# Voice Biometrics

Phan Trung Kiên - Phòng XLTN – Khối CNTLA

18/10/2024
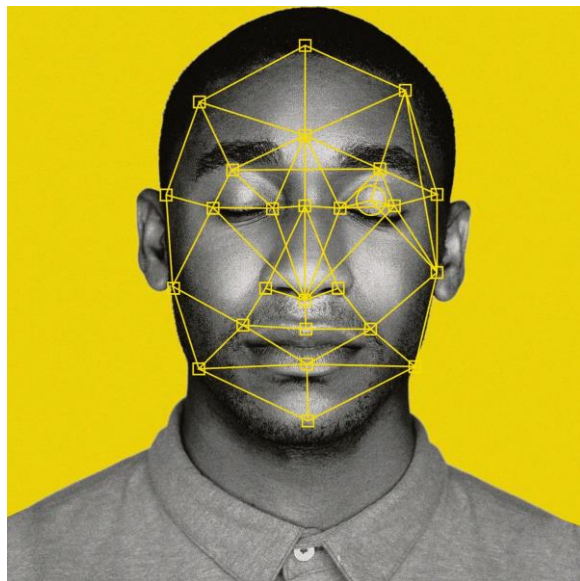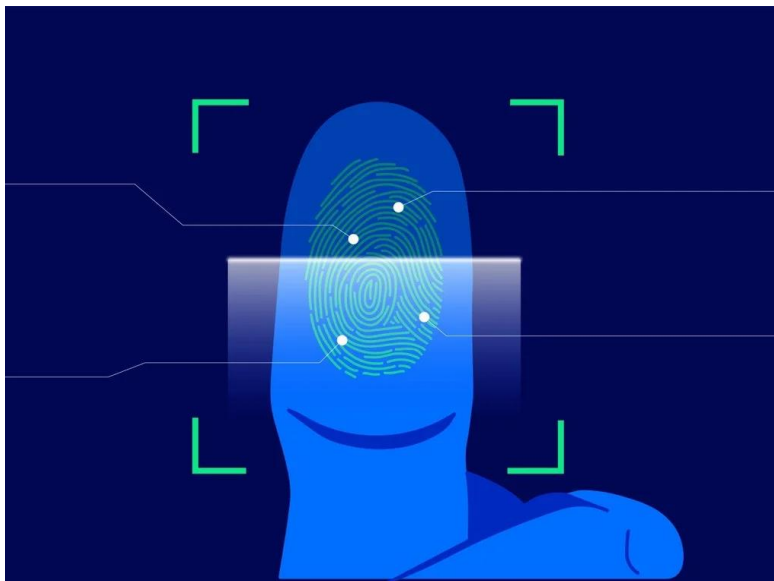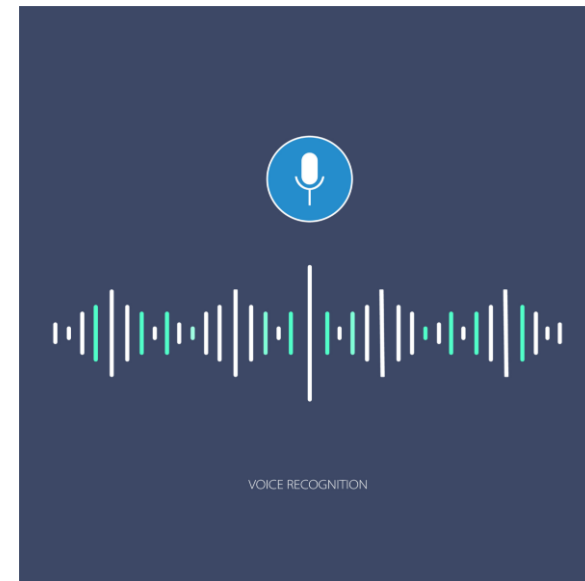
# Biometrics definition

Unique **physical** or **behavioral** human characteristics that can be used to digitally identify a person to grant access to systems, devices or data.
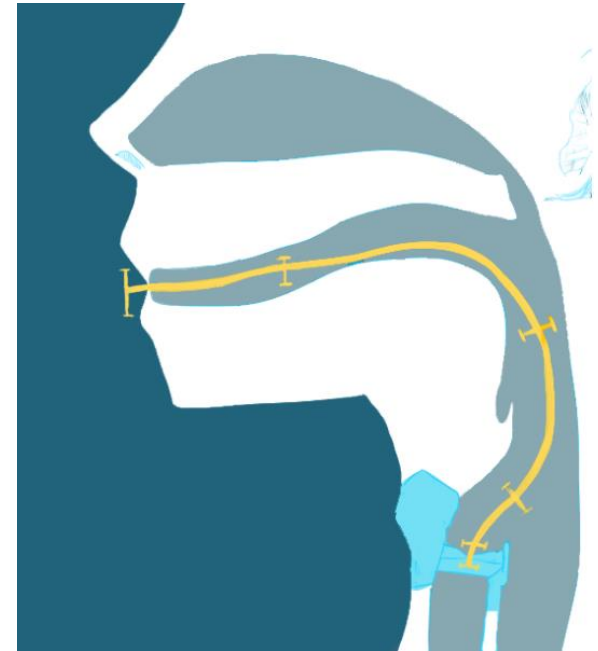
Face ID

