel-level adversarial learning^a

Model Robustness to Devices

Domain-specific Adapters

Instead of directly aligning domains with discrepancy measures, adding domain-specific adapters and related modules helps capture and mitigate domain variance, yielding domain-invariant embeddings.

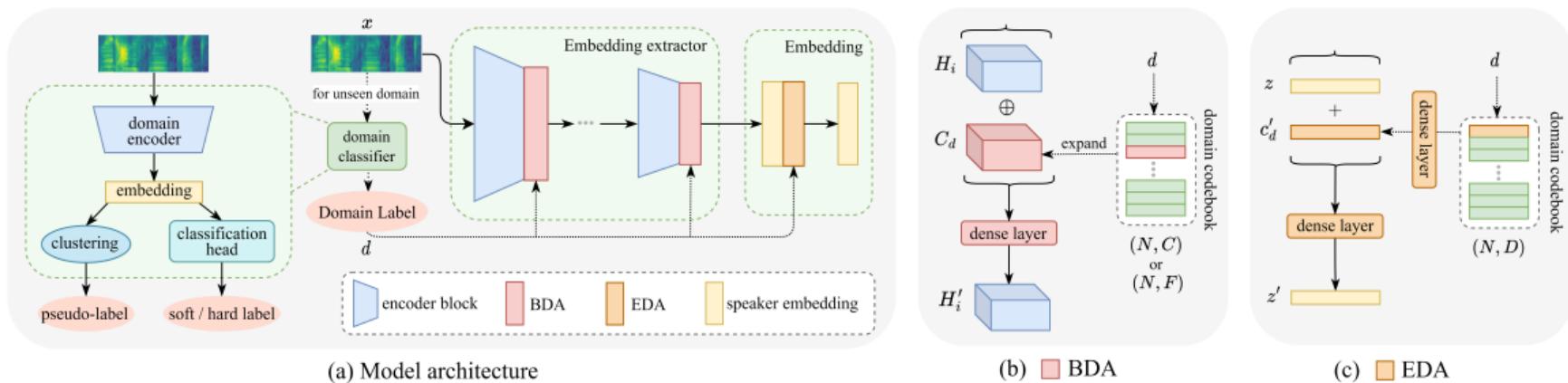


Figure: Framework with domain-specific adapters¹⁰

¹⁰Huang et al., "Enhancing Cross-Domain Speaker Verification through Multi-Level Domain Adapters".

Model Robustness to Language Mismatch

Language Mismatch across Datasets

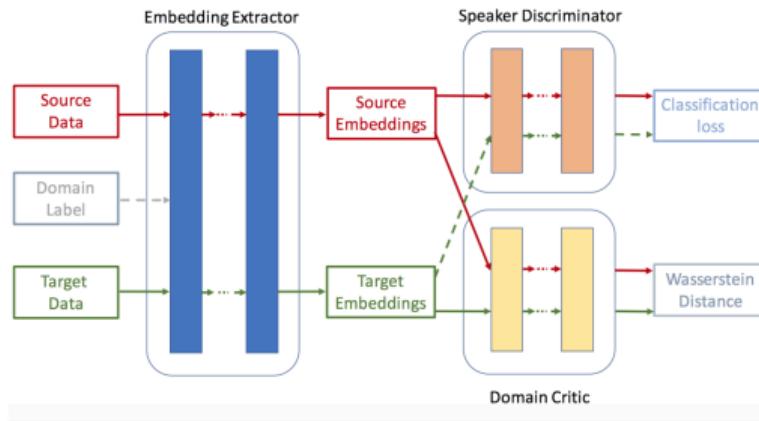


Observation: In real scenarios, SV systems trained in one language may degrade when tested in another language.

Model Robustness to Language Mismatch

Mismatch between Enrollment/Test

Over 40% of the world's population is bilingual. Mismatch happens when enrollment and test use different languages.



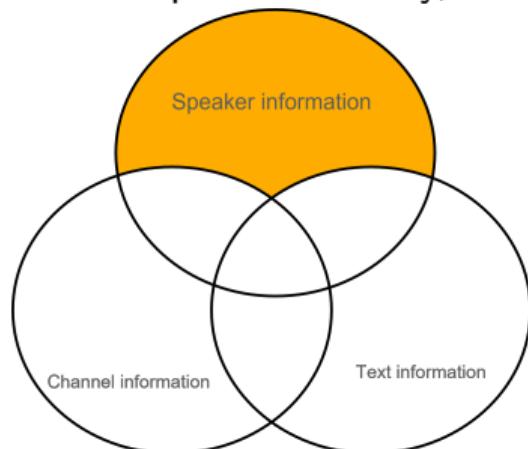
Adversarial learning uses a language classifier to remove discriminative language information from features. The min-max optimization minimizes language gaps and enforces language-invariant feature extraction^{a,b}.

^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".

Figure: Adversarial learning for language mismatch

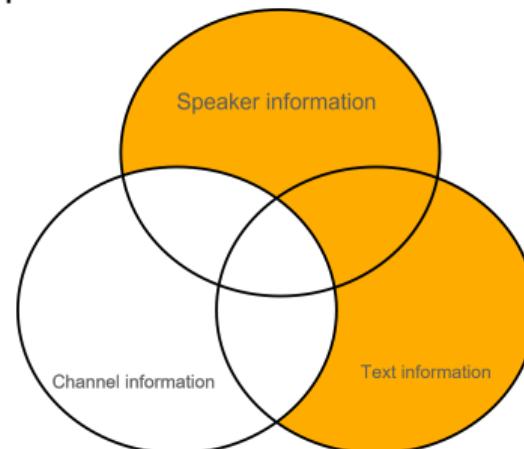
In addition to speaker identity, text/content is key to speech communication.



For

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

Enrollment: Hey Siri; Test: arbitrary phrase



For

text-dependent tasks, we also need **content information**.

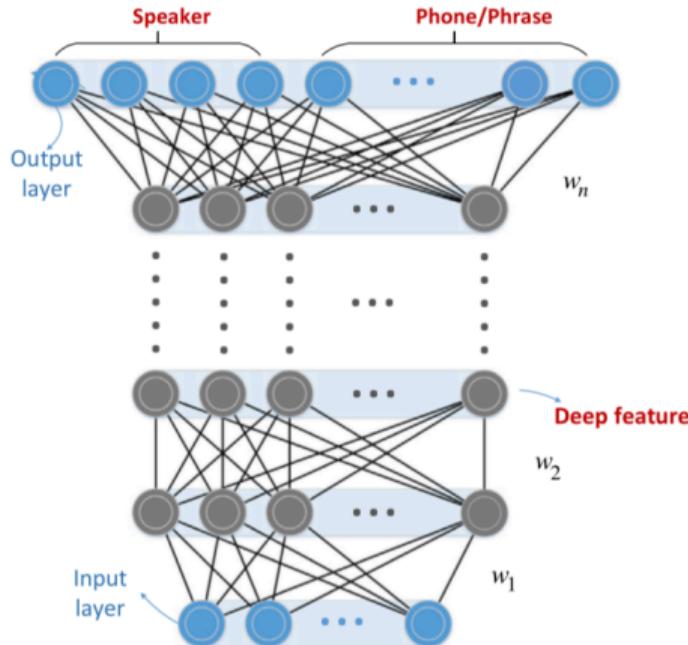
Enrollment: Hey Siri; Test: Hey Siri

Content Representations

- Phone indices
- ASR-predicted phone posteriors
- Hidden representations of ASR models
- Phrase IDs (fixed-phrase datasets)
- Normalized phone distribution

Text Robustness

Multi-task Learning in d-vector Framework¹¹



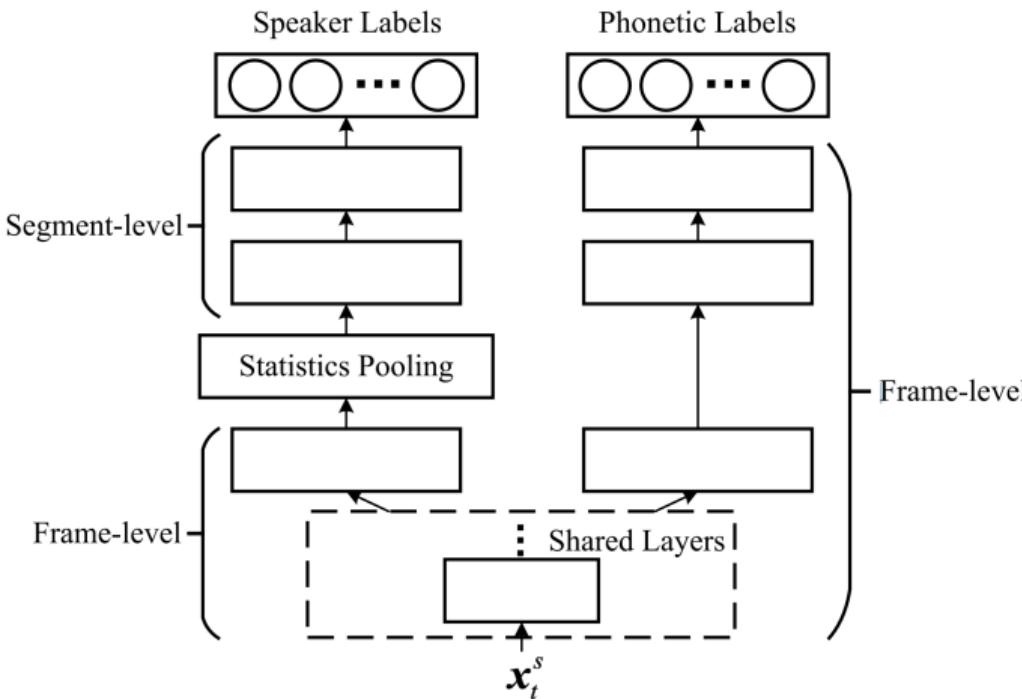
- **Text-dependent task**
- Multi-task at **frame-level**
- Performance improvement

Explicitly modeling phonetic information is intuitive for text-dependent SV.

¹¹ Liu, Yuan, et al. "Deep feature for text-dependent speaker verification." Speech Communication 73 (2015): 1-13.

Text Robustness

Multi-task Learning in x-vector Framework¹²

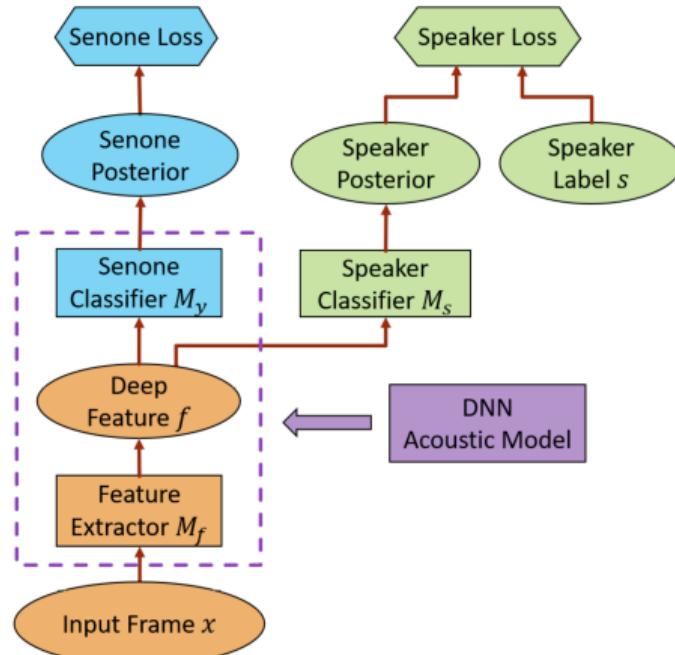


- **Text-independent task**
- Multi-task at **frame-level**
- Performance improvement

¹² Liu, Yi, et al. "Speaker Embedding Extraction with Phonetic Information." Proc. Interspeech 2018 (2018): 2247-2251.

Text Robustness

Speaker-invariant Training for ASR¹³

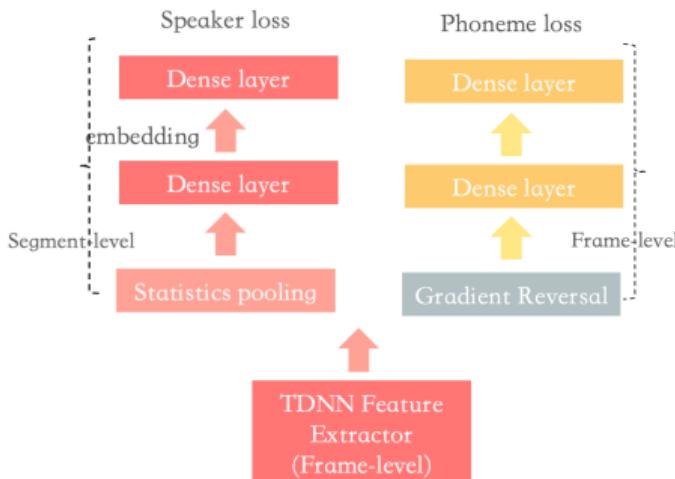
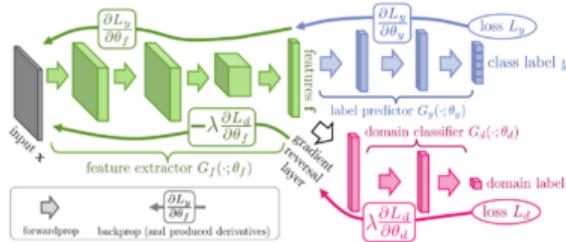


- Acoustic modeling
- Adversarial training to suppress speaker effect
- Performance improvement

¹³ Meng, Zhong, et al. "Speaker-invariant training via adversarial learning." ICASSP 2018.

