# Robust speaker recognition using unsupervised adversarial invariance

**Raghuveer Peri, Monisankha Pal, Arindam Jati, Krishna Somandepalli, Shrikanth Narayanan**

**Presented by**

Raghuveer Peri
Signal Analysis and Interpretation Laboratory
University of Southern California

ICASSP2020
Barcelona

1

**Goal**

Extract robust, low-dimensional, speaker-discriminative representations ("*speaker embeddings*")
from speech signal

### Goal

Extract robust, low-dimensional, speaker-discriminative representations ("*speaker embeddings*")
from speech signal

### Applications

- Automatic Speaker Verification (ASV): Verify identity of person from speech signal
- Speaker diarization: Determine who spoke when in multi-party conversations
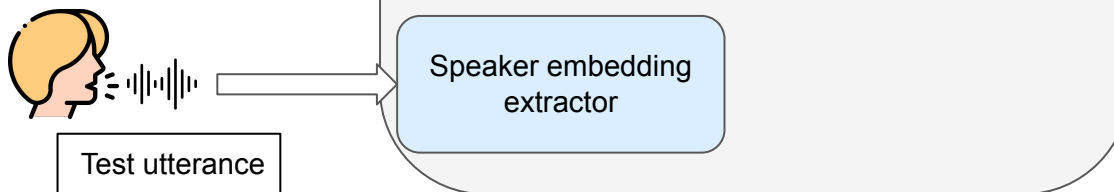- Automatic Speech Recognition: Speaker-adapted speech recognition models

## Goal

Extract robust, low-dimensional, speaker-discriminative representations ("*speaker embeddings*") from speech signal

## Applications

- Automatic Speaker Verification (ASV): Verify identity of person from speech signal
- Speaker diarization: Determine who spoke when in multi-party conversations
- Automatic Speech Recognition: Speaker-adapted speech recognition models

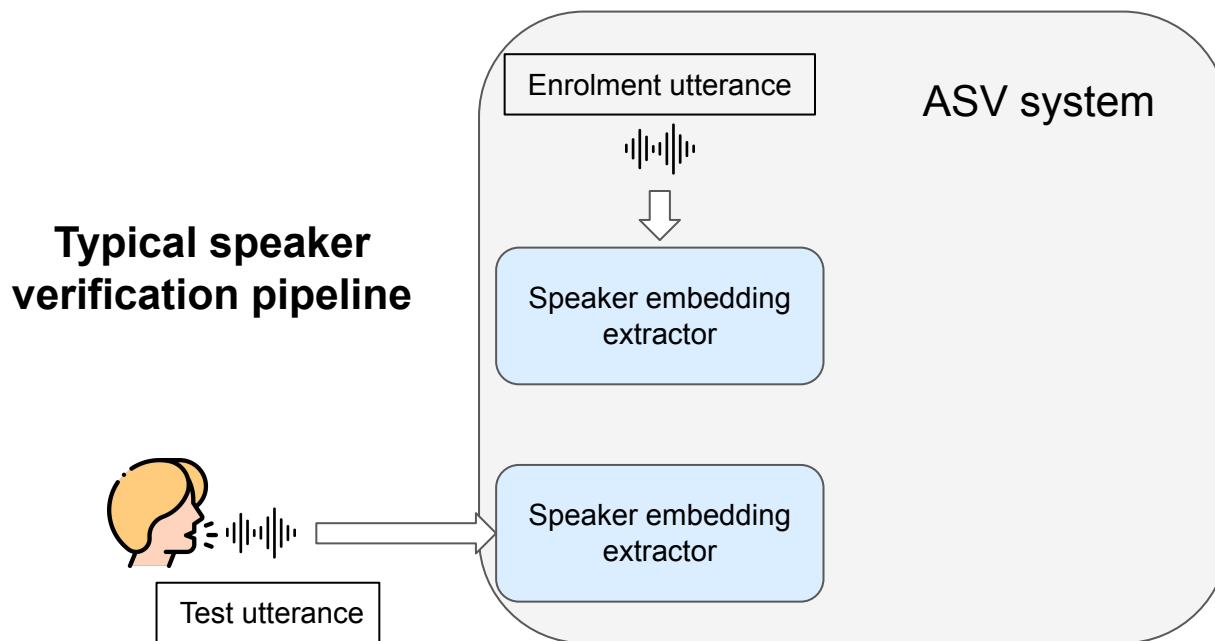**Typical speaker verification pipeline**