Fingerprint
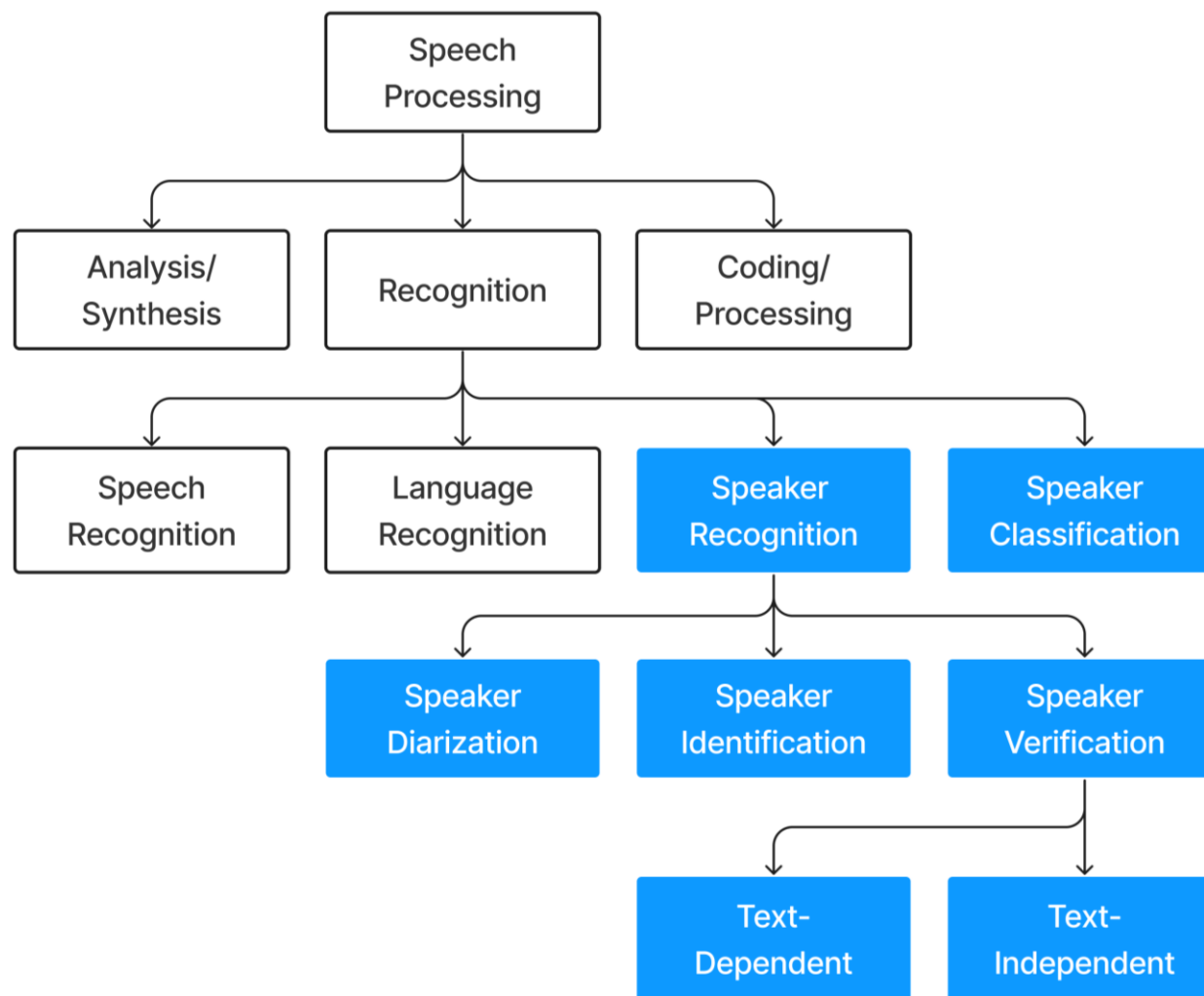
Voiceprint

# Voiceprint

Behavioral and Physical characteristic factors are combined to produce a unique voice pattern for each individual.

Biometric technology captures this pattern as a voiceprint.

- Speed of Speech

- Pronunciation and Emphasis

- Accents

- Unique Physical Traits
  of Vocal Tract

- Mouth Shape and Size
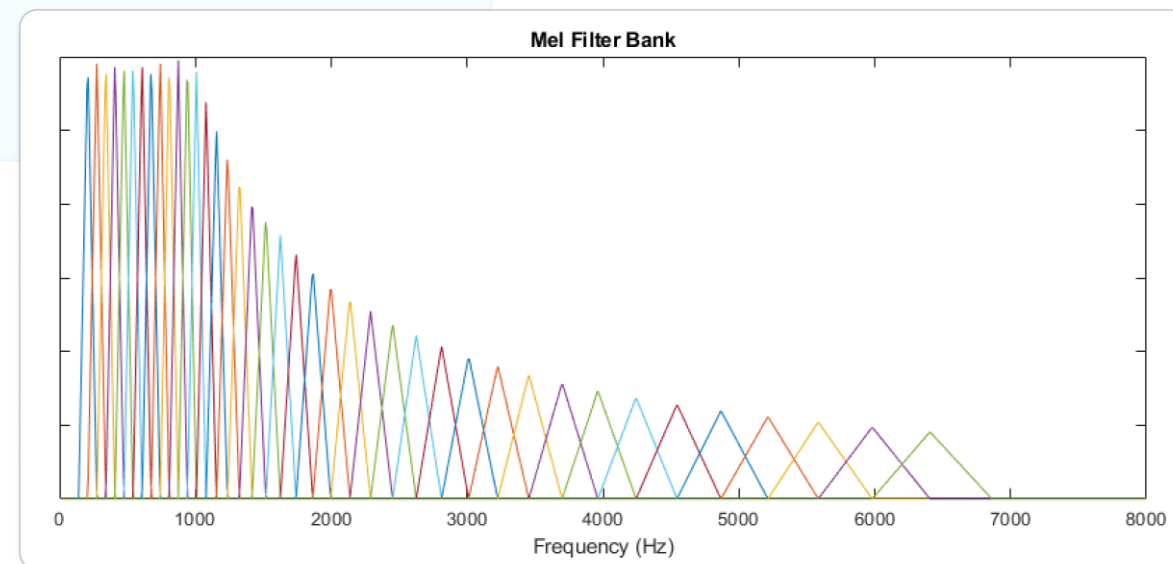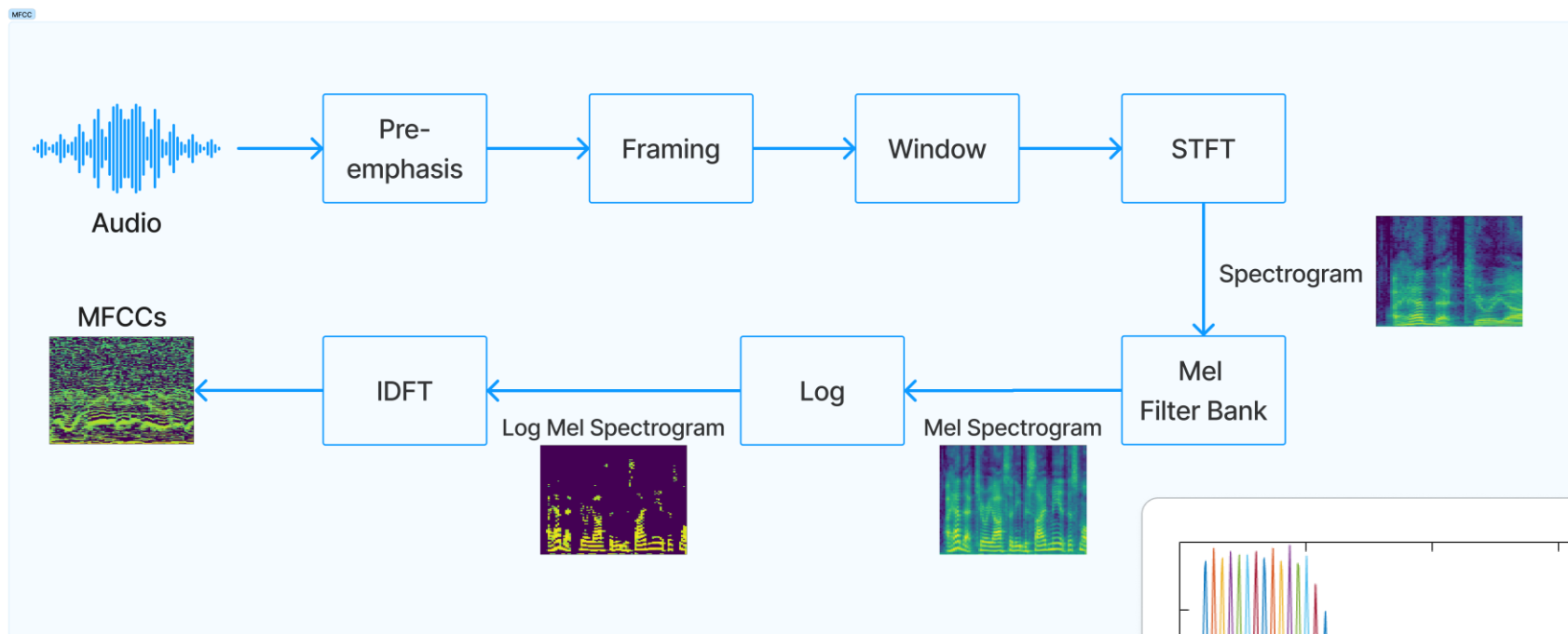
- Nasal Passages

# Voice Biometrics

# Speech preprocessing
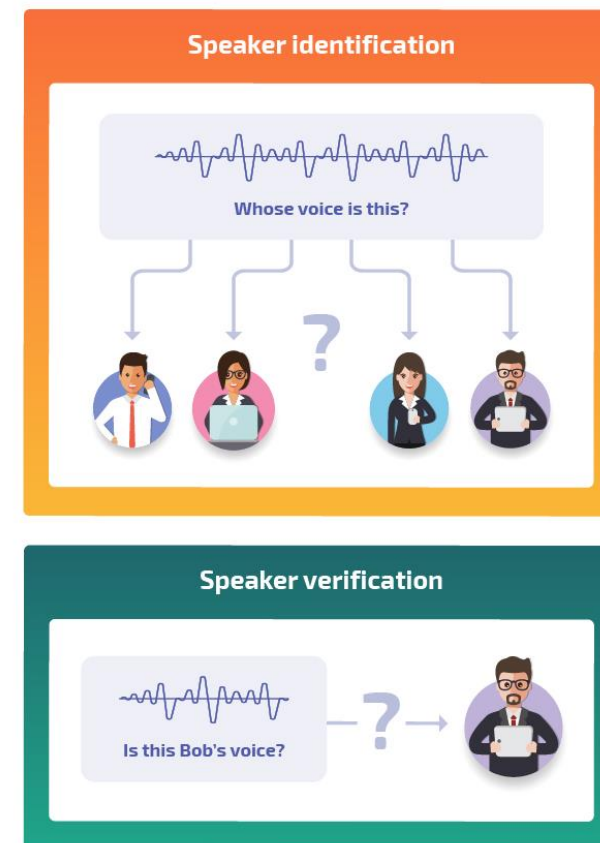
# Speaker Identification & Verification

**Speaker Verification:**

- The speaker claims to be of a certain identity
  and the voice is used to verify this claim.