Text Robustness

Frame-level Multi-task/Adversarial Training



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

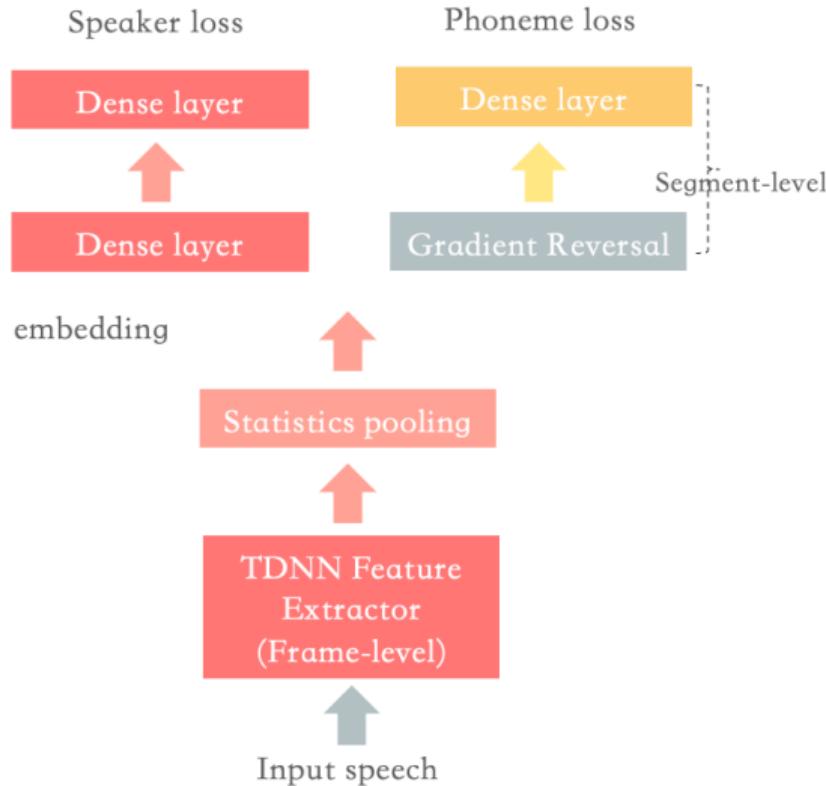
$$\mathcal{L}_p = \frac{1}{N} \sum_{i=1}^N \text{CE}(M_p(M_f(\mathbf{x}_i)), \mathbf{y}_i^p)$$

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

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

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 segment \mathbf{x} with N frames, the segment-level phonetic label \mathbf{y}^p is

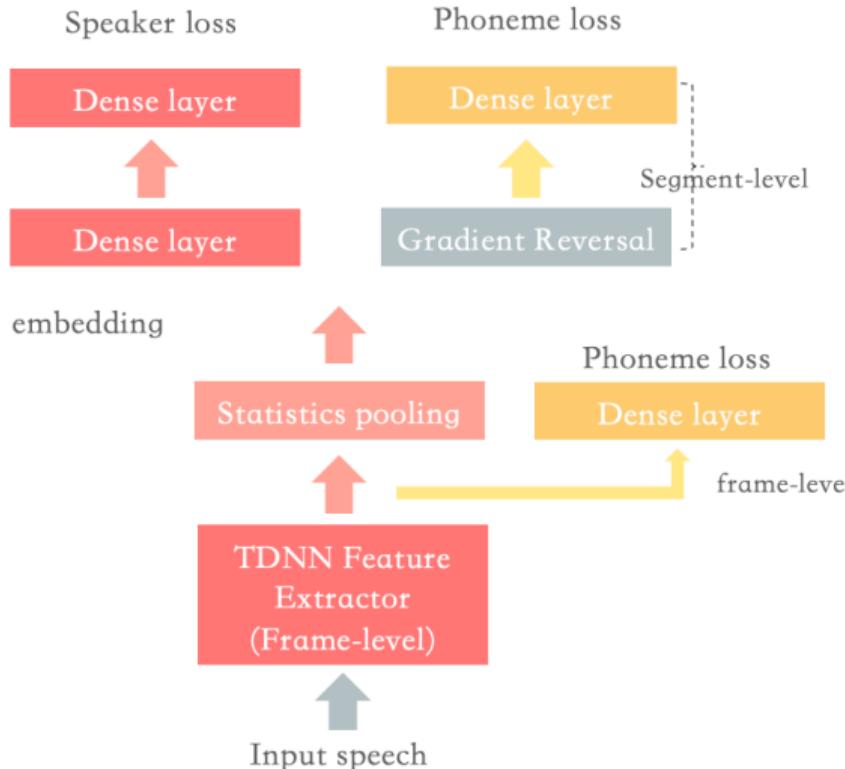
$$\mathbf{y}^p = \{y_1, y_2, \dots, y_C\}$$

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

where C is the size of the selected phone set and N_c is the number of occurrences of phone c in \mathbf{x} .

Text Robustness

Frame-level Multi-task + Segment-level Adversarial



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

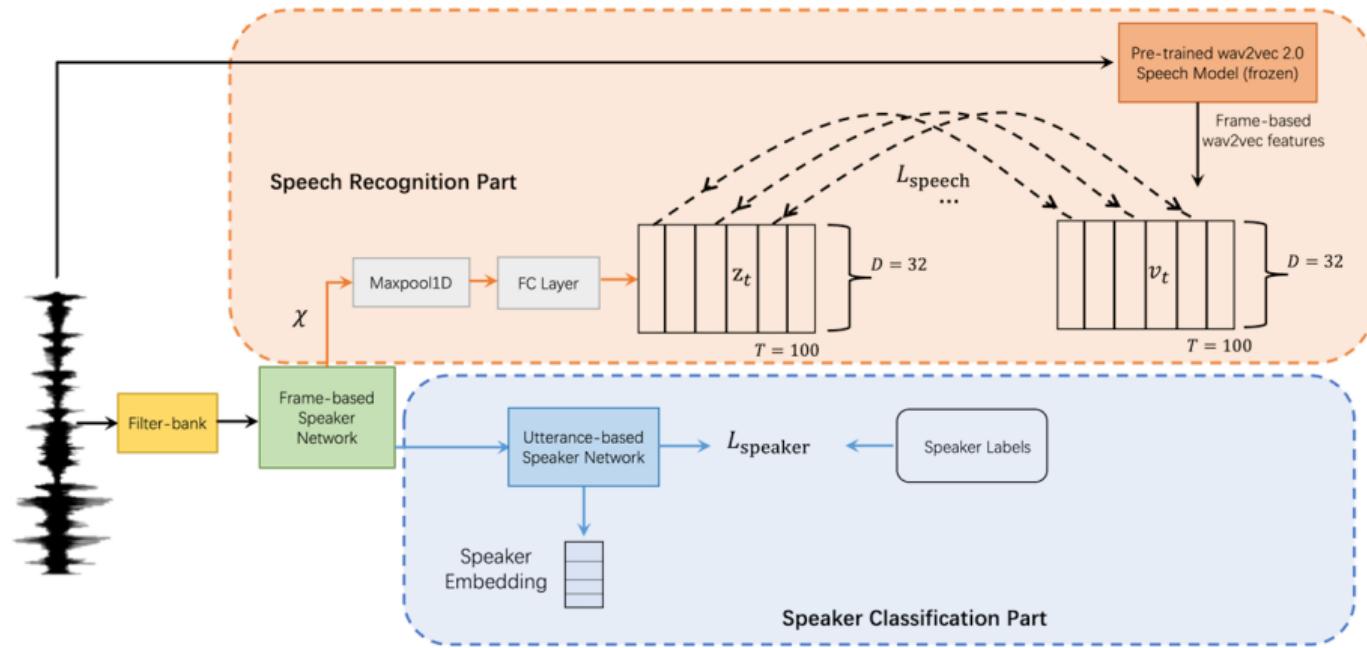
Figure: Segment-level multi-task/adversarial

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

Figure: Frame-level multi-task + segment-level adversarial

Text Robustness

Multi-task Training with Advanced Content Representations¹⁴



¹⁴Jin, Tu, and Mak, "Phonetic-aware speaker embedding for far-field speaker verification".

Text Robustness

Phonetic-aware Speaker Embedding Learning

Extract PBN from pretrained ASR and combine with filterbanks¹⁵

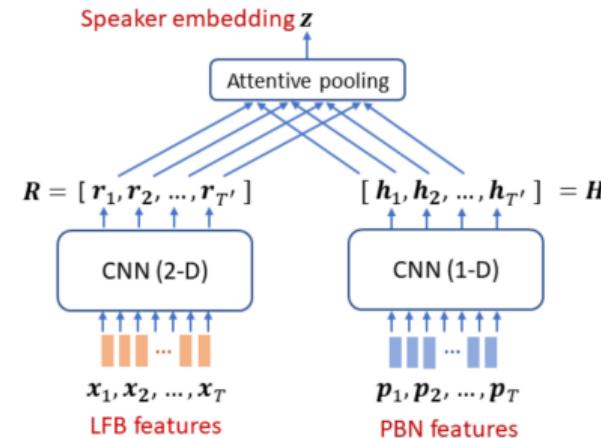
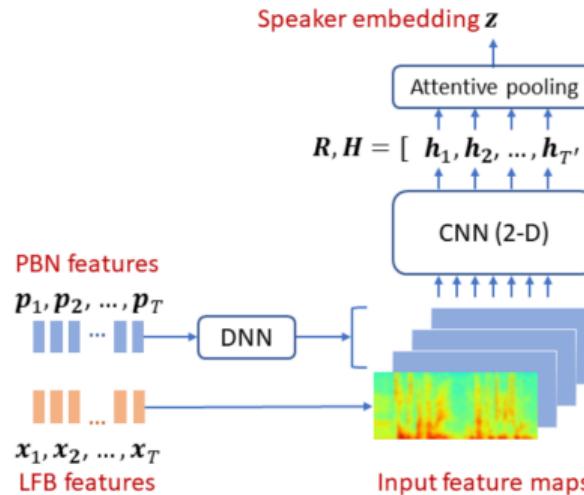


Fig. 1: Implicit phonetic attention by combining LFB and PBN features at the input layer (LFB: log filter bank; PBN: phonetic bottleneck).

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

¹⁵Zhou T, Zhao Y, Li J, et al. CNN with phonetic attention for text-independent speaker verification, *ASRU 2019*



Text Robustness

Phonetic-aware Speaker Embedding Learning

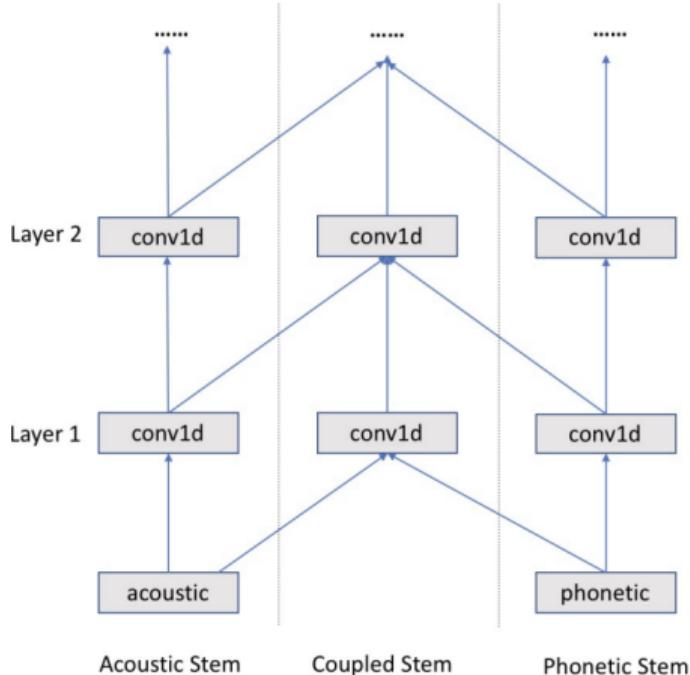


Table 1: *Network configurations of PacNet*

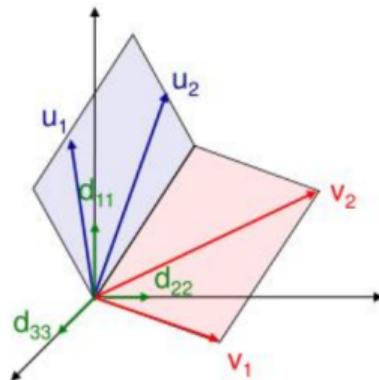
Layer 7	Linear	In=1024 Out=1000		
Layer 6	Pooling	In=1024 Out=1024		
Layer 5	Conv1d	In=2048 Out=1024		
Layer 4	Conv1d kernel=5	Out=512 In=512	Out=1024 In=2048	Out=512 In=512
Layer 3	Conv1d kernel=5	Out=512 In=512	Out=1024 In=2048	Out=512 In=512
Layer 2	Conv1d kernel=5	Out=512 In=512	Out=1024 In=2048	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

- Use Triplet loss instead of softmax

Back to the Past

Joint factor analysis for speaker representations¹⁷

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

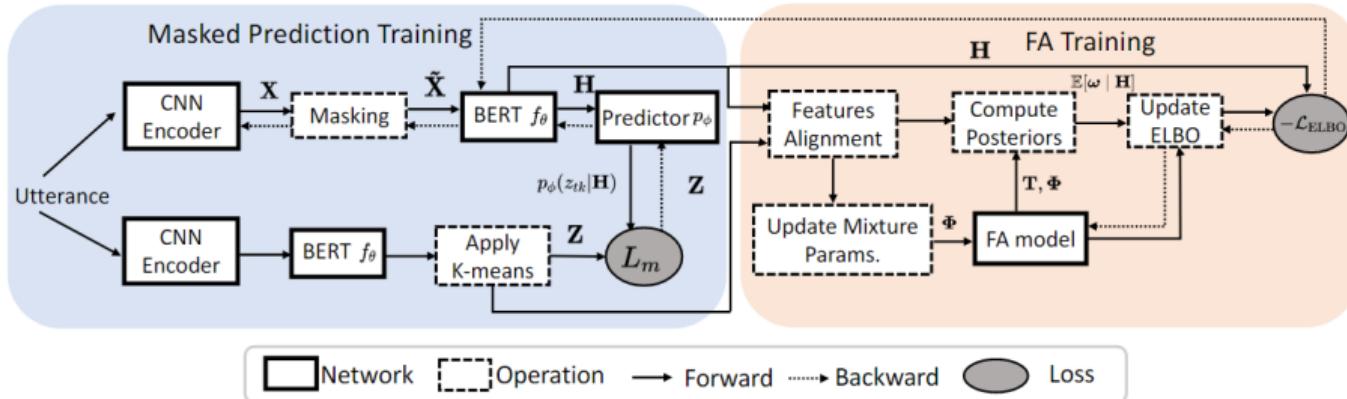


- Gaussian priors over factors \mathbf{y} , \mathbf{z} , \mathbf{x}
- Estimate \mathbf{M}^{UBM} , \mathbf{V} , \mathbf{D} , \mathbf{U} with EM
- \mathbf{V} captures primary speaker variability (speaker factors)
- \mathbf{D} captures channel variability
- \mathbf{U} captures residual variability

¹⁷Kenny, Patrick, et al. "Joint factor analysis versus eigenchannels in speaker recognition." TASLP 2007