ASV system

Speaker embedding extractor

Test utterance

4

## Goal

Extract robust, low-dimensional, speaker-discriminative representations ("*speaker embeddings*") from speech signal

## Applications

- Automatic Speaker Verification (ASV): Verify identity of person from speech signal
- Speaker diarization: Determine who spoke when in multi-party conversations
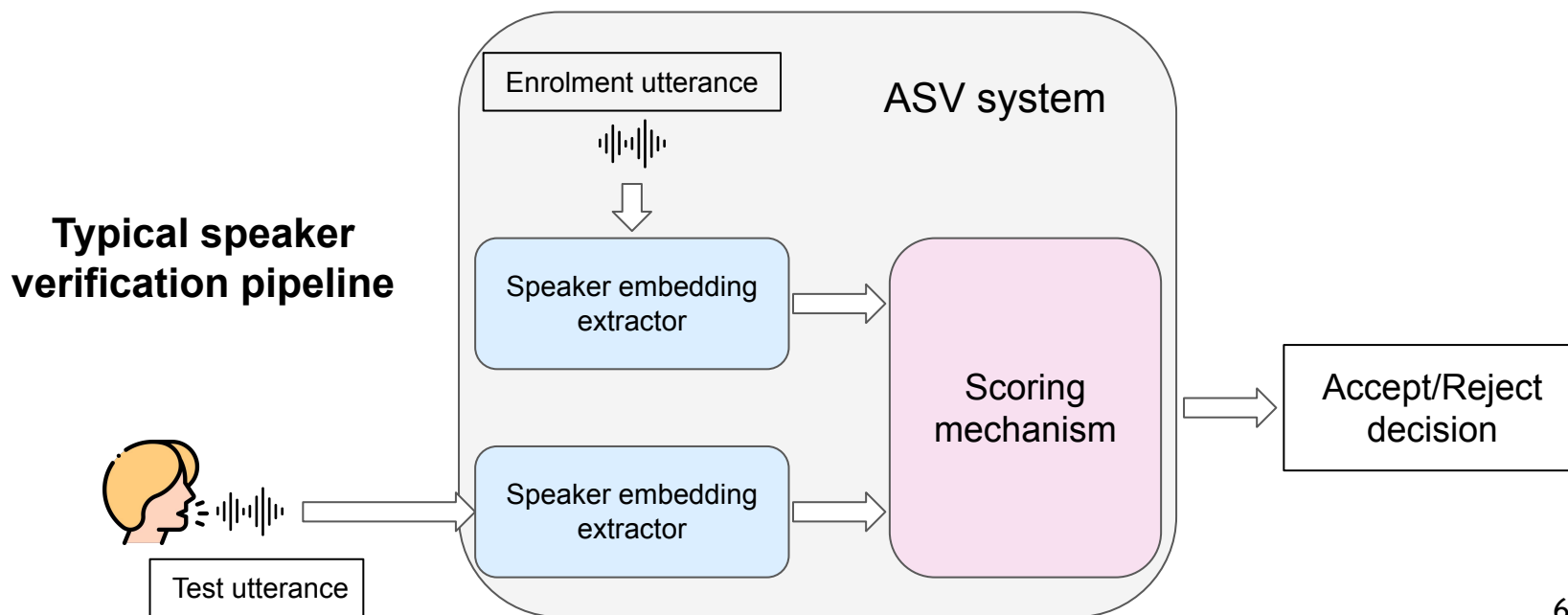- Automatic Speech Recognition: Speaker-adapted speech recognition models