- A 1:1 match where one speaker's voice is
  matched to a particular template.
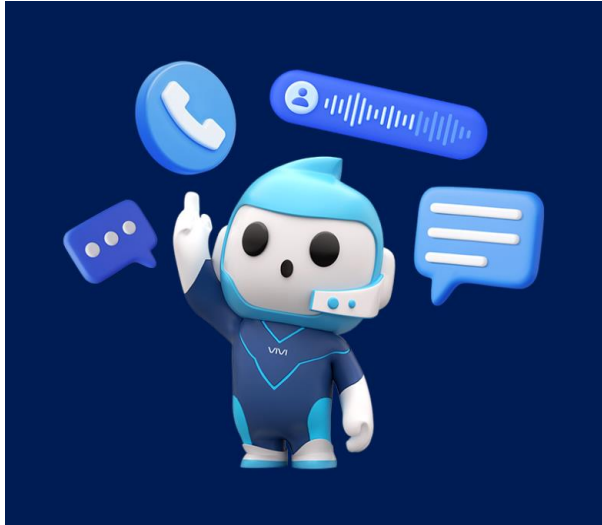
**Speaker Identification:**

- Determining an unknown speaker's identity.

- A 1:N match where the voice is compared
  against multiple templates.
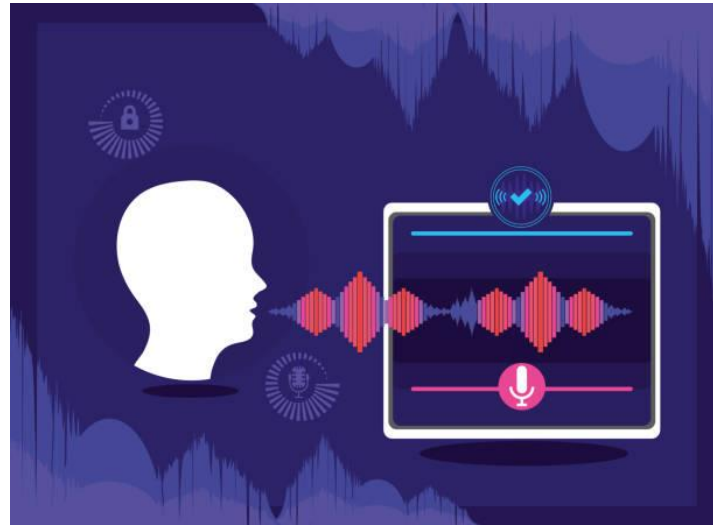


TYPES OF SPEAKER RECOGNITION

Speaker identification

Whose voice is this?

Speaker verification

Is this Bob's voice?