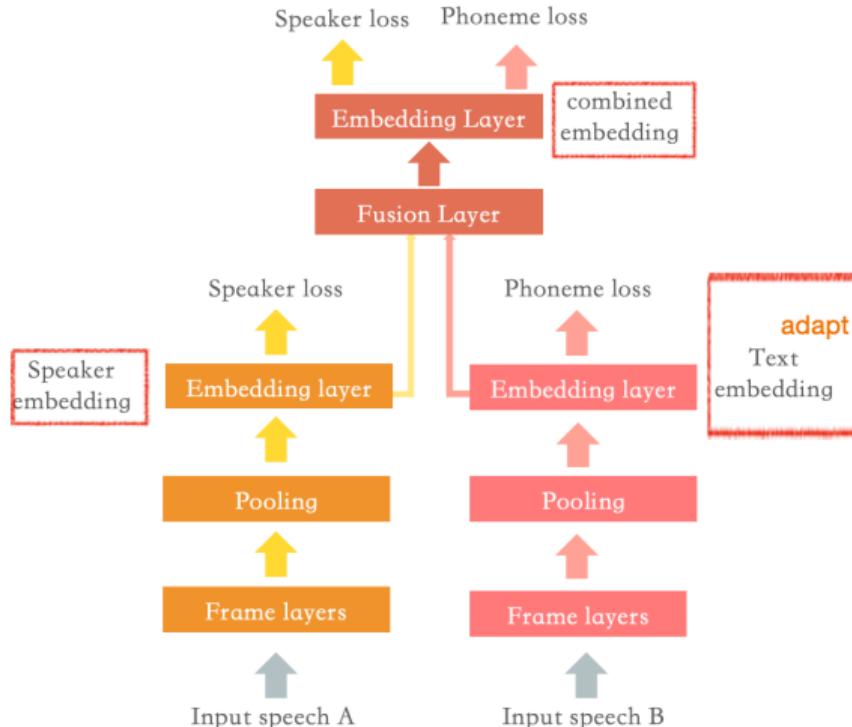
Disentangling Speech Representations

Neural Factor Analysis: ignore phonetic variability via extra alignment 18



Disentangling Speech Representations

Factorizing and Re-composing Phonetic Information¹⁹

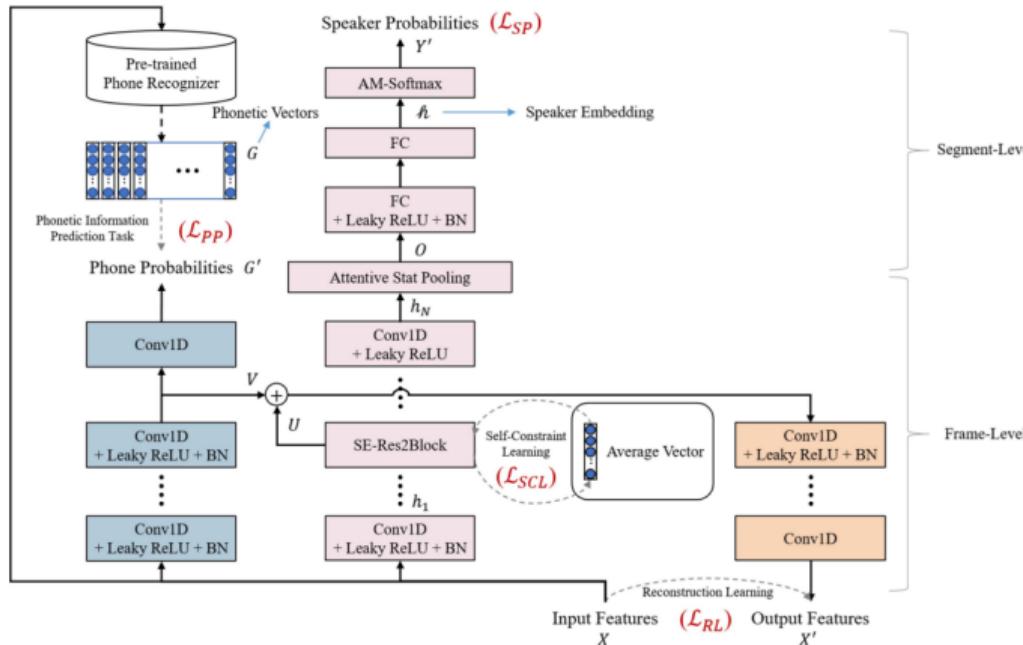


- Segment-level reconstruction
- Factorize speaker and text information
- For text-independent tasks, ignore text information
- For text-dependent tasks, use combined embeddings
- **Text adaptation:** modify text-related info in embeddings while preserving speaker identity (change enrolled keyword)

¹⁹Yang Y*, Wang S*, Gong X, et al. Text adaptation for speaker verification with speaker-text factorized embeddings. ICASSP 2020

Disentangling Speech Representations

Factorizing and Re-composing Phonetic Information²⁰



- Frame-level reconstruction
- Mean-center frame-level speaker reps
- Coarse phonetic classes (vowel, semivowel, 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.

²⁰Hong Q B, Wu C H, Wang H M. Decomposition and Reorganization of Phonetic Information for Speaker Embedding Learning. TASLP 2023

Disentangling Speech Representations

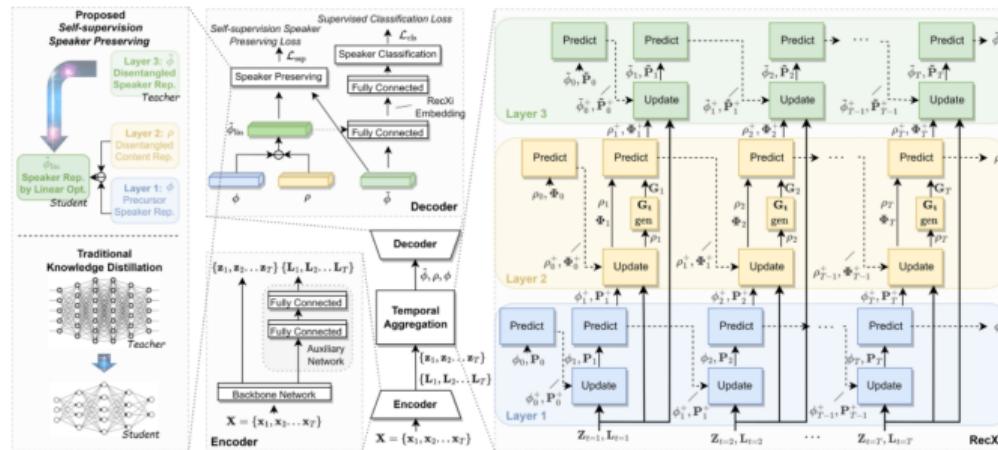
RecXi with Multi-Gaussian Inference²¹

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



by three Gaussian
inference layers

and a novel speaker
preserving self-supervision



²¹Liu T C, Disentangling Voice and Content with Self-Supervision for Speaker Recognition

Disentangling Speech Representations

Tokenizer-based Approach: SpeechTokenizer



Ensure the first-layer representation contains content-related information; residual layers naturally fill the remaining details—specifically, paralinguistic information.²²

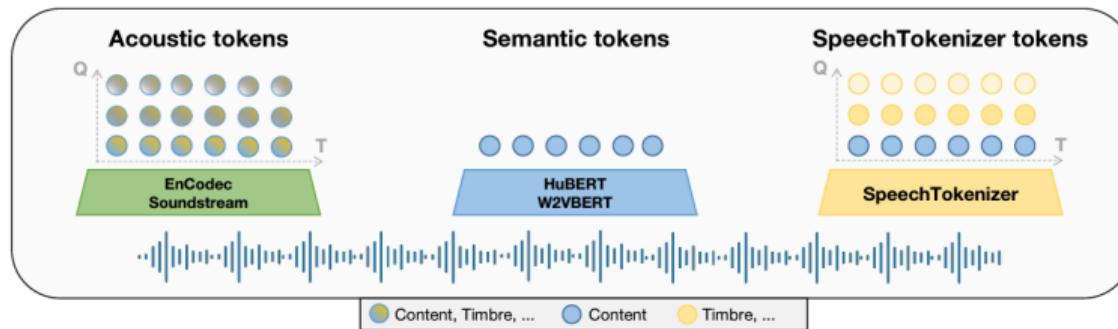


Figure 1: Illustration of information composition of different discrete speech representations. Speech tokens are represented as colored circles and different colors represent different information.

²²Zhang X, Zhang D, Li S, et al. SpeechTokenizer: Unified Speech Tokenizer for Speech Large Language Models, 2023.

Disentangling Speech Representations

Tokenizer-based Approach: SpeechTokenizer

Semantic distillation for disentanglement

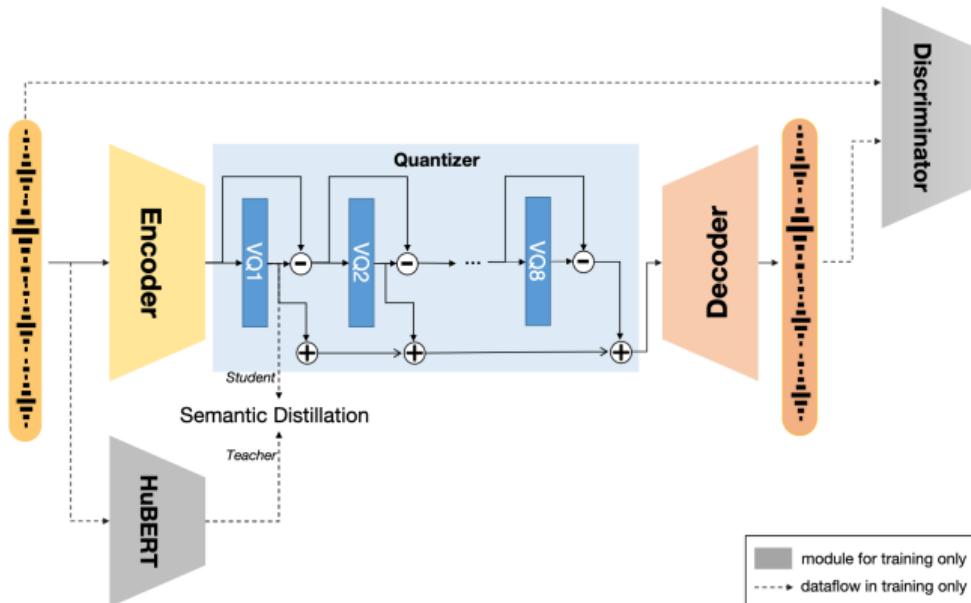


Figure 2: Illustration of SpeechTokenizer framework.

- **Continuous distillation** HuBERT layer-9 outputs / layer-mean

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

- **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



Assumption: If an attribute is encoded in speaker representations, a classifier predicting that attribute will achieve accuracy depending on how well it is embedded.^{23, 24, 25, 26}

- Speaker-related attributes: identity, gender, speaking rate.
- Text-related factors: spoken words, word order, utterance length.
- Channel-related elements: phone ID, noise type.

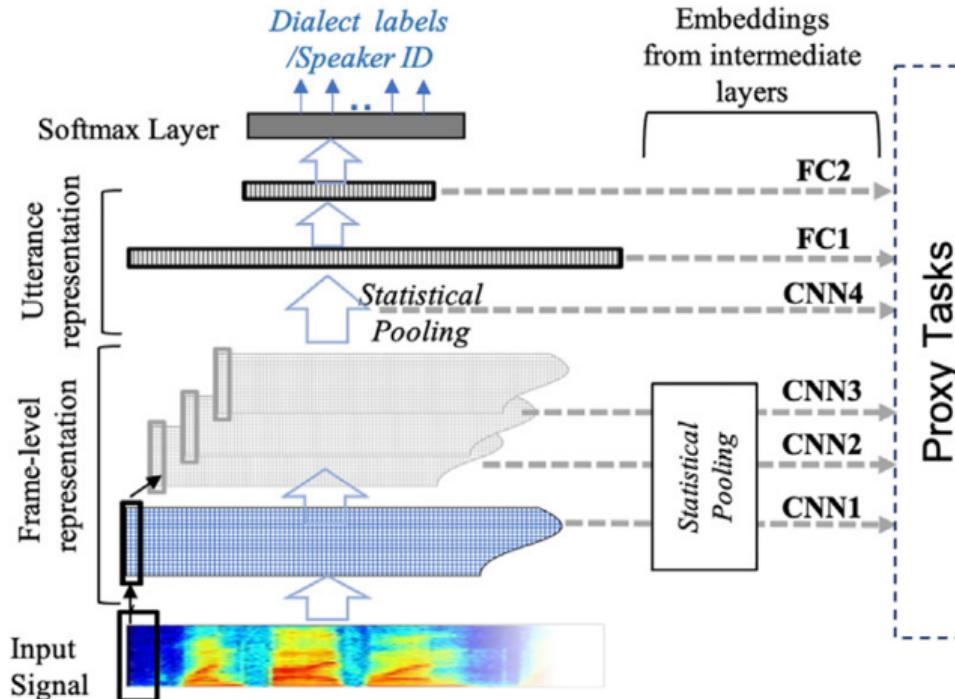
²³Wang, Qian, and Yu, "What does the speaker embedding encode?"

²⁴Belinkov and Glass, "Analyzing hidden representations in end-to-end automatic speech recognition systems".

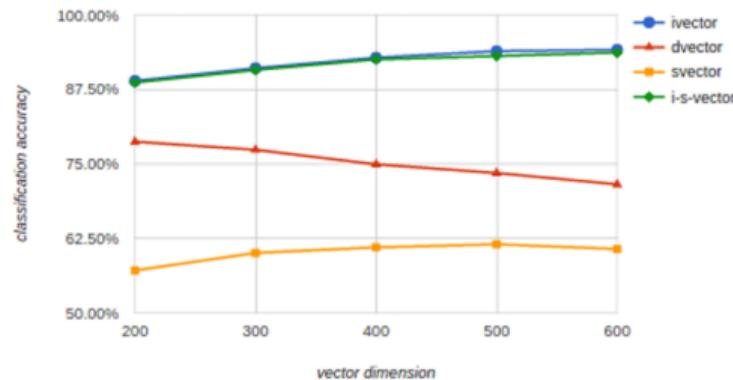
²⁵Raj et al., "Probing the information encoded in x-vectors".

²⁶Zhao et al., "Probing Deep Speaker Embeddings for Speaker-related Tasks".

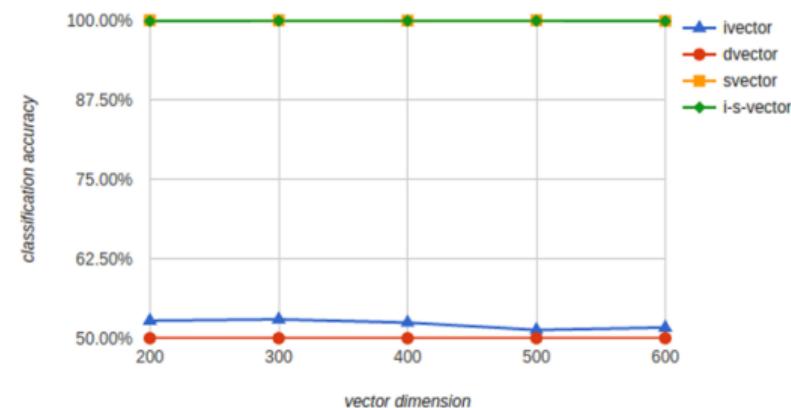
Paradigm: probe pretrained embeddings with proxy tasks²⁷



Examples of speaker identity and word order tasks²⁸



Speaker identity task



Word order task

²⁸Wang, Qian, and Yu, "What does the speaker embedding encode?"