**Typical speaker verification pipeline**

ASV system

Enrolment utterance

Speaker embedding extractor

Speaker embedding extractor

Test utterance

## Goal

Extract robust, low-dimensional, speaker-discriminative representations ("*speaker embeddings*") from speech signal
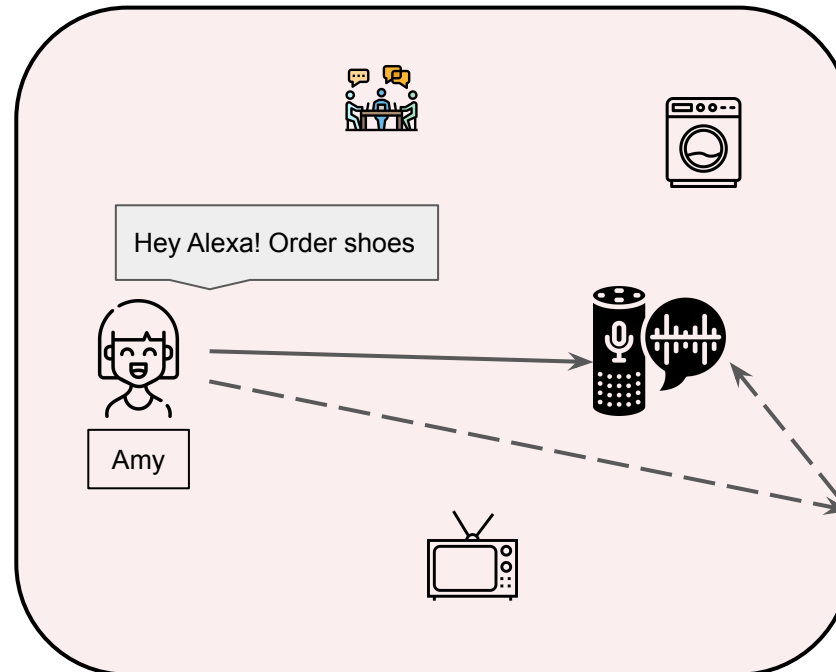
## Applications

- Automatic Speaker Verification (ASV): Verify identity of person from speech signal
- Speaker diarization: Determine who spoke when in multi-party conversations
- Automatic Speech Recognition: Speaker-adapted speech recognition models