www.apriorit.com

# Applications



Call-bot
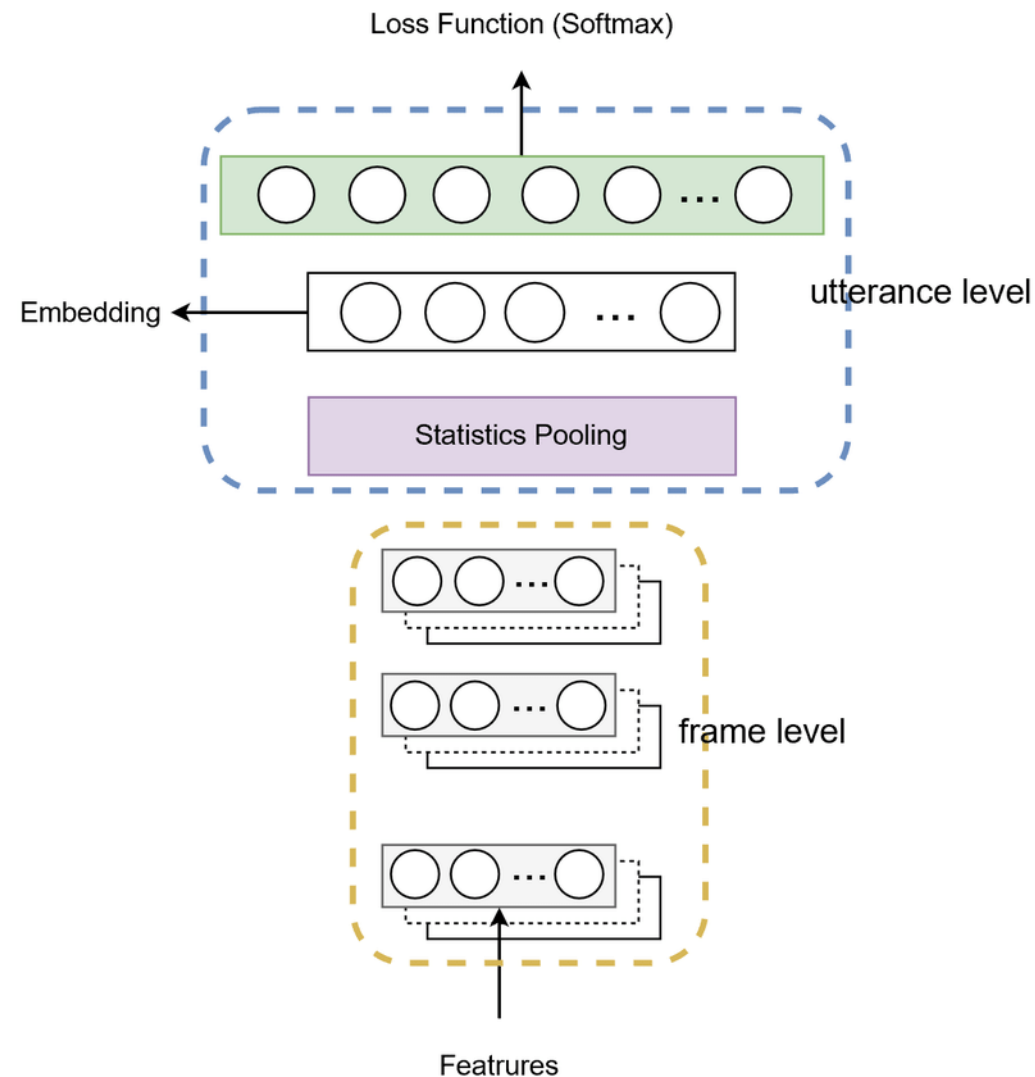


Authentication



Voice Assistant

# Speaker Embedding

**Model backbone:**
- ResNet
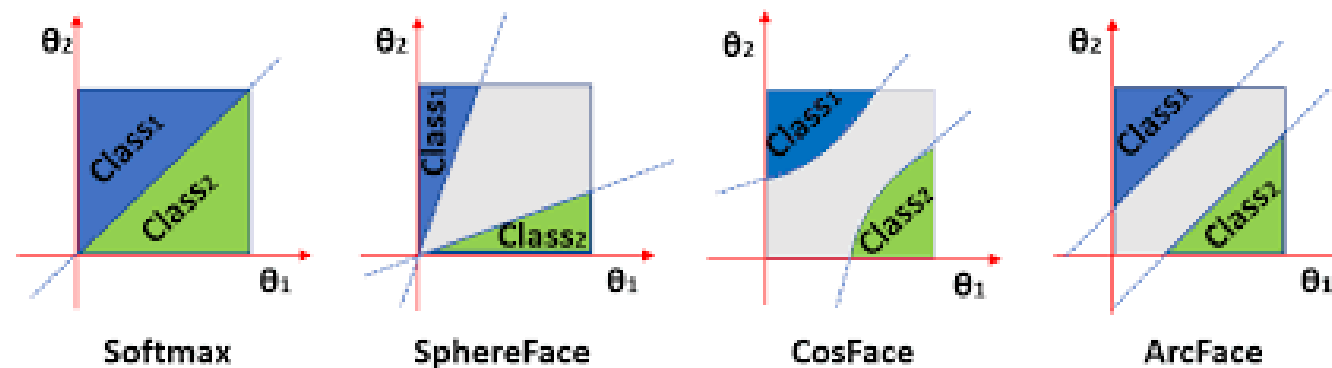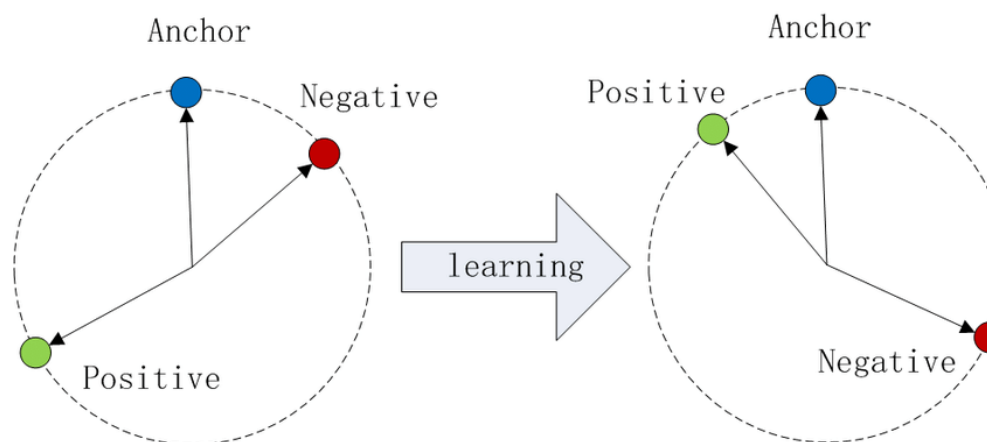- TDNN
- CNN-TDNN
- ECAPA-TDNN

**Pooling:**
- Statistics Pooling
- Attentive Statistics Pooling
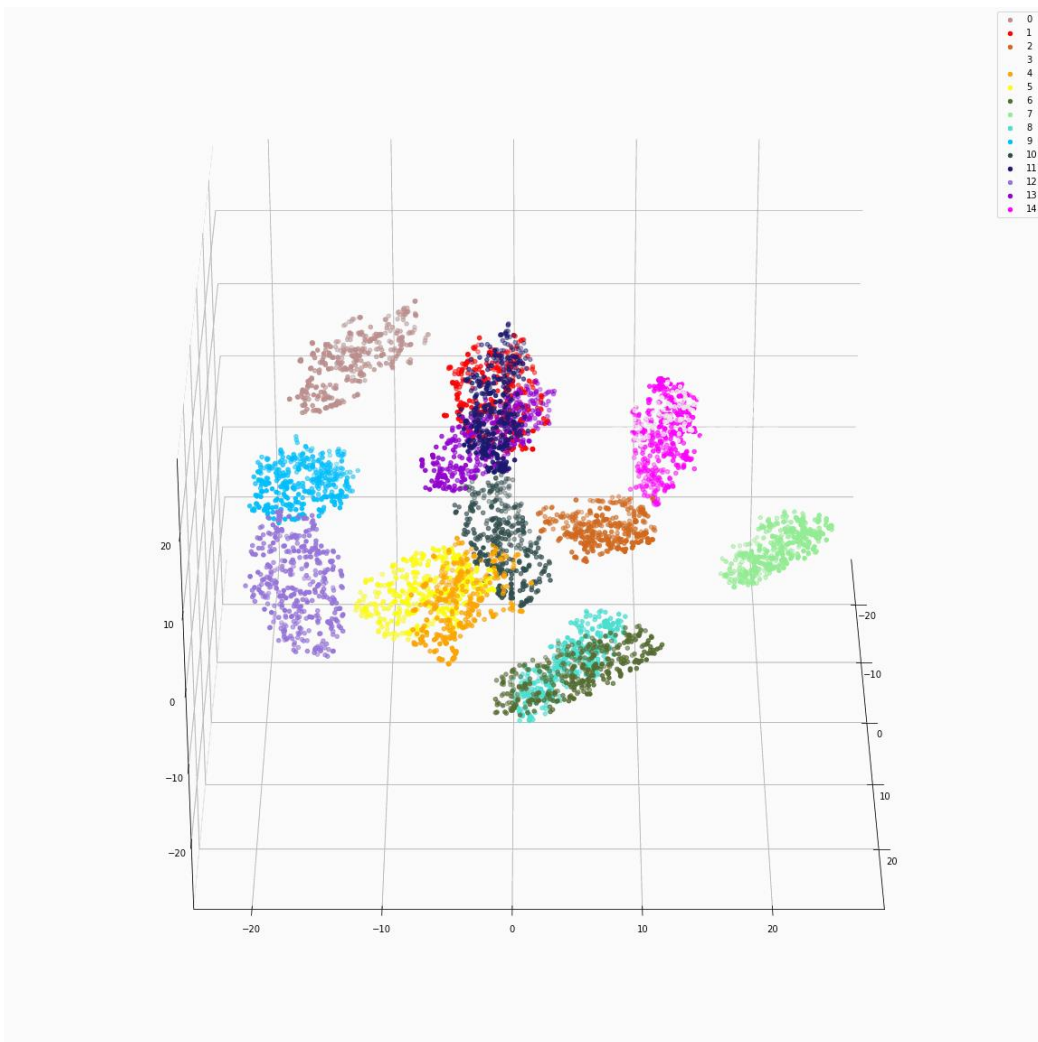- Multi-Head Attention Pooling