In speaker modeling, f is the speaker classifier, c denotes the class, and θ denotes trainable parameters.

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

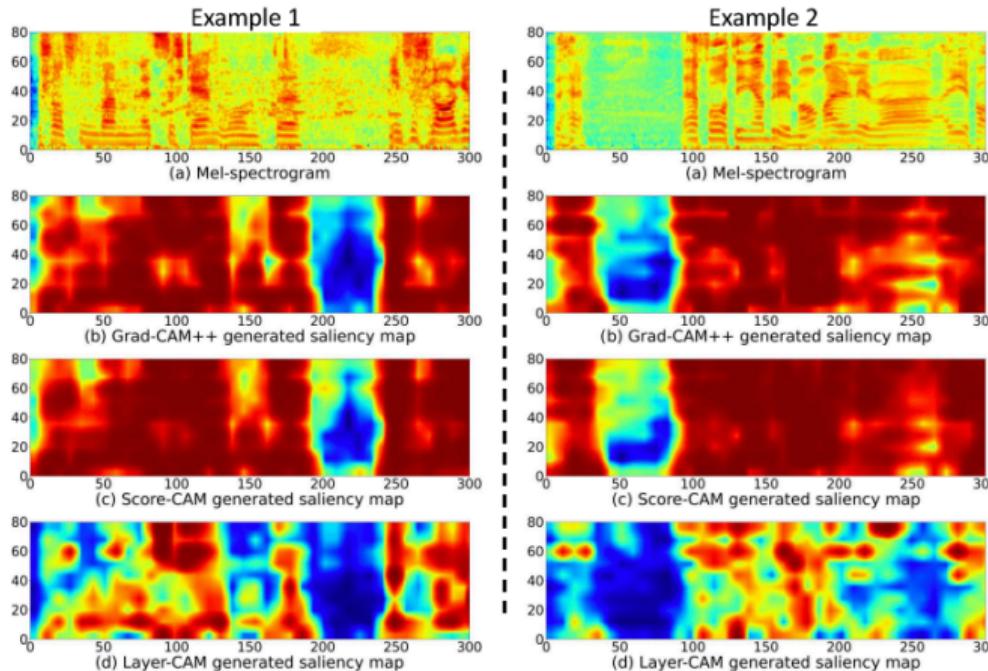
For the k -th activation map A^k (e.g., k indexes channels), define each entry w_{ij}^{kc} as

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

The saliency map is the linear combination

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

Interpretability: Visualizations in Speaker Recognition



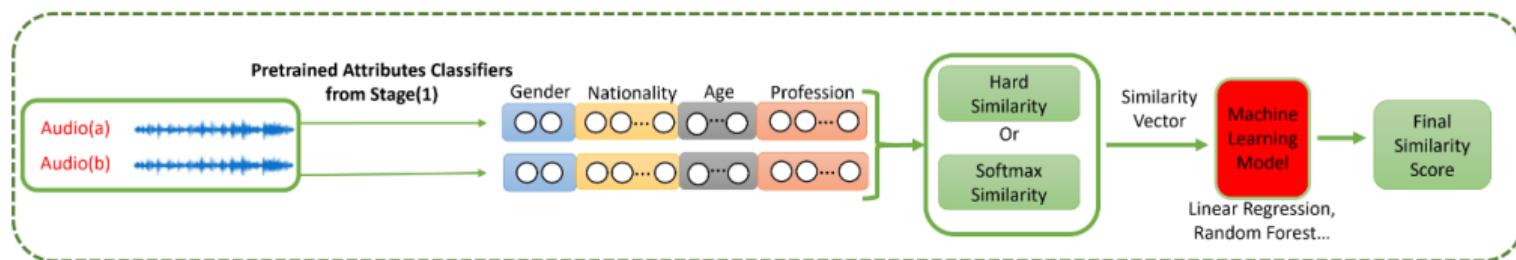
²⁹Li et al., "Reliable visualization for deep speaker recognition".

³⁰Li et al., "Visualizing data augmentation in deep speaker recognition".

Explainable Attribute-Based Speaker Verification³¹

Explainable attribute-based SV:

Use voice attributes (gender, age, nationality, profession) for transparent AI.



³¹Wu et al., "Explainable attribute-based speaker verification".

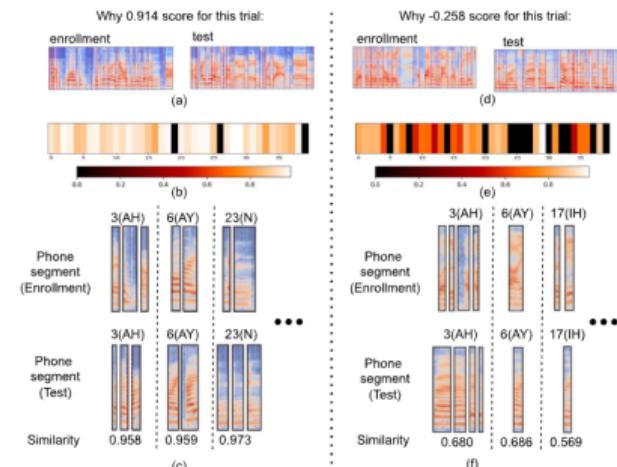
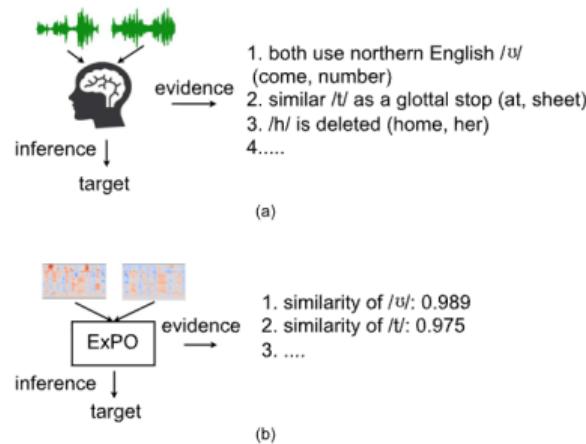
Interpretability: Decomposition of Speaker Information

Explainable Phonetic Trait-Oriented Network

ExPO^a

^aMa et al., "ExPO: Explainable Phonetic Trait-Oriented Network for Speaker Verification".

ExPO not only generates utterance-level speaker embeddings but also allows for fine-grained analysis and visualization of phonetic traits



- 1 Speaker Modeling: Background, Applications, and Trends
- 2 Discriminative Speaker Representation Learning
- 3 Self-supervised Speaker Representation Learning
- 4 Multi-modal Speaker Representation Learning
- 5 Robustness and Interpretability
- 6 Speaker Modeling in Related Tasks
- 7 Practice: WeSpeaker and WeSep

Speaker Modeling across Tasks

Different Tasks, Different Methods

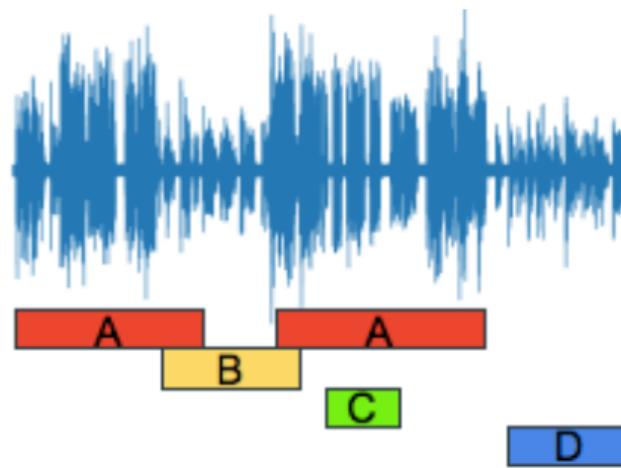


1. Pretrained speaker embeddings as extra input
2. Joint training to learn task-specific embeddings
3. Implicit speaker modeling

Speaker Modeling across Tasks

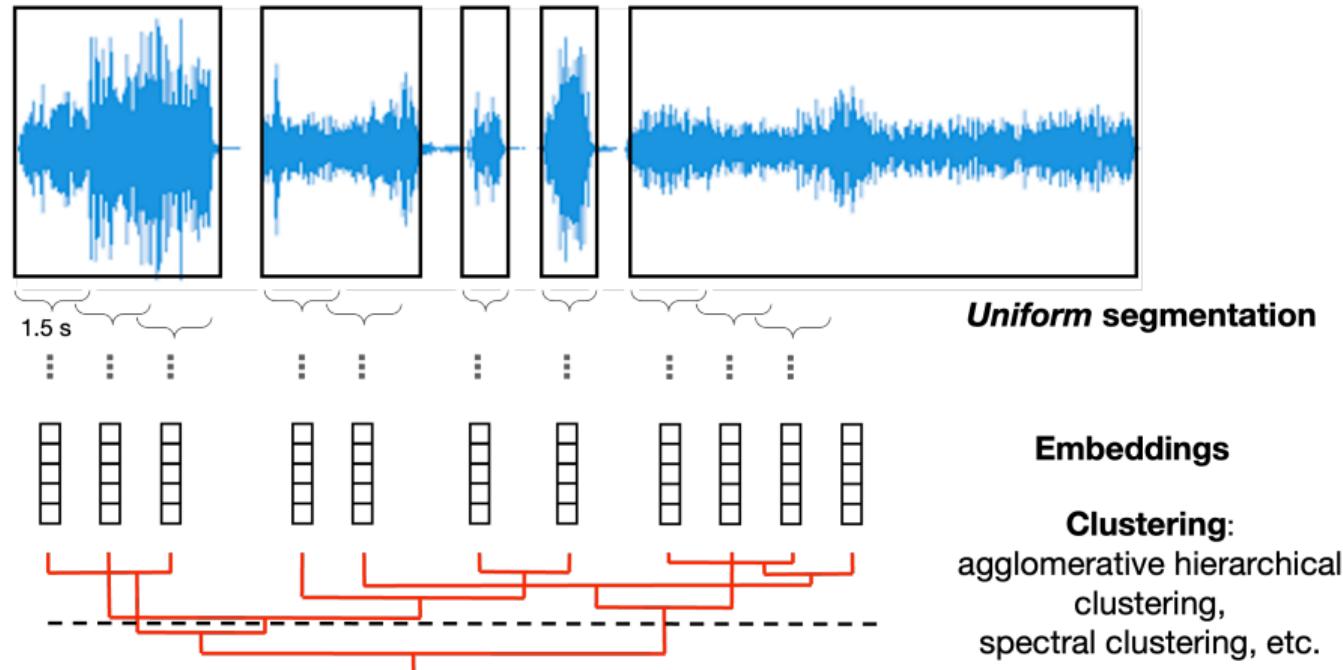
Speaker Diarization

- Speaker diarization is also known as speaker segmentation and clustering
- Answers the question: "Who spoke when?"
- In an audio, speakers A, B, C, D all speak; the goal is to:
 - Output the start and end time of each person's speech segments
- After diarization, per-speaker segments are easier for downstream processing



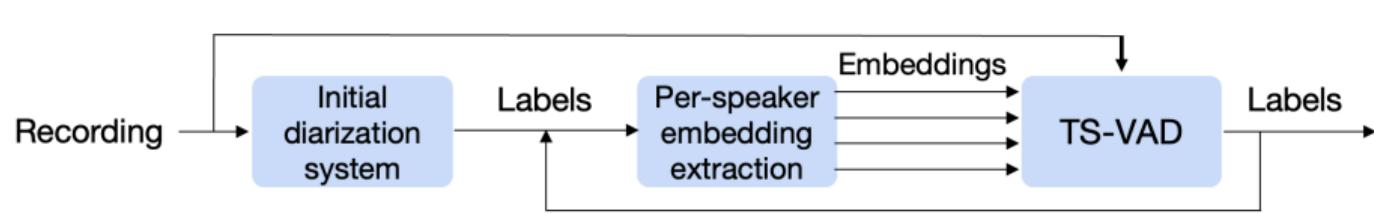
Speaker Modeling Across Tasks

Clustering based Speaker Diarization

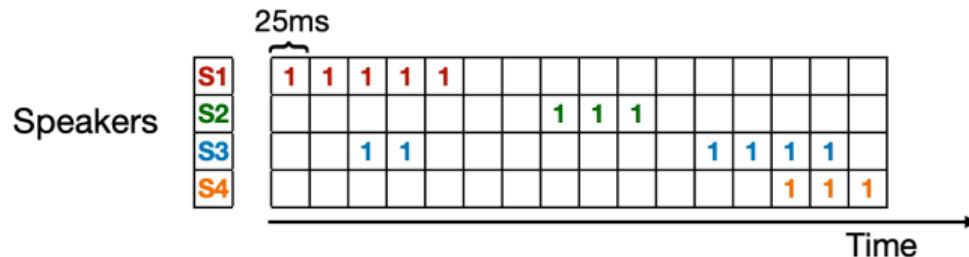


Speaker Modeling Across Tasks

TS-VAD based Speaker Diarization

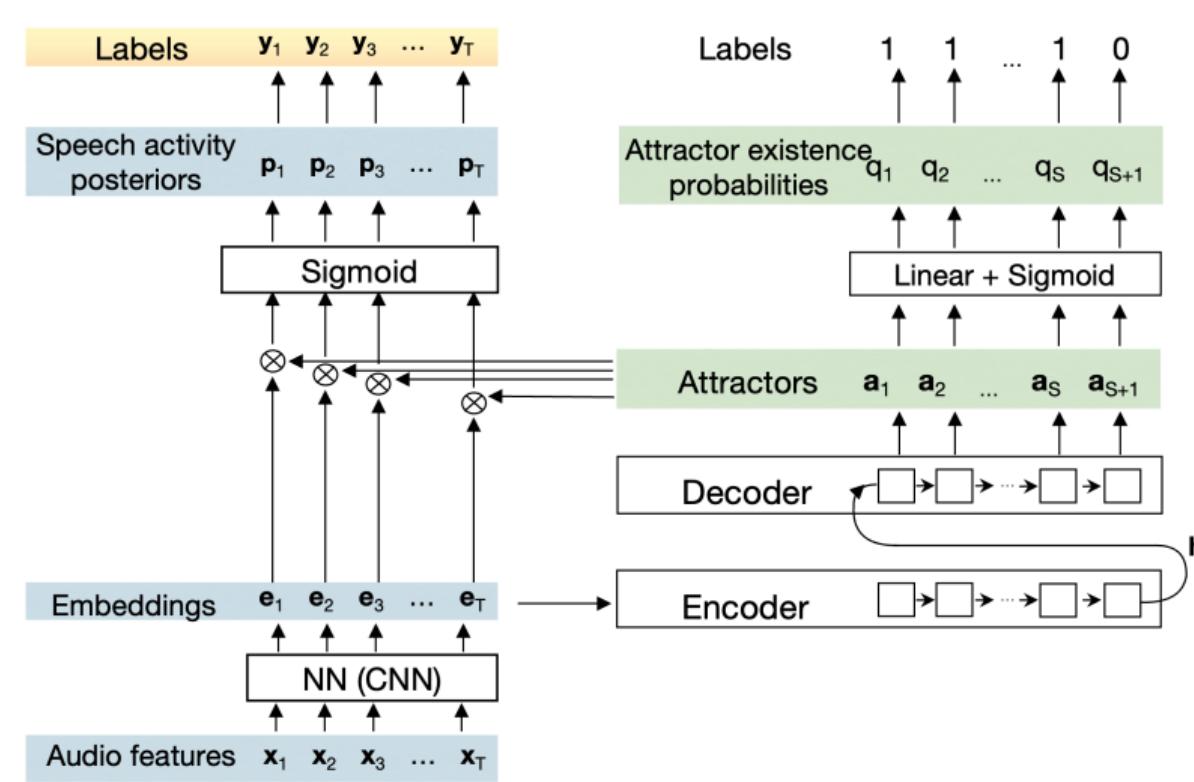


Binary VAD decisions are made independently for each speaker per-frame with the same NN