**Typical speaker verification pipeline**

Enrolment utterance

ASV system

Speaker embedding extractor

Scoring mechanism

Accept/Reject decision

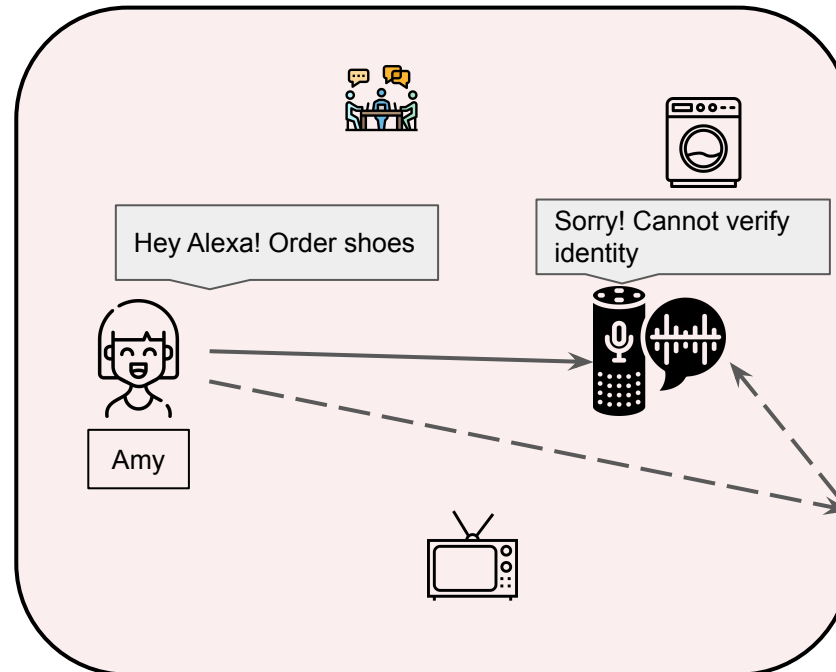Speaker embedding extractor

Test utterance

## Challenges

- Speech is an information-rich signal

- *Nuisance factors* unrelated to speaker identity entangled in signal
  - Channel factors
    - Acoustic noise (TV, babble etc.)
    - Reverberation
  - Content factors
    - Affective state (happy, angry etc.)
    - Linguistic content
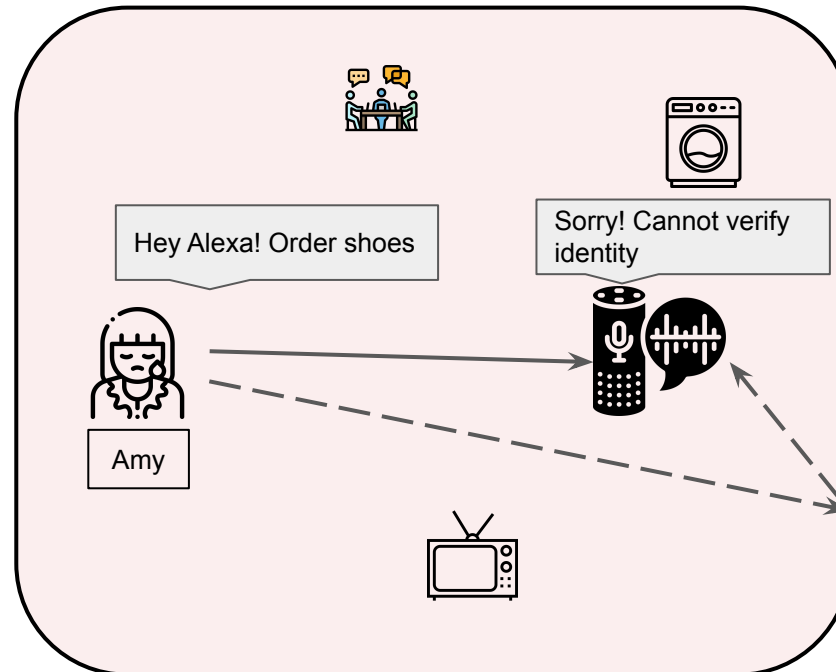
**USC** Viterbi
School of Engineering

## Challenges

- Speech is an information-rich signal

- *Nuisance factors* unrelated to speaker identity entangled in signal
  - Channel factors
    - Acoustic noise (TV, babble etc.)
    - Reverberation
  - Content factors
    - Affective state (happy, angry etc.)
    - Linguistic content

Hey Alexa! Order shoes

Sorry! Cannot verify identity

Amy

## Challenges

- Speech is an information-rich signal

- *Nuisance factors* unrelated to speaker identity entangled in signal
  - Channel factors
    - Acoustic noise (TV, babble etc.)
    - Reverberation
  - Content factors
    - Affective state (happy, angry etc.)
    - Linguistic content

Hey Alexa! Order shoes

Sorry! Cannot verify identity

Amy

**Prior work**

- Total Variability Modeling (i-vectors - *Dehak et al., 2011*)
  - Capture all factors of variability in total variability space
  - Perform additional channel compensation steps, such as length normalization

- Deep learning methods (x-vectors - *Snyder et al., 2017*)
  - Train deep models on artificially augmented audio using various noise and reverberation.
  - Extract hidden layer representations as utterance-level features.