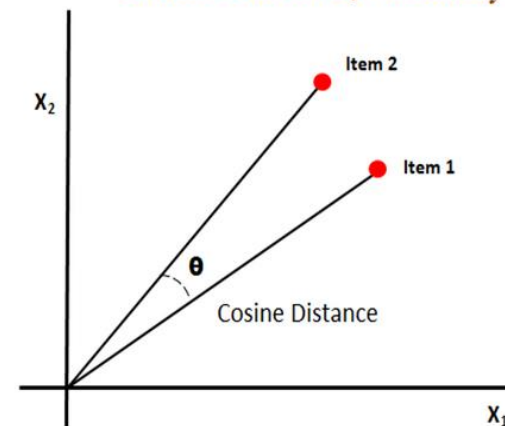# Speaker Embedding

**Loss Function:**
- Metric Loss:
  - Triplet Loss
  - Contrastive Loss
- Classification Loss:
  - Softmax Loss
  - A-Softmax Loss
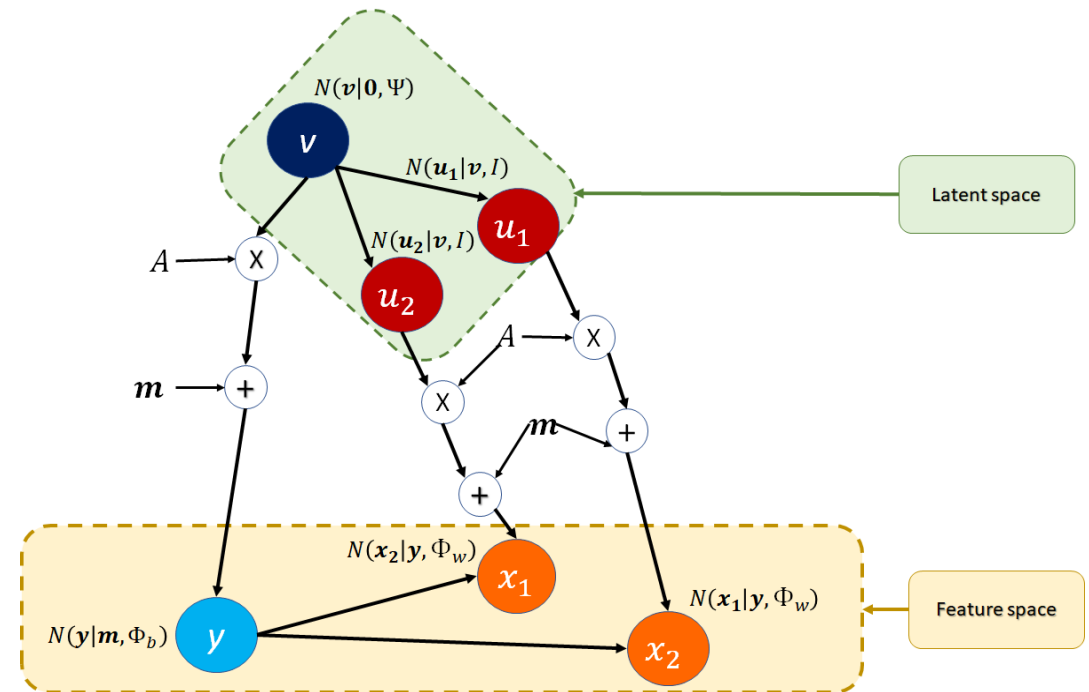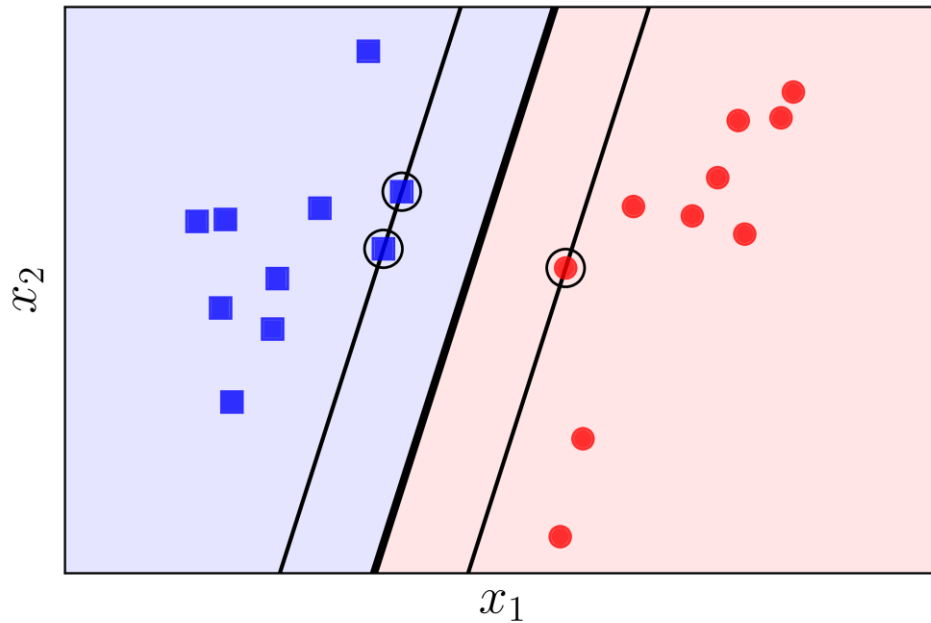  - AM-Softmax Loss
  - AAM-Softmax Loss

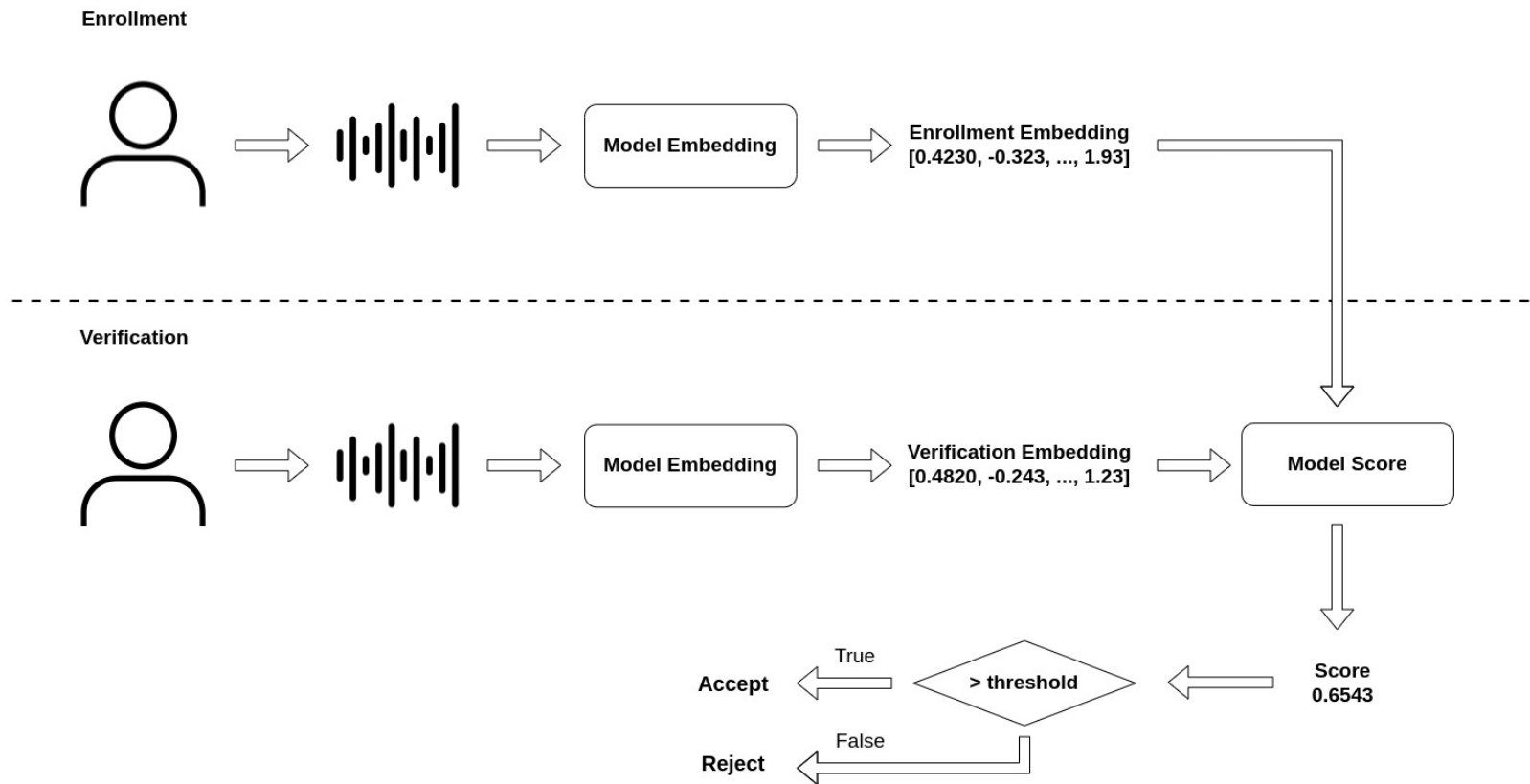# Speaker Embedding





*Cosine Distance/Similarity*

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\|\|\mathbf{B}\|} = \frac{\sum\limits_{i=1}^{n} A_i B_i}{\sqrt{\sum\limits_{i=1}^{n} A_i^2}\sqrt{\sum\limits_{i=1}^{n} B_i^2}},$$