Speaker Modeling Across Tasks

EEND-EDA based Speaker Diarization



Speaker Modeling across Tasks

Multi-speaker ASR



SCHOOL OF
INTELLIGENCE SCIENCE AND TECHNOLOGY

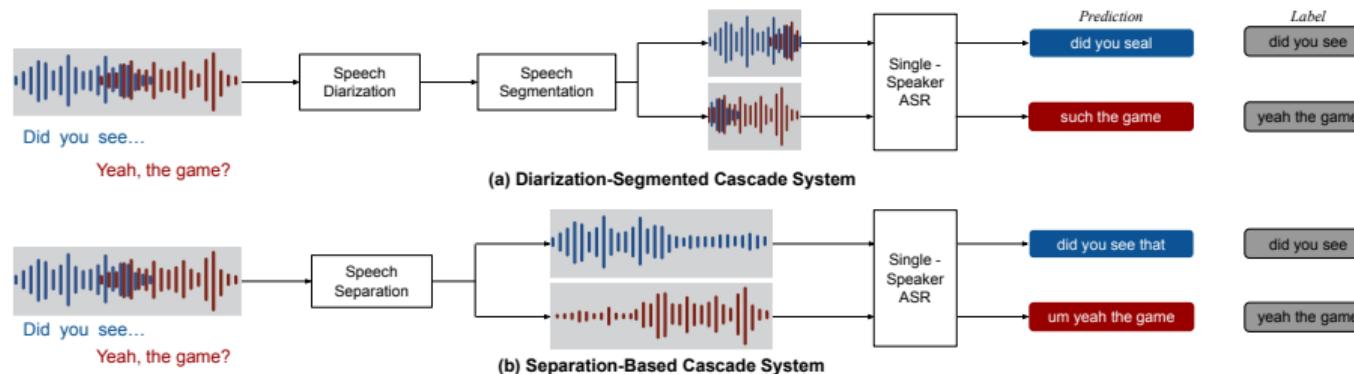
- Multi-speaker ASR: in the same audio, recognize **content** from all speakers and **attribute** it (who says what)
- Typical scenarios: meetings, group discussions, call-center logs, classroom interactions, podcast interviews
- Related to the "cocktail party effect": extracting target info from overlapping speech

Challenges

- **Overlapped speech**: overlapping speakers cause mixing
- **Speaker attribution**: not only "what was said" but also "who said it"
- **Data scarcity**: scarce large-scale, fine-grained multi-speaker annotations
- **Coupled subtasks**: recognition, separation, segmentation, attribution, overlap/boundary detection

Speaker Modeling across Tasks

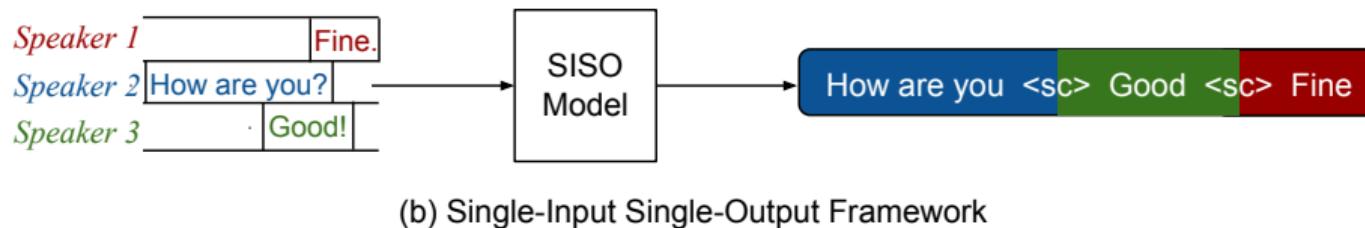
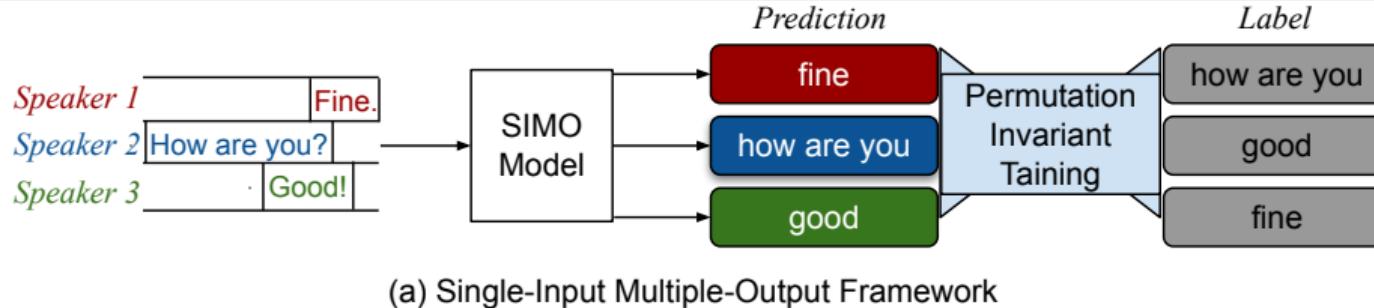
Multi-speaker ASR: Cascaded Systems



- Diarization-Segmented: first “speaker segmentation”, then single-speaker ASR
- Separation-based: first “speech separation”, then single-speaker ASR

Speaker Modeling across Tasks

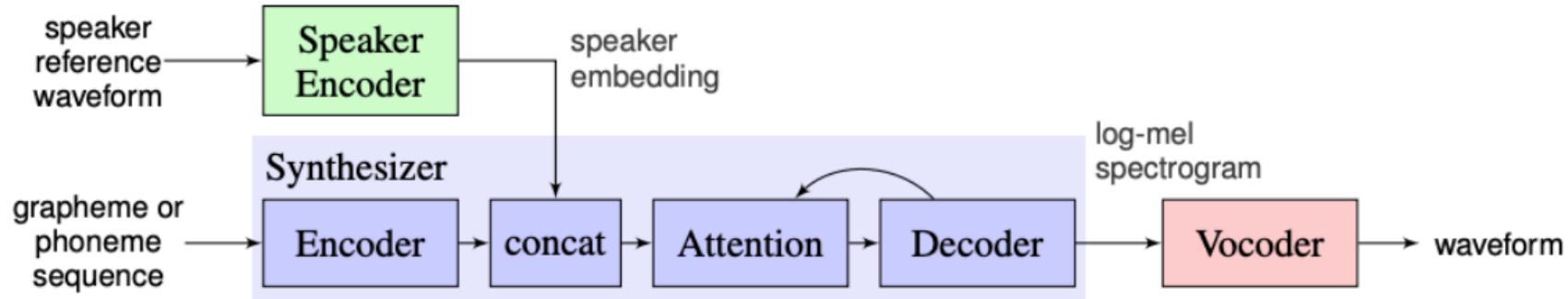
Multi-speaker ASR: End-to-end



- Directly map from **mixtures** to **speaker-attributed transcripts**, avoiding cascaded errors
- Support **joint optimization** of "who spoke" and "what was said"
- Two paradigms: **SIMO** (single-input multi-output) and **SISO** (single-input single-output)

Speaker Modeling across Tasks

Example: Explicit Speaker Modeling for Zero-shot TTS^{32, 33, 34}



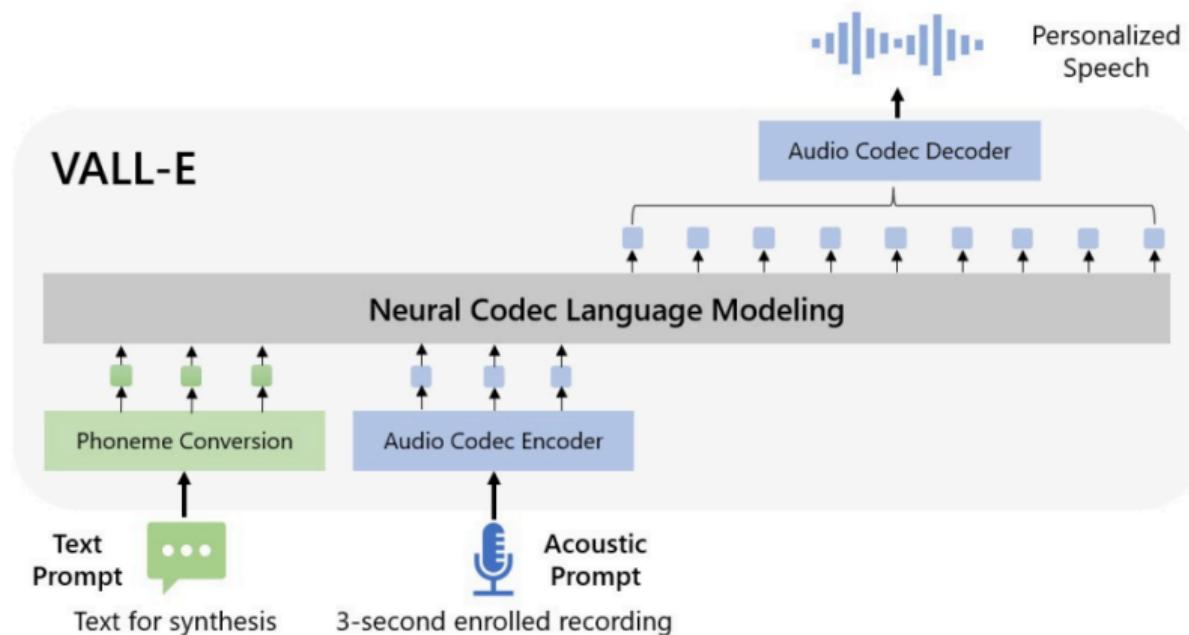
³²Jia et al., “Transfer learning from speaker verification to multispeaker text-to-speech synthesis”.

³³Casanova et al., “Your tts: Towards zero-shot multi-speaker tts and zero-shot voice conversion for everyone”.

³⁴Wu et al., “Adaspeech 4: Adaptive text to speech in zero-shot scenarios”.

Speaker Modeling across Tasks

Example: Implicit Speaker Modeling for Zero-shot TTS^{35, 36, 37}



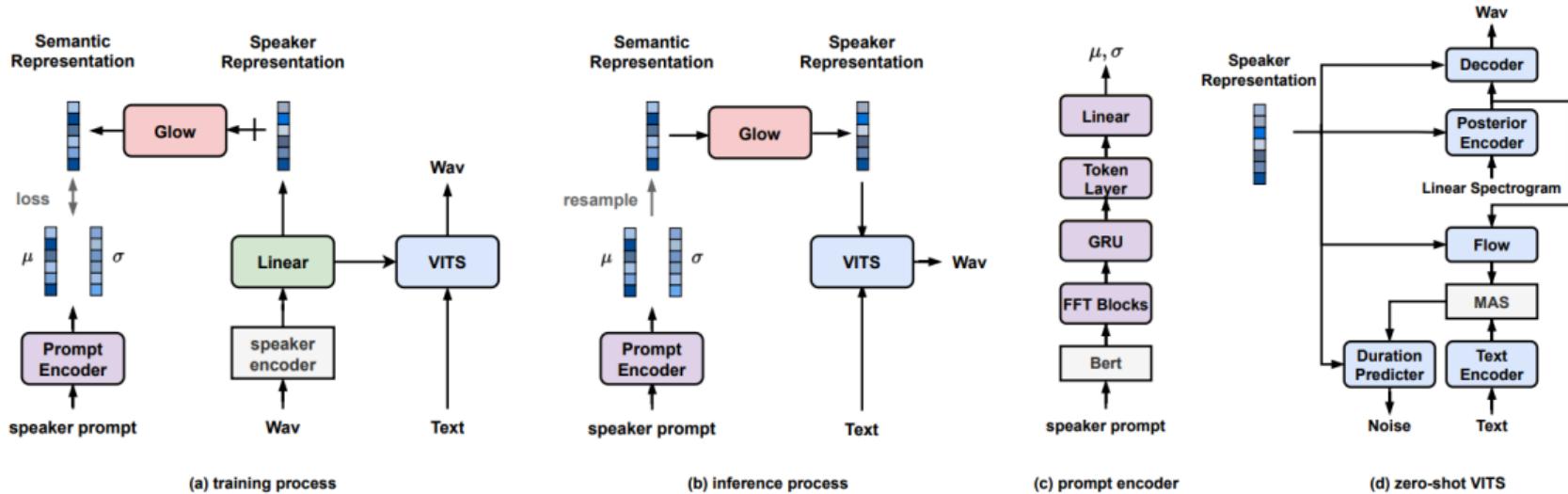
³⁵Wang et al., "Neural codec language models are zero-shot text to speech synthesizers".

³⁶Du et al., "UniCATS: A Unified Context-Aware Text-to-Speech Framework with Contextual VQ-Diffusion and Vocoding".

³⁷Le et al., "Voicebox: Text-guided multilingual universal speech generation at scale".

Speaker Modeling across Tasks

Example: Towards Controllability and Novel Voice Generation^{38, 39, 40, 41}



³⁸Zhang et al., "PromptSpeaker: Speaker Generation Based on Text Descriptions".