- More recent supervised domain adversarial training techniques (*Bhattacharya et al., 2019*)
  - Train models to discriminate speakers
  - Simultaneously made robust to "specific" factors of variability by training adversarially, such as known noise type or channel conditions.
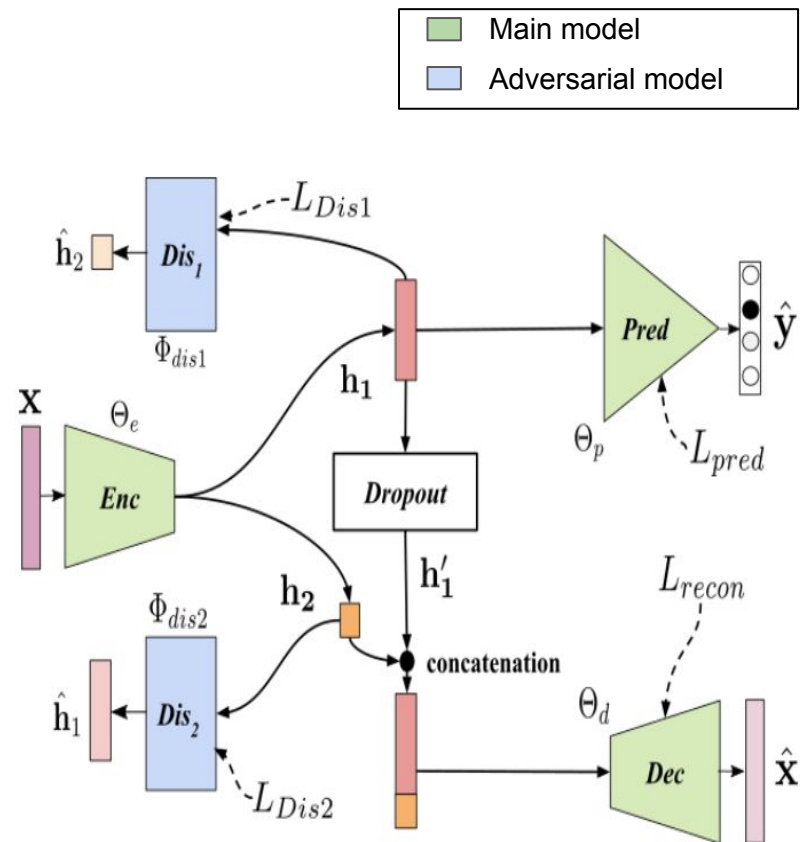
## Prior work

- Total Variability Modeling (i-vectors - Dehak et al., 2011)
  - Capture all factors of variability in total variability space
  - Perform additional channel compensation steps, such as length normalization

- Deep learning methods (x-vectors - Snyder et al., 2017)
  - Train deep models on artificially augmented audio using various noise and reverberation.
  - Extract hidden layer representations as utterance-level features.

- More recent supervised domain adversarial training techniques (Bhattacharya et al., 2019)
  - Train models to discriminate speakers
  - Simultaneously made robust to "specific" factors of variability by training adversarially, such as known noise type or channel conditions.

## Proposed work

- Disentangle speech representations into two embeddings
  - Speaker factors
  - Nuisance factors