# Back-end
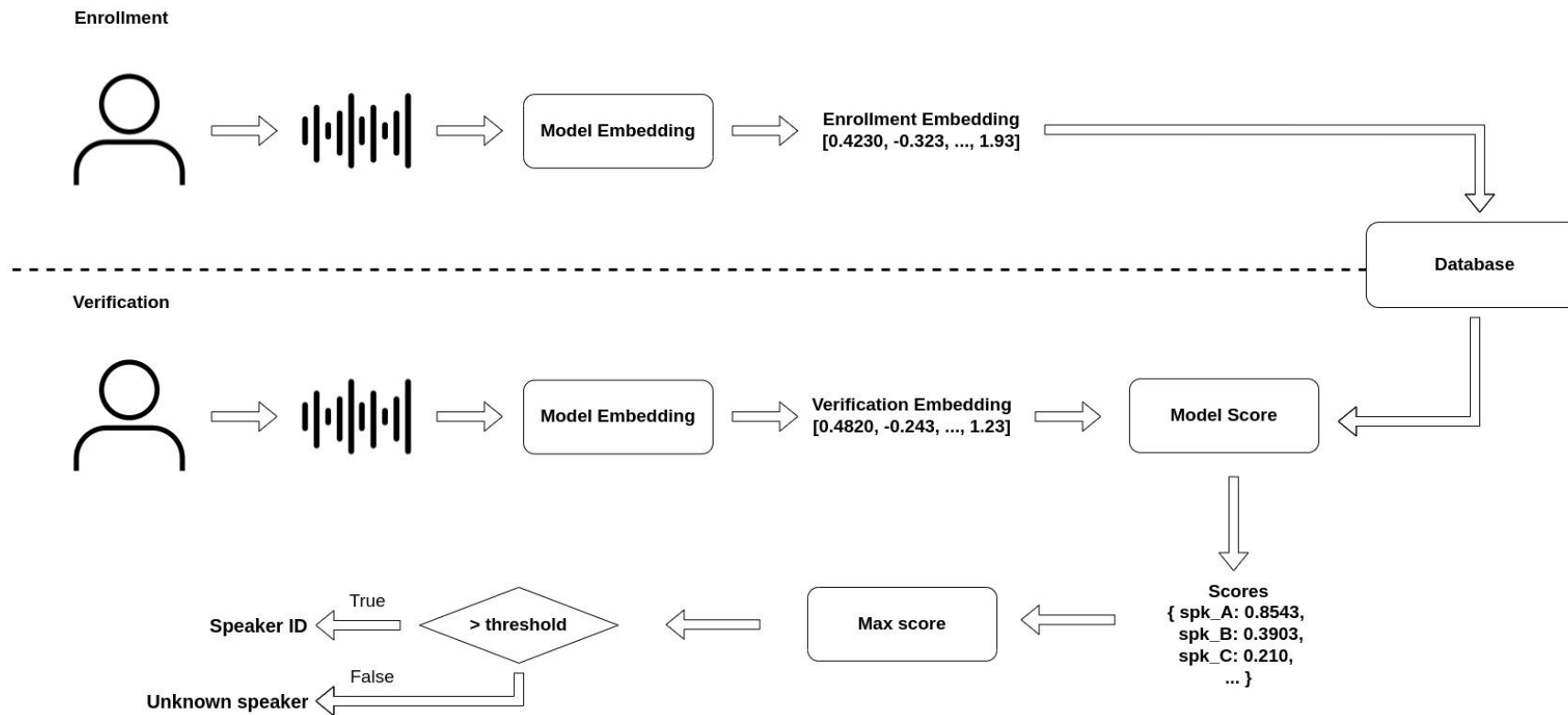
- Cosine Similarity

- Basic Classifiers: SVM, GMM, Logistic Regression (LR)

- PLDA Classifiers: PLDA, APLDA, CORAL, …

# Speaker Verification

Enrollment
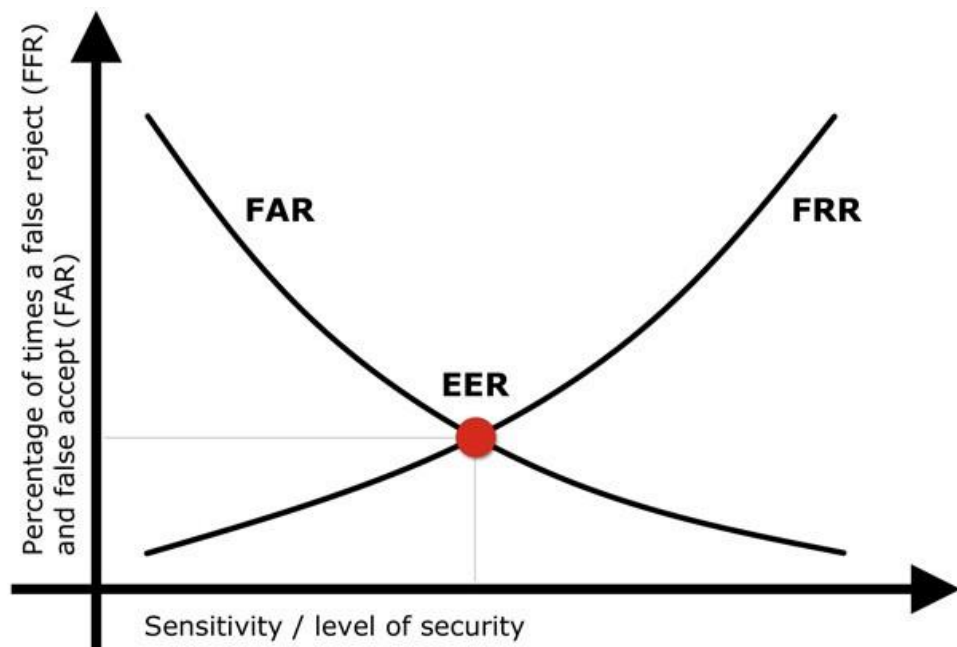
Enrollment Embedding
[0.4230, -0.323, ..., 1.93]

Verification

Model Embedding

Verification Embedding
[0.4820, -0.243, ..., 1.23]

Model Score

Model Embedding

Score
0.6543

> threshold

True → Accept

False → Reject

# Speaker Identification

# Metrics

Equal Error Rate & Minimum Detection Cost Function

## EER

## Min DCF

$$C_{\det}(P_{\mathrm{miss}}, P_{\mathrm{FA}}) = C_{\mathrm{miss}}P_{\mathrm{miss}}P_{\mathrm{tar}} + C_{\mathrm{FA}}P_{\mathrm{FA}}(1 - P_{\mathrm{tar}})$$

$C_{miss}$ - cost of a miss target (false reject)
$C_{FA}$ – cost of a false alarm (false accept)
$P_{miss}$ and $P_{FA}$ are determined by the evaluator by counting errors.
$P_{tar}$ is the prior probability that a target speaker event occurs in the application.
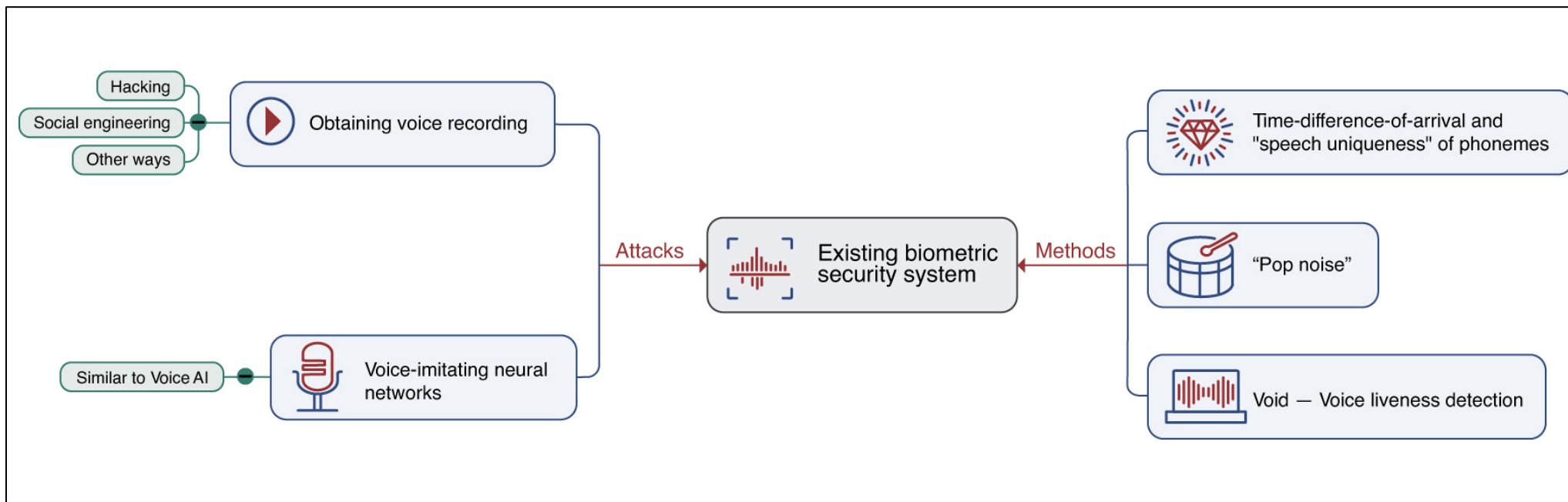
# Challenges