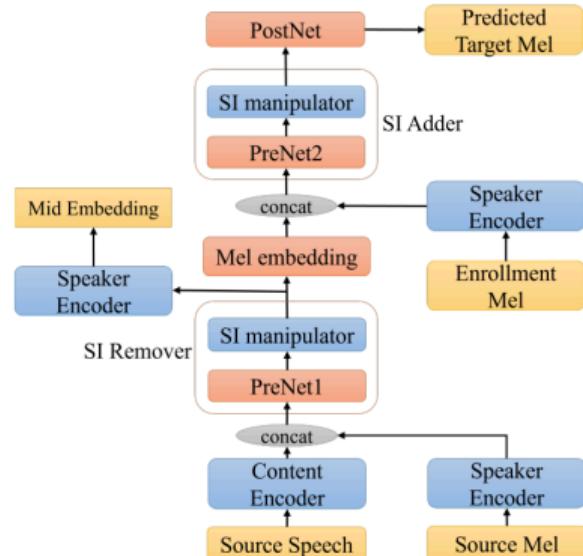
³⁹Stanton et al., "Speaker generation".

⁴⁰Shimizu et al., "PromptTTS++: Controlling Speaker Identity in Prompt-Based Text-to-Speech Using Natural Language Descriptions".

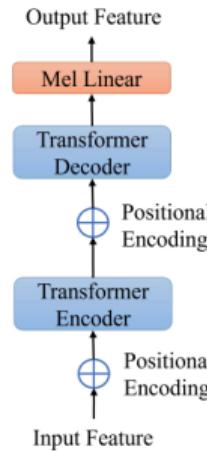
⁴¹Bilinski et al., "Creating new voices using normalizing flows".

Speaker Modeling across Tasks

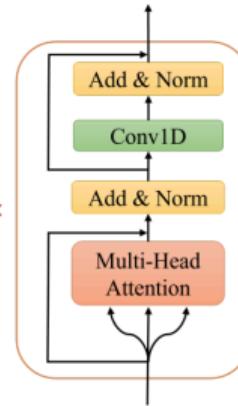
Example: Explicit Speaker Modeling for Zero-shot VC⁴²⁴³⁴⁴



(a) Proposed System



(b) SI Manipulator



(c) Encoder and decoder architecture

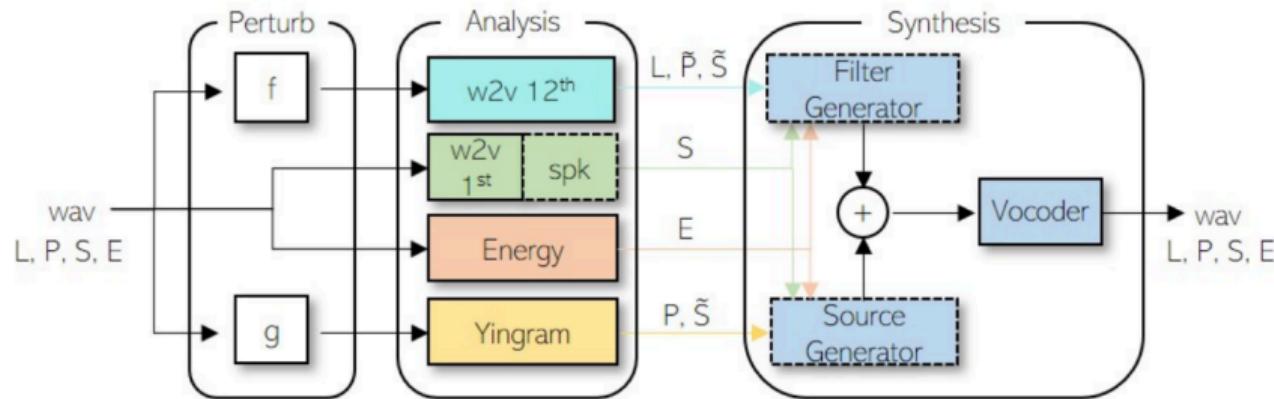
⁴²Zhang et al., "SIG-VC: A Speaker Information Guided Zero-Shot Voice Conversion System for Both Human Beings and Machines".

⁴³Chen and Duan, "ControlVC: Zero-Shot Voice Conversion with Time-Varying Controls on Pitch and Rhythm".

⁴⁴Hussain et al., "ACE-VC: Adaptive and Controllable Voice Conversion Using Explicitly Disentangled Self-Supervised Speech Representations"

Speaker Modeling across Tasks

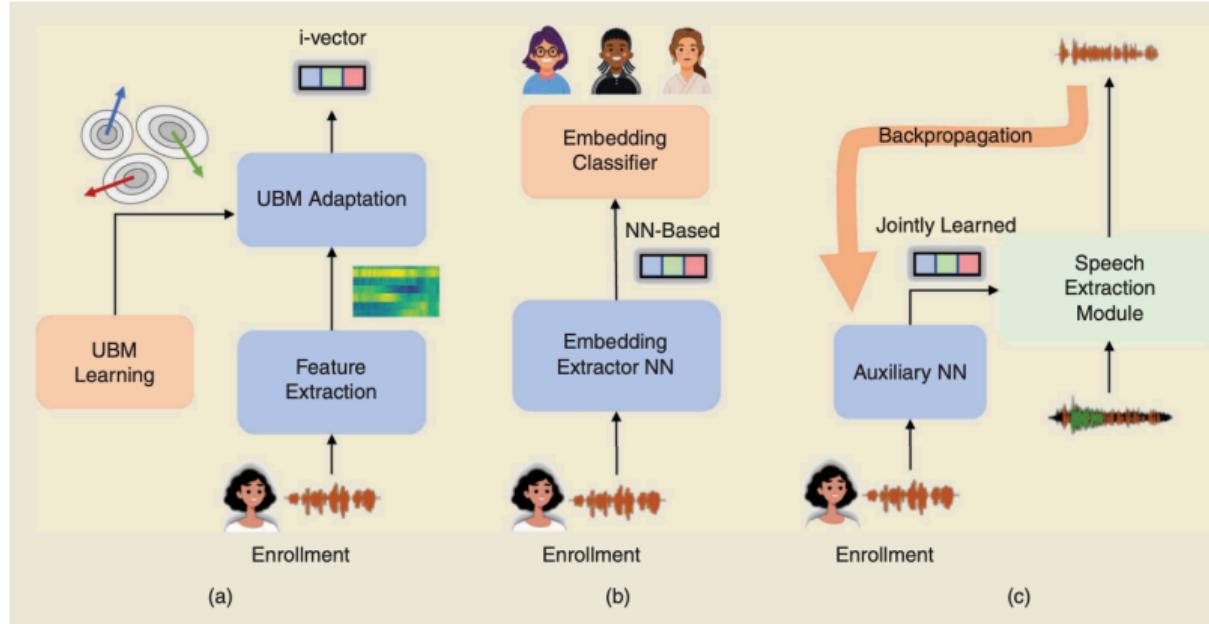
Example: Implicit Speaker Modeling for Zero-shot VC⁴⁵



⁴⁵Choi et al., 2021; Wu et al., 2020; Wu et al., 2020

Speaker Modeling across Tasks

Example: Explicit Speaker Modeling for Target Speaker Extraction⁴⁶⁴⁷⁴⁸



⁴⁶Zmolikova et al., "Neural Target Speech Extraction: An overview".

⁴⁷Delcroix et al., "Single channel target speaker extraction and recognition with speaker beam".

⁴⁸Ge et al., "Spex+: A complete time domain speaker extraction network".

Speaker Modeling across Tasks

Example: Implicit Speaker Modeling for Target Speaker Extraction⁴⁹

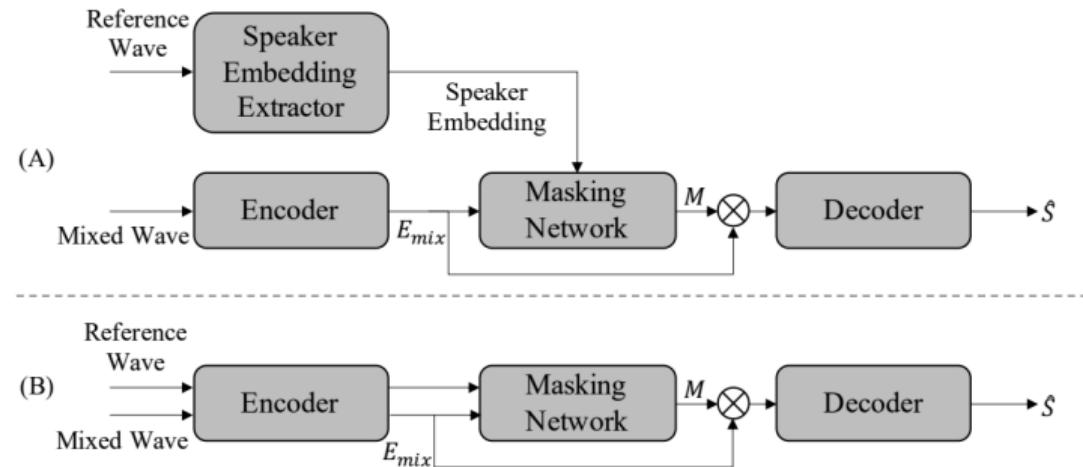


Figure 1: (A) is the diagram of a typical time-domain target speaker extraction method. (B) is the diagram of our proposed method. \otimes is an operation for element-wise product.

⁴⁹e.g., Zeng; Yang 2023

Speaker Modeling across Tasks

Myth: "One Embedding Fits All"



Common Assumption

Speaker embeddings optimized for recognition work for all tasks

Status quo:

- x-vector, ECAPA-TDNN, ResNet widely used
- Assumed universal applicability
- Directly transferred to other tasks

Problem

**$SV \neq TTS \neq VC \neq$
Target-Speaker Processing**

Key Question

What constitutes an "ideal" speaker representation?

Speaker Recognition

- Maximize discriminability
- Minimize within-speaker variance
- Robust to noise/channel
- Compact representations

Generative Tasks

- Capture fine-grained details
- Preserve prosody/emotion
- Natural synthesis
- Rich, expressive representations

Target-Speaker Processing

- Relative discriminability
- Maximize correlation with target in *mixtures*
- Compact representations

Conflicts

Discriminative optimization vs. generative richness
Absolute discriminability vs. relative discriminability

- SV/verification: suppress within-speaker variability
- TTS/VC: preserve and leverage such variability
- TSE, target-speaker ASR, speaker separation: relative discriminability within small sets

Speaker Modeling across Tasks

Speech Synthesis: Beyond Simple Embeddings



Challenge

Generate natural, expressive speech sounding like the target speaker

Key Requirements:

- Timbre accuracy
- Prosodic naturalness
- Emotional expressiveness
- Speaking style preservation

Evidence:

- Customized encoders > SV embeddings
- Prompt-based methods dominate zero-shot TTS

GAP

SV embeddings miss:

- Dynamics
- Prosodic details
- Emotional cues
- Style variations

Speech Synthesis: Embedding Comparison⁵⁰

Findings by Adriana STAN et al. (2023)

① Embedding choice does not affect the learning process

- Networks adapt to speaker conditioning regardless of embedding choice
- Similar synthesis quality can be achieved

② Speaker leakage is inevitable

- Core modules contain speaker information under standard training
- Simple conditioning cannot ensure perfect disentanglement

③ Inconsistency under zero-conditioning

- Core networks learn similar representations
- Speaker identity is unstable in zero-conditioning

⁵⁰Stan and O'Mahony, "An analysis on the effects of speaker embedding choice in non auto-regressive TTS".

Fundamental Issue

$$\text{Speech} = \text{Speaker} + \text{Content} + \text{Prosody} + \text{Accent} + \text{Emotion} + \dots$$

Current methods:

- Adversarial training
- Gradient reversal layer
- Multi-encoder architectures
- Self-supervised representations

Challenges:

- Perfect disentanglement is impossible
- Residual speaker information
- Content-speaker entanglement
- Prosody control complexity

Open Question

Can we ever achieve perfect disentanglement?

Paradigm Shift

From pre-computed embeddings to adaptive, context-aware modeling

Traditional:

- Pretrained speaker embeddings
- Fixed representations
- Limited context awareness
- Performance bottlenecks

Recent:

- USEF-TSE (embedding-free)
- Attention mechanisms
- Multi-level representations
- SSL-based features

Key Insight

Direct acoustic matching can outperform abstract embeddings

Limitations in Current Speaker-related Tasks



Limitation 1: Over-reliance on SV-optimized embeddings

Convenience Trap

Easy to use, but often suboptimal

Reasons:

- Availability of pretrained models
- Early success in SV
- Convenience

Evidence

USEF-TSE outperforms embedding-based methods

YourTTS custom encoders > SV embeddings

Costs:

- Suboptimal performance
- Information bottleneck
- Limited innovation
- Task mismatch

Limitation 2: Insufficient Dynamic Feature Capture

"Averaging" Problem

SV embeddings compress diverse acoustic expressions into one point

Lost Content:

- Emotional variation
- Prosody patterns
- Speaking rate changes
- Style nuances
- Contextual features

Impact on Applications:

- Monotonic TTS output
- Limited VC expressiveness
- Poor emotion control
- Unnatural prosody

Key Question

How to retain intra-speaker variability while preserving discriminability?

Limitation 3: Disentanglement Challenges

Intrinsic Complexity

Speech factors are inherently entangled rather than independently encoded

Entanglements:

- F0 contour: emotion + linguistic structure
- Spectral features: timbre + content
- Prosody: speaker + emotion + content
- No clear boundaries

Current Solutions:

- Adversarial learning
- Mutual information minimization
- Multi-encoder architectures
- Specialized losses

Reality

Perfect disentanglement remains open

Audio LLMs (ALLMs):

- Inspired by text LLMs (GPT, Qwen)
- Strong on diverse audio tasks:
 - ASR
 - Audio captioning
 - Music QA
- Excellent generalization
- Can identify speaker attributes (gender, age, accent)

Key Question

Can ALLMs perform speaker verification?

Observation

ALLM-based systems are still largely insensitive to speaker identity in dialogue.

Key idea: reformulate speaker verification as audio QA.

Four prompting strategies:

- ① **Separate:** feed two utterances independently
 - Prompt: “Audio1: [audio1], Audio2: [audio2]. Same speaker?”
- ② **Concat:** concatenate two utterances
 - Prompt: “How many speakers are in this audio?”
- ③ **Concat+Silence:** concatenate with 1s silence
 - Prompt: same as Concat
 - **Hypothesis: silence helps distinguish speakers**
- ④ **Mix:** overlap two utterances
 - Prompt: “This audio mixes two tracks. Same speaker?”