- **No assumptions on specific factors of variability**

## Input

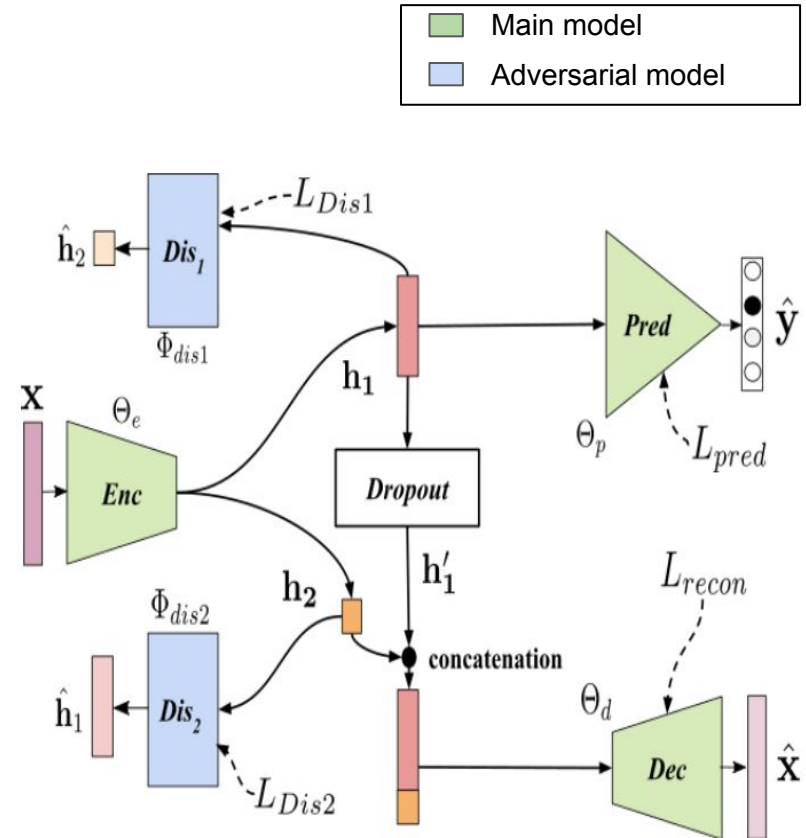- Speech representations (MFCC, x-vectors etc)
- Speaker labels



Legend:
- Main model
- Adversarial model

### Input

- Speech representations (MFCC, x-vectors etc)
- Speaker labels

### Main model

- Predictor (*Pred*): Predicts speakers

- Decoder (*Dec*): Reconstruct input
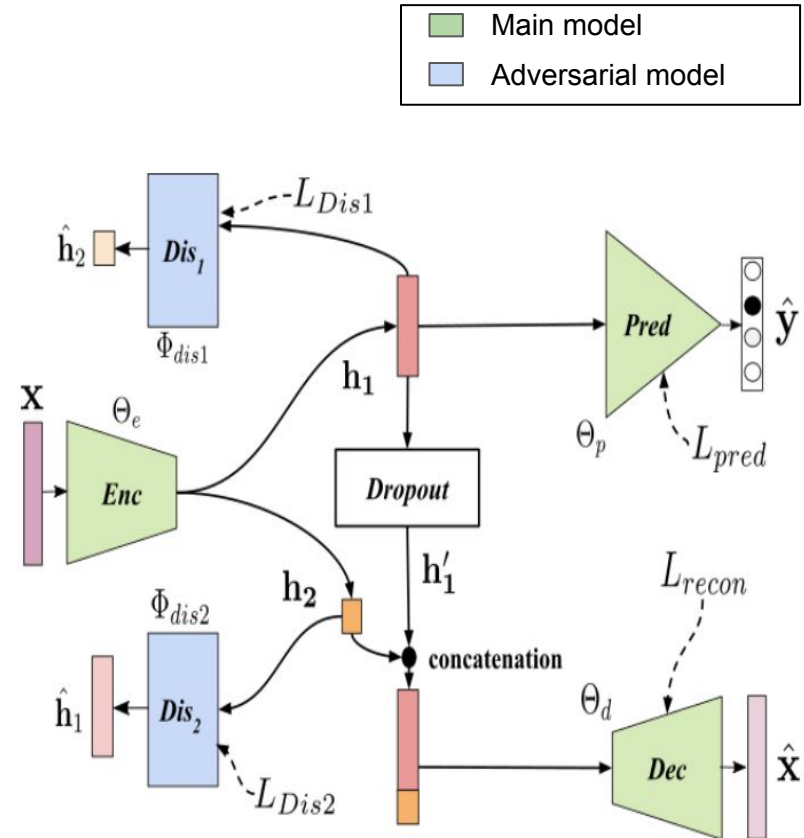


Main model
Adversarial model

## Input

- Speech representations (MFCC, x-vectors etc)
- Speaker labels

## Main model

- Predictor (*Pred*): Predicts speakers

- Decoder (*Dec*): Reconstruct input

## Adversarial model

- Disentanglers (*Dis$_1$* and *Dis$_2$*): Make h$_1$ and h$_2$ poor predictors of each other