Data



Voice Quality



Spoofing



Cross Device

# Voice anti-spoofing



Logical access:
- TTS
- Voice Conversion

Physical access:
- Replay
- Voice Parody

# Speaker Diarization

**Speaker Diarization** is the task of segmenting and co-indexing audio recordings by speaker.
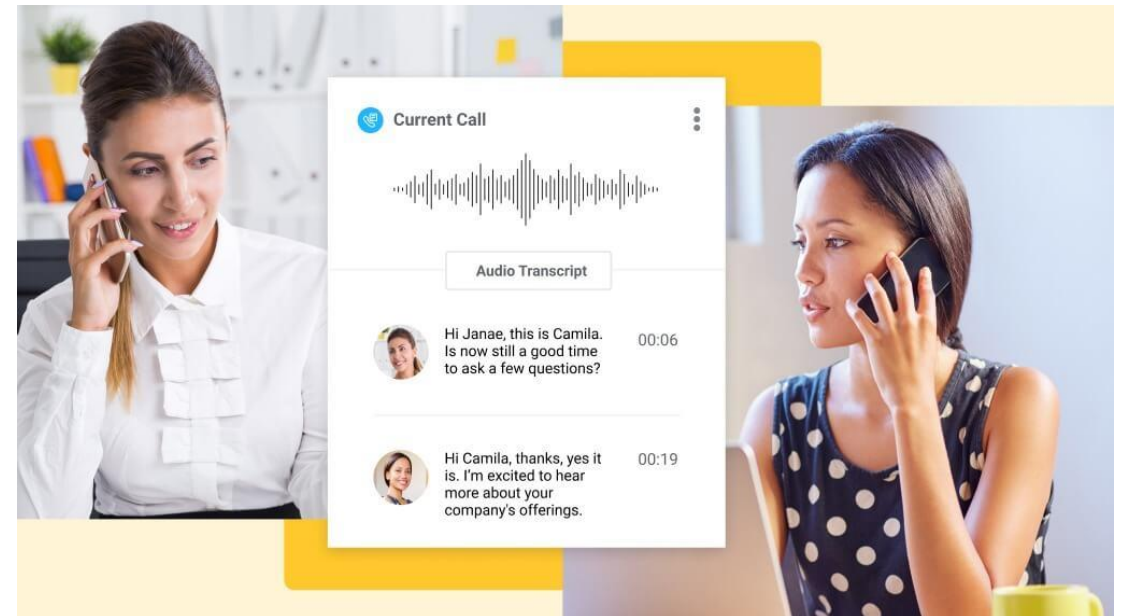
Diarization implies finding speaker boundaries and grouping segments that belong to the same speaker, and, as a by-product, determining the number of distinct speakers.
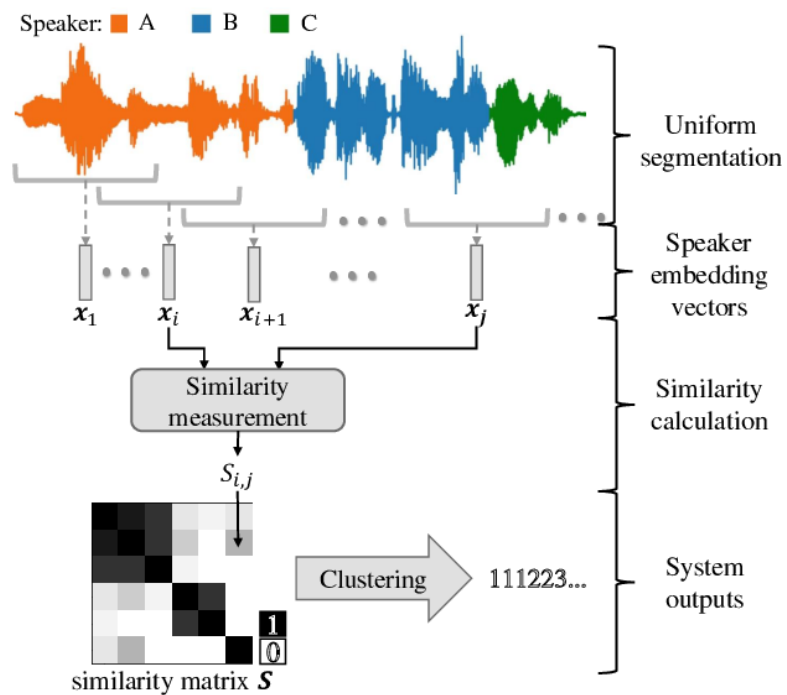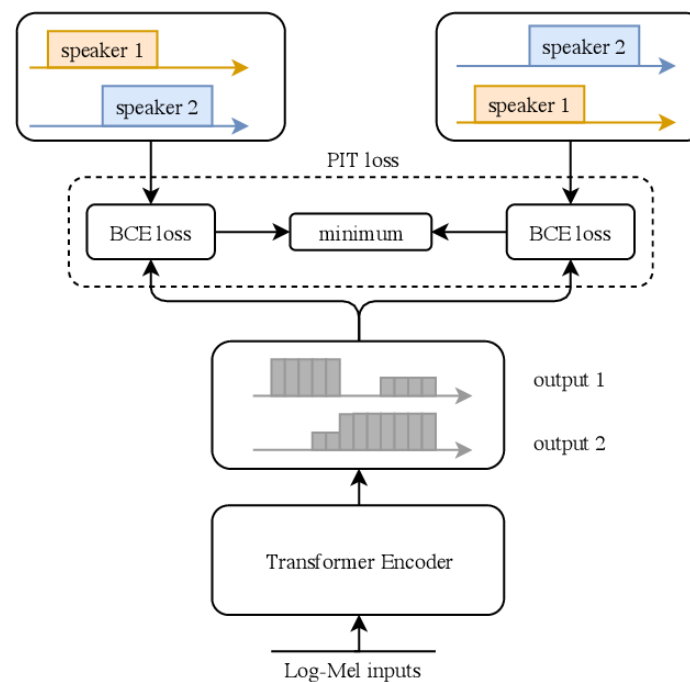
# Applications