Rise of Audio LLMs

Zero-shot Results: Cross-dimension Performance



Model	Gender	Lang	Age	Device	Dur<2s	Dur>6s
Kimi (C+S)	70.20	68.40	63.40	52.67	53.70	73.60
Qwen2 (C+S)	59.40	58.60	53.87	52.20	50.60	59.10
Step (C+S)	64.20	60.40	56.80	57.47	54.60	71.80

Observations:

- 70% accuracy on long utterances
- Significant drop under challenging conditions:
 - Cross-device: 52–57%
 - Short duration: 50–55%
- Choose Kimi-Audio with Concat + Silence for fine-tuning

Conclusion

Zero-shot ALLMs have limited SV ability → motivates fine-tuning

Rise of Audio LLMs

Fine-tuning Effects: Significant Improvements

Model	Gender	Lang	Age	Device	Dialog	Dur<2s	Dur>6s
Kimi (zero-shot)	70.20	68.40	63.40	52.67	55.00	53.70	73.60
Kimi (fine-tuned)	95.07	97.00	92.40	88.20	89.00	80.90	89.50
Kimi (random sampling)	94.80	93.07	92.27	85.60	80.53	77.00	89.00
ECAPA-TDNN	99.33	99.27	94.13	94.67	93.00	78.80	95.60

Findings:

- ① **Huge improvements:** e.g., gender 70% → 95%
- ② **Hard negative sampling matters:** consistently better than random sampling
- ③ **ALLM surpasses ECAPA-TDNN on short duration!** (80.90% vs 78.80%)
- ④ **Still a gap on easy conditions** (e.g., gender: 95% vs 99%)

Insight

ALLMs show **stronger robustness** in challenging scenarios, indicating potential in noisy real-world settings.

Joint verification formulation:

- Enrollment: [audio1], Test: [audio2], Target text: “Hello world”
- Question: “Same speaker as enrollment? Does test match the target text?”
- Answer: “Speaker: yes/no, Content: yes/no”

Evaluation on LibriSpeech:

Model	Spk Acc (%)	Text Acc (%)	Overall (%)
Kimi (zero-shot)	62.09	89.61	52.31
Kimi (fine-tuned)	98.92	99.95	98.87
Whisper + ECAPA	99.08	99.75	98.83

- 1 Speaker Modeling: Background, Applications, and Trends
- 2 Discriminative Speaker Representation Learning
- 3 Self-supervised Speaker Representation Learning
- 4 Multi-modal Speaker Representation Learning
- 5 Robustness and Interpretability
- 6 Speaker Modeling in Related Tasks
- 7 Practice: WeSpeaker and WeSep

Toolkit	Speaker-specific	SSL	Pre-trained Models	Deployment
Kaldi	No	No	No	No
VoxCeleb_Trainer	Yes	No	No	No
ASV-Subtools	Yes	No	No	Yes
SpeechBrain	No	No	No	No
NeMo	No	No	No	Yes
EspNet	No	No	Yes	No
3D-Speaker	Yes	Yes	No	No
Wespeaker	Yes	Yes	Yes	Yes

Table: Common toolkits for speaker modeling

Dataset	Year	Speakers	Utterances	Duration
VoxCeleb1	2017	1,251	153,516	351h
VoxCeleb2	2018	6,112	1,128,246	2,442h
CN-Celeb1	2020	1,000	130,109	274h
CN-Celeb2	2020	2,000	529,485	1,090h
3D-Speaker	2023	10,000	579,013	1,124h
VoxBlink	2023	38,065	1,455,190	2,135h

Table: Representative Speaker Recognition Datasets

Performance Trends:

- VoxCeleb is approaching saturation
- Need more challenging scenarios
- Cross-genre and far-field datasets
- Large-scale unlabeled datasets for SSL

Wespeaker is a speaker embedding toolkit for both **research** and **production**, featuring

- Lightweight codebase
- SOTA performance
- Discriminative and SSL paradigms
- Runtime/deployment support
- Adopted by many groups in industry and academia:

The screenshot shows a PyTorch Model Zoo page for the "wespeaker-voxceleb-resnet34-LM" model. At the top, there's a navigation bar with "pyannote" and a search bar containing "wespeaker-voxceleb-resnet34-LM". Below the search bar are several buttons for different frameworks: "pyannote.audio" (selected), "PyTorch", "voxceleb", "pyannote", "pyannote-audio-model", "wespeaker", "audio", "voice", "speech", and "speaker". There are also buttons for "speaker-recognition", "speaker-verification", "speaker-identification", "speaker-embedding", and a "License: cc-by-4.0" button. A "Model card" tab is selected, followed by "Files" and "Community". On the right, there's a "Use this model" button. Below the tabs, there's a section asking if the user is using the model in production and suggesting switching to "pyannoteAI". To the right, there's a chart showing "Downloads last month" at 14,362,520 and a line graph showing the trend over time.

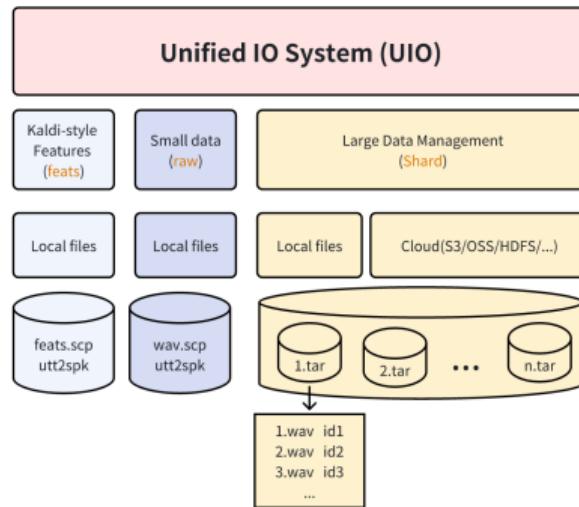


Figure: Unified I/O system

Unified I/O

- Also adopted in WeNet ASR
- Inspired by webdataset and tfrecord

Idea

- Raw: load wav and label files from disk (small-scale)
- Shard:
 - Pack many small files into larger shards
 - Read and decompress shards on the fly
- Feat: compatible with Kaldi-style features
- Efficiently load 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: Convert train/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
```

Step 3: Start training

```

if [ ${stage} -le 3 ] && [ ${stop_stage} -ge 3 ]; then
    echo "Start training ..."
    num_gpus=$(echo $gpus | awk -F ',' '{print NF}')
    torchrn --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}/rir2/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'])

```

Architectures:

- ResNet family
- TDNN
- ECAPA-TDNN
- RepVGG
- CAM++
- ReDimNet
- Pretrained frontend (e.g., WavLM)

Pooling:

- TSTP
- ASTP
- MQMHASTP

Losses:

- 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-ends:

- Cosine
- LDA
- PLDA
- PSDA
- 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
```

Model	Params	vox1-O-clean	vox1-E-clean	vox1-H-clean
ReDimNetB0	1.0M	1.128	1.181	2.008
ReDimNetB3	3.2M	0.537	0.790	1.433
XVEC	4.61M	1.590	1.641	2.726
Res2Net34_Base	4.68M	1.234	1.232	2.162
ECAPA_TDNN_GLOB_c512	6.19M	0.782	1.005	1.824
RepVGG_TINY_A0	6.26M	0.824	0.953	1.709
Gemini_DFResNet114	6.53M	0.638	0.839	1.427
ResNet34	6.63M	0.659	0.821	1.437
ERes2Net34_Base	7.88M	0.744	0.896	1.603
CAM++	7.18M	0.659	0.803	1.569
ECAPA_TDNN_GLOB_c1024	14.6M	0.707	0.894	1.615
ResNet221	23.8M	0.505	0.676	1.213
SimAM_ResNet34 (VoxBlink2 Pretrain)	25.2M	0.372	0.559	0.997
ResNet293	28.6M	0.425	0.641	1.146
SimAM_ResNet100 (VoxBlink2 Pretrain)	50.2M	0.202	0.421	0.795
WavLM+EcapaTDNN		0.415	0.551	1.118

Figure: List of supported models

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
```

Wespeaker

Deployment and Product-oriented Setup

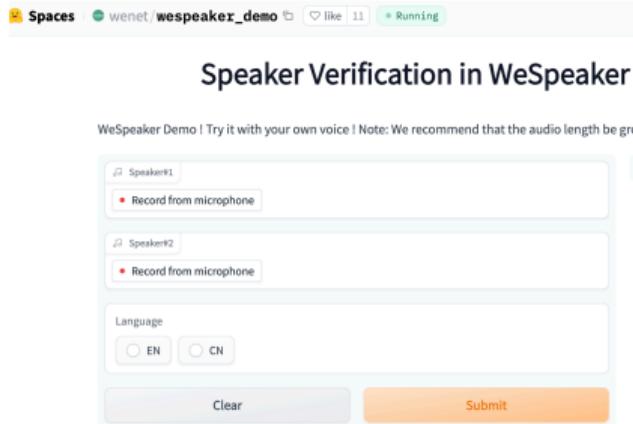


Figure: Wespeaker demo

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 API:

```
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')
```

The first open-source toolkit for Target Speaker Extraction.⁵¹

Main Contributions

- Propose WeSep, focusing on target speaker extraction
- Provide versatile speaker modeling capabilities
- Implement online data simulation and scalability
- Offer end-to-end training and deployment support

Technical Highlights

- Seamless integration with WeSpeaker
- Unified I/O data management
- Dynamic speaker mixing strategies
- Multiple fusion method support

⁵¹Wang et al., "Wesep: A scalable and flexible toolkit towards generalizable target speaker extraction".

UIO framework

- Efficient for both lab-scale and production-scale datasets
- Supports tens of thousands of hours of data
- Handles massive small-file shards
- Integrated in WeNet and WeSpeaker

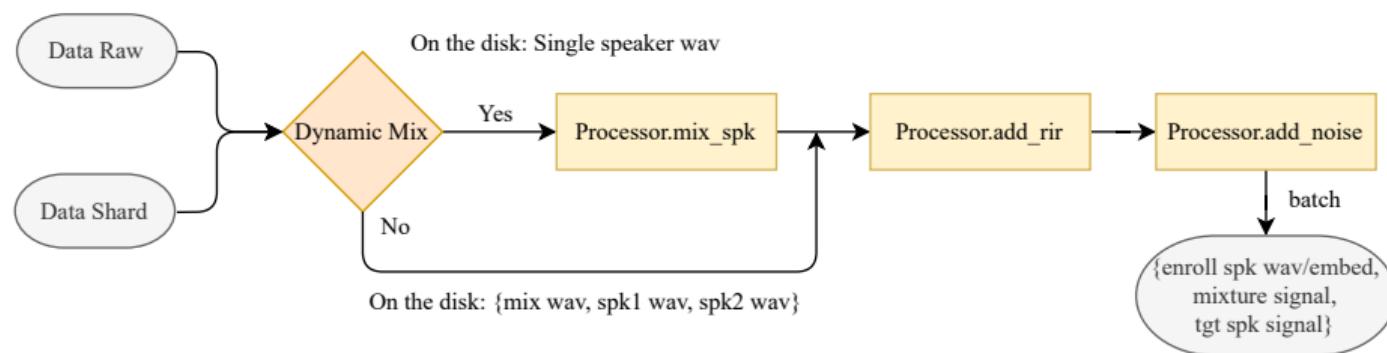


Figure: WeSep online data preparation pipeline (2-speaker example)

Issues of traditional offline simulation:

- Store preprocessed data
- Large disk footprint
- Limited diversity

Advantages of online simulation:

- Save storage
- Create diverse training data
- Improve robustness
- Flexible data generation

Supported functions

- Online noise addition
- Reverb generation (standard RIR + fast random approx.)
- Dynamic speaker mixing

```
pip install git+https://github.com/wenet-e2e/wespeaker.git
```

```
# pseudo-codes for integrating wespeaker models
from wespeaker import get_speaker_model
# TDNN/ECAPA/ResNet/CAM++/WavLM...
s = get_speaker_model(spk_model_name) (**spk_args)
m = BSRNN(**sep_args) # Or other backbones
m.speaker_model = s
if use_pretrain_spk_encoder:
    m.spk_model.load_state_dict(pretrain_path)
    m.speaker_model.freeze()
```