Meeting transcription



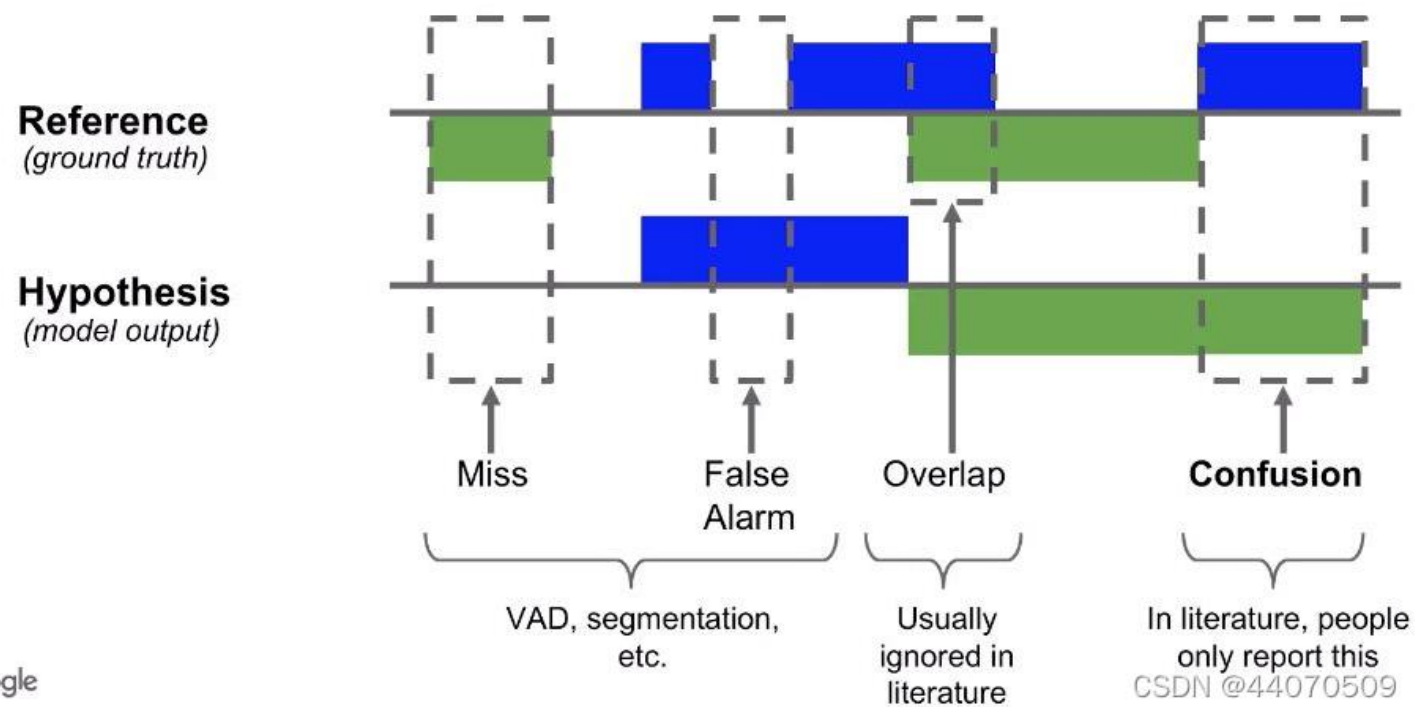Call transcription

# Speaker Diarization

**Clustering**



**End to end**

# Metrics

Diarization Error Rate:
$$\mathrm{DER} = \frac{T_{FA} + T_{MISS} + T_{SPKR}}{T_{SPEECH}}$$

# Speaker Classification
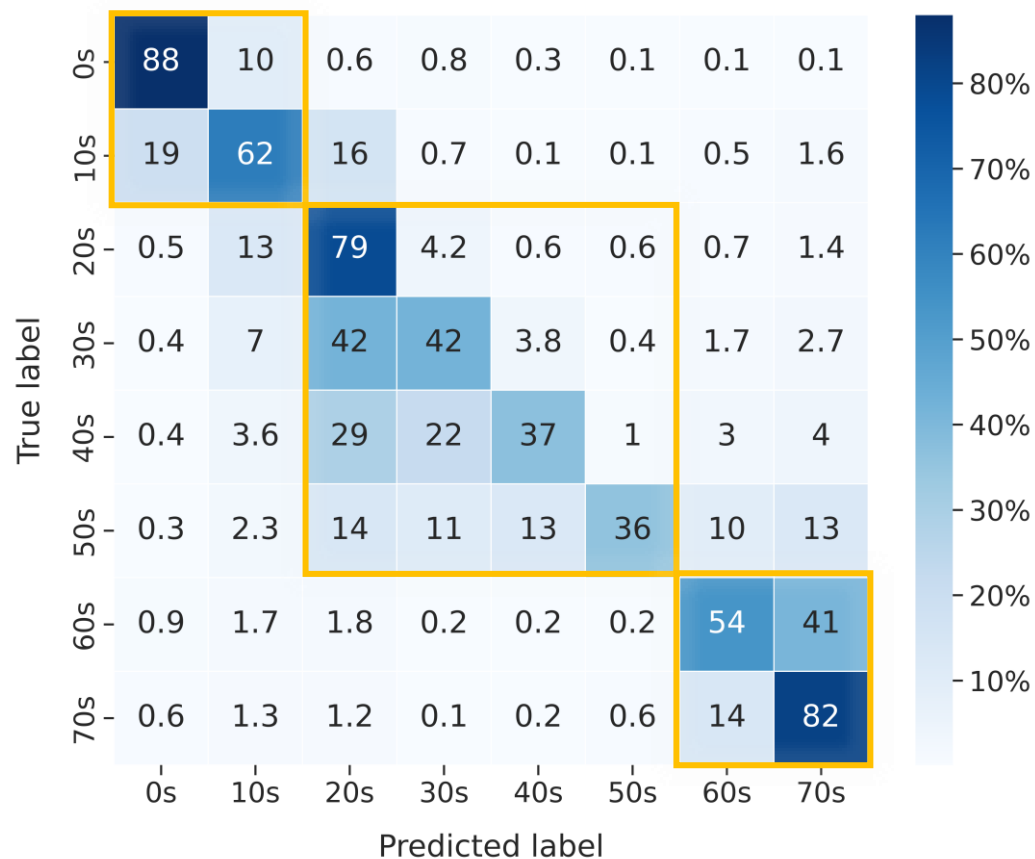
Classify a speaker by:

- Language, accent

- Age

- Gender

- Emotions

Customer benefit :

- Easy UI Language switch

- Gender/Age specific UI

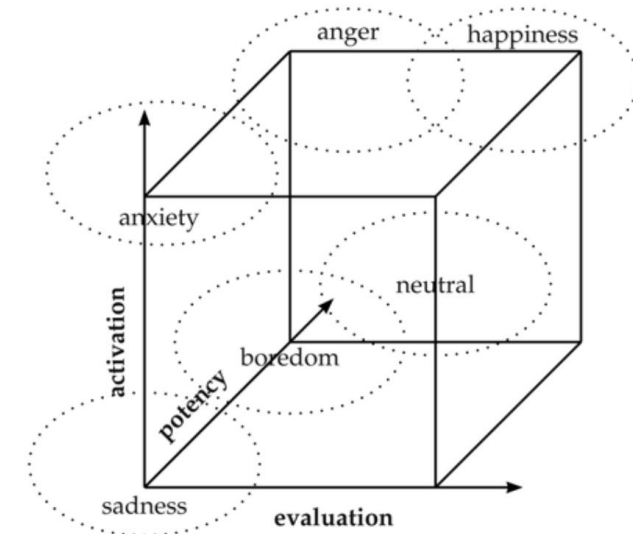- Enables parental control

# Age Estimation



Age range

- 0-20 years old
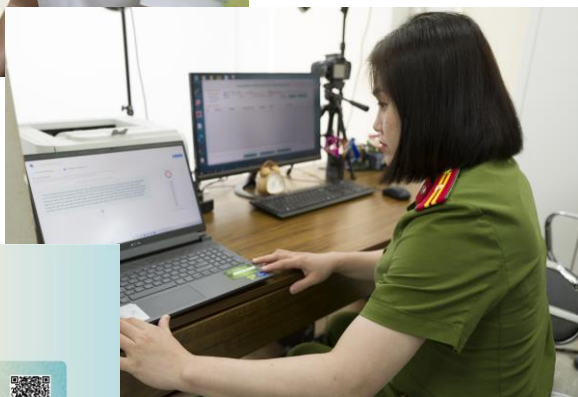- 20-60 years old
- > 60 years old

# Emotions



Emotions describe subjective feelings in short periods of time that are related to events, persons, or objects

Approaches:
- Categorical emotion approach
- Dimensional emotion approach

# Our Project

**VINBIGDATA**

# Thank you

Tầng 9, Century Tower, Times City, 458 Minh Khai, Phường
Vĩnh Tuy, Quận Hai Bà Trưng, Hà Nội.

info@vinbigdata.org

product.vinbigdata.org