```
spk_fuse_type: 'multiply'
use_spk_transform: False
multi_fuse: False
joint_training: True
##### ResNet
spk_model: ResNet34
spk_model_init: False
./wespeaker_models/model.pt
```

① ConvTasNet

- Convolutional network in time domain
- Learn and estimate separation masks
- Support Spex+ variants

② BSRNN

- Band-split RNN
- Explicitly split spectrogram into bands
- Fine-grained modeling

③ DPCCN

- Dense-connected pyramid complex CNN
- Combine DenseUNet, TCN and DenseNet
- Improved separation

④ TF-GridNet

- Operate in T-F domain
- Stack multi-path blocks
- Leverage local/global T-F information

Given \mathbf{e}_s and intermediate outputs $\mathbf{H} = \{\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_T\}$:

- ① **Concat**: replicate \mathbf{e}_s and concatenate
- ② **Add**: project and add element-wise
- ③ **Multiply**: project and multiply element-wise
- ④ **FiLM**: feature-wise linear modulation

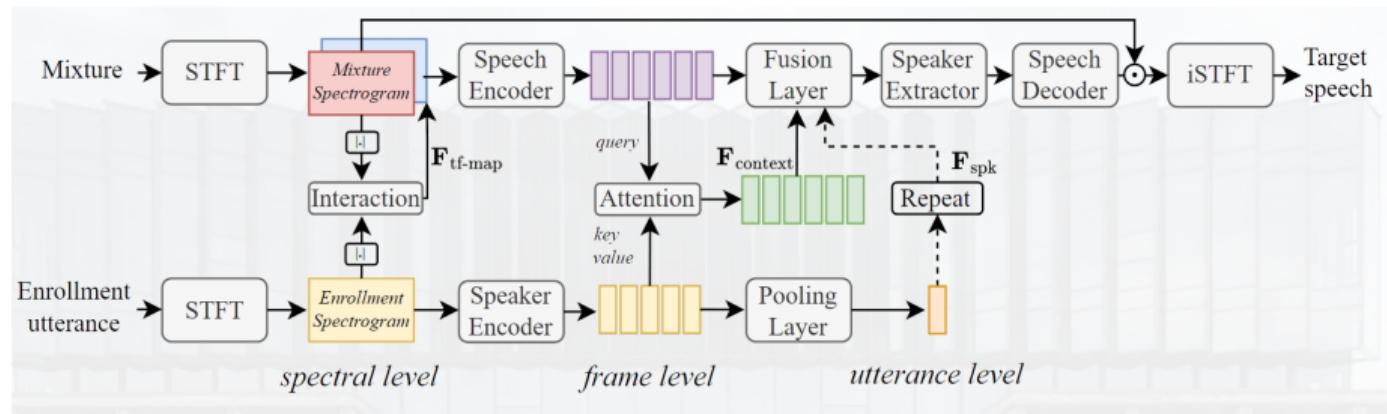
$$\mathbf{h}'_t = \gamma(\mathbf{e}_s) \odot \mathbf{h}_t + \beta(\mathbf{e}_s) \quad (3)$$

Why FiLM

- γ and β are functions of \mathbf{e}_s
- \odot is element-wise multiplication
- Learn an affine transform conditioned on \mathbf{e}_s

Beyond embedding-level guidance, add finer-grained context guidance^a

^aZhang et al., "Multi-level speaker representation for target speaker extraction".



Multi-level Speaker Modeling

Speaker embedding (Baseline)



TF-Map (Proposed)



Model	Speaker Model	Speaker Model	SI-SDRi on Libri2mix	Accuracy / %	Pub.
TD-SpeakerBeam	ResNet	Joint	13.03	95.2	ICASSP, 2020
SpEx+*	ResNet	Joint	13.41	-	Interspeech, 2020
sDPCCN	ConvNet	Joint	11.61	-	ICASSP, 2022
Target-Confusion*	ResNet	Joint	13.88	-	Interspeech, 2022
MC-SpEx*	ResNet	Joint	14.61	-	Interspeech, 2023
X-T-TasNet	d-vector	Pretrained	13.48	95.3	Interspeech, 2024
SSL-TD-SpeakerBeam	ResNet + WavLM	Pretrained	14.01	96.1	ICASSP, 2024
		Pretrained + Fine-tuning	14.65	97.0	
BSRNN	Campplus + SHuBERT	Pretrained	15.39	-	SPL, 2024
BSRNN	Ecapa-TDNN	Pretrained	15.91	97.0	Proposed
BSRNN			17.99	98.6	

- ✓ SOTA Performance, with simple but effective multi-level speaker modeling
- ✓ The generalization ability is largely enhanced
(The gap between training and validation error)

Overall Summary

A holistic view of speaker modeling

- **Evolution:** from traditional GMM-UBM to deep learning
- **Multi-task applications:** speaker recognition, separation, synthesis, conversion, etc.
- **Challenges:** robustness, efficiency, interpretability, multi-modal fusion
- **Trends:** self-supervised learning, joint modeling, tooling

Key Takeaways

Rethinking speaker modeling



- **Beyond recognition:** speaker modeling is more than verification
- **Beyond embeddings:** speaker modeling is more than embedding learning
- **Task-oriented:** tailor methods to the target application
- **Evaluation-oriented:** use task-specific metrics and protocols

-  Baevski, Alexei et al. "wav2vec 2.0: A framework for self-supervised learning of speech representations". In: [Advances in neural information processing systems 33 \(2020\)](#), pp. 12449–12460.
-  Belinkov, Yonatan and James Glass. "Analyzing hidden representations in end-to-end automatic speech recognition systems". In: [Advances in Neural Information Processing Systems 30 \(2017\)](#).
-  Bilinski, Piotr et al. "Creating new voices using normalizing flows". In: (2022).
-  Cai, Danwei, Weiqing Wang, and Ming Li. "An iterative framework for self-supervised deep speaker representation learning". In: [ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing](#) IEEE. 2021, pp. 6728–6732.

-  Cai, Danwei et al. "Pretraining Conformer with ASR for Speaker Verification". In: [ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing](#). IEEE. 2023, pp. 1–5.
-  Caron, Mathilde et al. "Emerging properties in self-supervised vision transformers". In: [Proceedings of the IEEE/CVF international conference on computer vision](#). 2021, pp. 9650–9660.
-  Casanova, Edresson et al. "Yourtts: Towards zero-shot multi-speaker tts and zero-shot voice conversion for everyone". In: [International Conference on Machine Learning](#). PMLR. 2022, pp. 2709–2720.
-  Chen, Meiying and Zhiyao Duan. "ControlVC: Zero-Shot Voice Conversion with Time-Varying Controls on Pitch and Rhythm". In: [arXiv preprint arXiv:2209.11866](#) (2022).

-  Chen, Sanyuan et al. "Unispeech-sat: Universal speech representation learning with speaker aware pre-training". In: [ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing](#) IEEE. 2022, pp. 6152–6156.
-  Chen, Sanyuan et al. "Wavlm: Large-scale self-supervised pre-training for full stack speech processing". In: [IEEE Journal of Selected Topics in Signal Processing](#) 16.6 (2022), pp. 1505–1518.
-  Chen, Ting et al. "A simple framework for contrastive learning of visual representations". In: [International conference on machine learning](#). PMLR. 2020, pp. 1597–1607.
-  Chen, Zhengyang et al. "A comprehensive study on self-supervised distillation for speaker representation learning". In: [2022 IEEE Spoken Language Technology Workshop \(SLT\)](#). IEEE. 2023, pp. 599–604.

-  Chen, Zhengyang et al. "Channel invariant speaker embedding learning with joint multi-task and adversarial training". In: ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing IEEE. 2020, pp. 6574–6578.
-  Chen, Zhengyang et al. "Large-scale self-supervised speech representation learning for automatic speaker verification". In: ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing IEEE. 2022, pp. 6147–6151.
-  Chowdhury, Shammur Absar, Nadir Durrani, and Ahmed Ali. "What do end-to-end speech models learn about speaker, language and channel information? a layer-wise and neuron-level analysis". In: Computer Speech Language 83 (2023), p. 101539.
-  Chung, Joon Son et al. "In defence of metric learning for speaker recognition". In: arXiv preprint arXiv:2003.11982 (2020).

-  Dehak, Najim et al. "Front-end factor analysis for speaker verification". In: [IEEE Transactions on Audio, Speech, and Language Processing 19.4 \(2010\)](#), pp. 788–798.
-  Delcroix, Marc et al. "Single channel target speaker extraction and recognition with speaker beam". In: [2018 IEEE international conference on acoustics, speech and signal processing \(ICASSP\)](#). IEEE. 2018, pp. 5554–5558.
-  Desplanques, Brecht, Jenthe Thienpondt, and Kris Demuynck. "Ecapa-tdnn: Emphasized channel attention, propagation and aggregation in tdnn based speaker verification". In: [arXiv preprint arXiv:2005.07143 \(2020\)](#).
-  Du, Chenpeng et al. "UniCATS: A Unified Context-Aware Text-to-Speech Framework with Contextual VQ-Diffusion and Vocoding". In: [arXiv preprint arXiv:2306.07547 \(2023\)](#).
-  Ge, Meng et al. "Spex+: A complete time domain speaker extraction network". In: [arXiv preprint arXiv:2005.04686 \(2020\)](#).

-  Han, Bing, Zhengyang Chen, and Yanmin Qian. "Self-Supervised Learning with Cluster-Aware-DINO for High-Performance Robust Speaker Verification". In: [arXiv preprint arXiv:2304.05754](#) (2023).
-  — . "Self-supervised speaker verification using dynamic loss-gate and label correction". In: [arXiv preprint arXiv:2208.01928](#) (2022).
-  Hsu, Wei-Ning et al. "Hubert: Self-supervised speech representation learning by masked prediction of hidden units". In: [IEEE/ACM Transactions on Audio, Speech, and Language Processing](#) 29 (2021), pp. 3451–3460.
-  Huang, Wen et al. "Enhancing Cross-Domain Speaker Verification through Multi-Level Domain Adapters". In: (2023).
-  Huh, Jaesung et al. "Augmentation adversarial training for self-supervised speaker recognition". In: [arXiv preprint arXiv:2007.12085](#) (2020).

-  Hussain, Shehzeen et al. "ACE-VC: Adaptive and Controllable Voice Conversion Using Explicitly Disentangled Self-Supervised Speech Representations". In: ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing. IEEE. 2023, pp. 1–5.
-  Jia, Ye et al. "Transfer learning from speaker verification to multispeaker text-to-speech synthesis". In: Advances in neural information processing systems 31 (2018).
-  Jin, Yufeng et al. "Cross-modal distillation for speaker recognition". In: Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 37. 11. 2023, pp. 12977–12985.
-  Jin, Zezhong, Youzhi Tu, and Man-Wai Mak. "Phonetic-aware speaker embedding for far-field speaker verification". In: arXiv preprint arXiv:2311.15627 (2023).
-  Le, Matthew et al. "Voicebox: Text-guided multilingual universal speech generation at scale". In: arXiv preprint arXiv:2306.15687 (2023).

-  Li, Jianchen, Jiqing Han, and Hongwei Song. "CDMA: Cross-Domain Distance Metric Adaptation for Speaker Verification". In: [ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing](#) IEEE. 2022, pp. 7197–7201.
-  Li, Pengqi et al. "Reliable visualization for deep speaker recognition". In: [arXiv preprint arXiv:2204.03852](#) (2022).
-  — . "Visualizing data augmentation in deep speaker recognition". In: [arXiv preprint arXiv:2305.16070](#) (2023).
-  Li, Rongjin, Weibin Zhang, and Dongpeng Chen. "The coral++ algorithm for unsupervised domain adaptation of speaker recognition". In: [ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing](#) IEEE. 2022, pp. 7172–7176.

-  Li, Ze, Ming Cheng, and Ming Li. "Enhancing Speaker Verification with w2v-BERT 2.0 and Knowledge Distillation guided Structured Pruning". In: [arXiv preprint arXiv:2510.04213 \(2025\)](#).
-  Liao, Dexin et al. "Towards a unified conformer structure: from asr to asv task". In: [ICASSP 2023. IEEE. 2023, pp. 1–5](#).
-  Ma, Yi et al. "ExPO: Explainable Phonetic Trait-Oriented Network for Speaker Verification". In: [IEEE Signal Processing Letters \(2025\)](#).
-  Peng, Junyi et al. "Parameter-efficient transfer learning of pre-trained Transformer models for speaker verification using adapters". In: [ICASSP 2023. IEEE. 2023, pp. 1–5](#).
-  Qian, Yanmin, Zhengyang Chen, and Shuai Wang. "Audio-visual deep neural network for robust person verification". In: [IEEE/ACM Transactions on Audio, Speech, and Language Processing 29 \(2021\), pp. 1079–1092](#).

-  Raj, Desh et al. "Probing the information encoded in x-vectors". In: [2019 IEEE Automatic Speech Recognition and Understanding Workshop \(ASRU\)](#). IEEE. 2019, pp. 726–733.
-  Reynolds, Douglas A et al. "Gaussian mixture models.". In: [Encyclopedia of biometrics](#) 741.659-663 (2009).
-  Rohdin, Johan et al. "Speaker verification using end-to-end adversarial language adaptation". In: [ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing](#). IEEE. 2019, pp. 6006–6010.
-  Shimizu, Reo et al. "PromptTTS++: Controlling Speaker Identity in Prompt-Based Text-to-Speech Using Natural Language Descriptions". In: [arXiv preprint arXiv:2309.08140](#) (2023).

-  Snyder, David et al. "X-vectors: Robust dnn embeddings for speaker recognition". In: [2018 IEEE international conference on acoustics, speech and signal processing \(ICASSP\)](#). IEEE. 2018, pp. 5329–5333.
-  Stan, Adriana and Johannah O'Mahony. "An analysis on the effects of speaker embedding choice in non auto-regressive TTS". In: [arXiv preprint arXiv:2307.09898](#) (2023).
-  Stanton, Daisy et al. "Speaker generation". In: [ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing](#). IEEE. 2022, pp. 7897–7901.
-  Tao, Ruijie et al. "Self-supervised speaker recognition with loss-gated learning". In: [ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing](#). IEEE. 2022, pp. 6142–6146.

-  Variani, Ehsan et al. "Deep neural networks for small footprint text-dependent speaker verification". In: [2014 IEEE international conference on acoustics, speech and signal processing \(ICASSP\)](#). IEEE. 2014, pp. 4052–4056.
-  Wan, Li et al. "Generalized end-to-end loss for speaker verification". In: [2018 IEEE International Conference on Acoustics, Speech and Signal Processing \(ICASSP\)](#). IEEE. 2018, pp. 4879–4883.
-  Wang, Chengyi et al. "Neural codec language models are zero-shot text to speech synthesizers". In: [arXiv preprint arXiv:2301.02111 \(2023\)](#).
-  Wang, Shuai, Yanmin Qian, and Kai Yu. "What does the speaker embedding encode?" In: [arXiv preprint arXiv:2301.02111 \(2023\)](#).
-  Wang, Shuai et al. "Wesep: A scalable and flexible toolkit towards generalizable target speaker extraction". In: [arXiv preprint arXiv:2409.15799 \(2024\)](#).

-  Wu, Xiaoliang et al. "Explainable attribute-based speaker verification". In: arXiv preprint arXiv:2405.19796 (2024).
-  Wu, Yihan et al. "Adaspeech 4: Adaptive text to speech in zero-shot scenarios". In: arXiv preprint arXiv:2204.00436 (2022).
-  Xia, Wei, Jing Huang, and John HL Hansen. "Cross-lingual text-independent speaker verification using unsupervised adversarial discriminative domain adaptation". In: ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing. IEEE. 2019, pp. 5816–5820.
-  Zhang, Haoran, Yuexian Zou, and He-lin Wang. "Contrastive self-supervised learning for text-independent speaker verification". In: ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing. IEEE. 2021, pp. 6713–6717.

-  Zhang, Haozhe et al. "SIG-VC: A Speaker Information Guided Zero-Shot Voice Conversion System for Both Human Beings and Machines". In: [ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing](#) IEEE. 2022, pp. 6567–65571.
-  Zhang, Ke et al. "Multi-level speaker representation for target speaker extraction". In: [ICASSP 2025-2025 IEEE International Conference on Acoustics, Speech and Signal Processing](#) IEEE. 2025, pp. 1–5.
-  Zhang, Leying, Zhengyang Chen, and Yanmin Qian. "Knowledge Distillation from Multi-Modality to Single-Modality for Person Verification". In: [Proc. Interspeech 2021 \(2021\)](#), pp. 1897–1901.
-  Zhang, Yongmao et al. "PromptSpeaker: Speaker Generation Based on Text Descriptions". In: [arXiv preprint arXiv:2310.05001](#) (2023).

-  Zhao, Zifeng et al. "Probing Deep Speaker Embeddings for Speaker-related Tasks". In: [arXiv preprint arXiv:2212.07068 \(2022\)](#).
-  Zheng, Rong, Shuwu Zhang, and Bo Xu. "Text-independent speaker identification using gmm-ubm and frame level likelihood normalization". In: [2004 International Symposium on Chinese Spoken Language Processing](#). IEEE. 2004, pp. 289–292.
-  Zheng, Siqi, Yun Lei, and Hongbin Suo. "Phonetically-Aware Coupled Network For Short Duration Text-Independent Speaker Verification.". In: [INTERSPEECH](#). 2020, pp. 926–930.
-  Zmolikova, Katerina et al. "Neural Target Speech Extraction: An overview". In: [IEEE Signal Processing Magazine](#) 40.3 (2023), pp. 8–29.