14

## Input

- Speech representations (MFCC, x-vectors etc)
- Speaker labels

## Main model

- Predictor (*Pred*): Predicts speakers

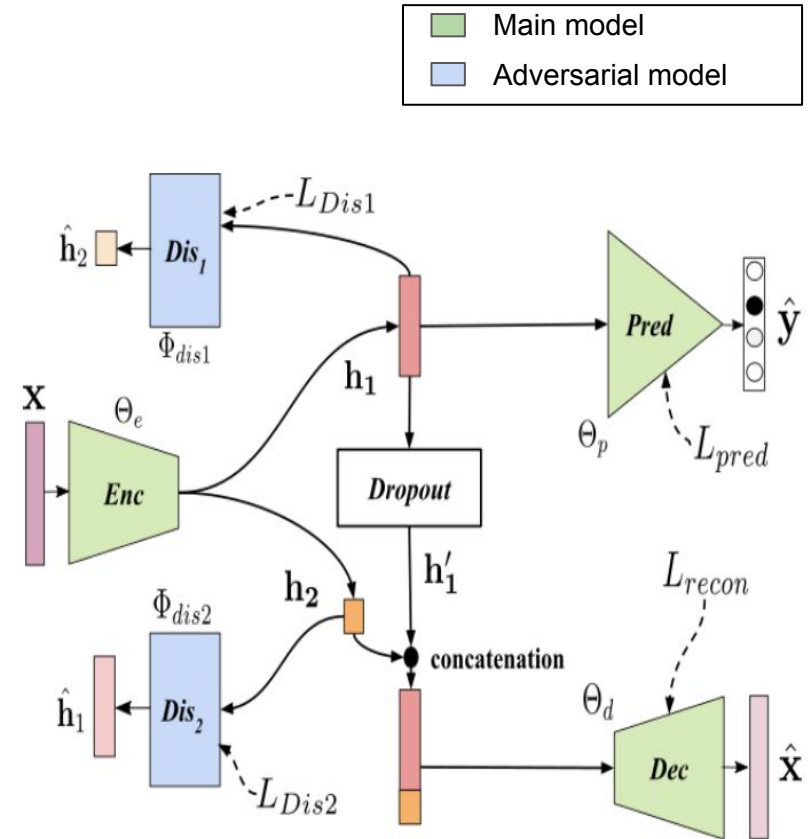- Decoder (*Dec*): Reconstruct input

## Adversarial model

- Disentanglers (*Dis$_1$* and *Dis$_2$*): Make h$_1$ and h$_2$ poor predictors of each other

**Legend:**
- Main model
- Adversarial model



## Adversarial Training[*]

$$L_{main} = \alpha L_{pred}\left(\mathbf{y}, \hat{\mathbf{y}}\right) + \beta L_{recon}\left(\mathbf{x}, \hat{\mathbf{x}}\right)$$

$$L_{adv} = L_{Dis1}(\mathbf{h}_2, \hat{\mathbf{h}}_2) + L_{Dis2}(\mathbf{h}_1, \hat{\mathbf{h}}_1)$$

$$\min_{\Theta_e, \Theta_d, \Theta_p} \max_{\Phi_{dis1}, \Phi_{dis2}} L_{main} + \gamma L_{adv}$$

*Jaiswal, A., Wu, R.Y., Abd-Almageed, W. and Natarajan, P., 2018. Unsupervised adversarial invariance. In Advances in Neural Information Processing Systems (pp. 5092-5102).

### Input

- Speech representations (MFCC, x-vectors etc)
- Speaker labels

### Main model

- Predictor (*Pred*): Predicts speakers

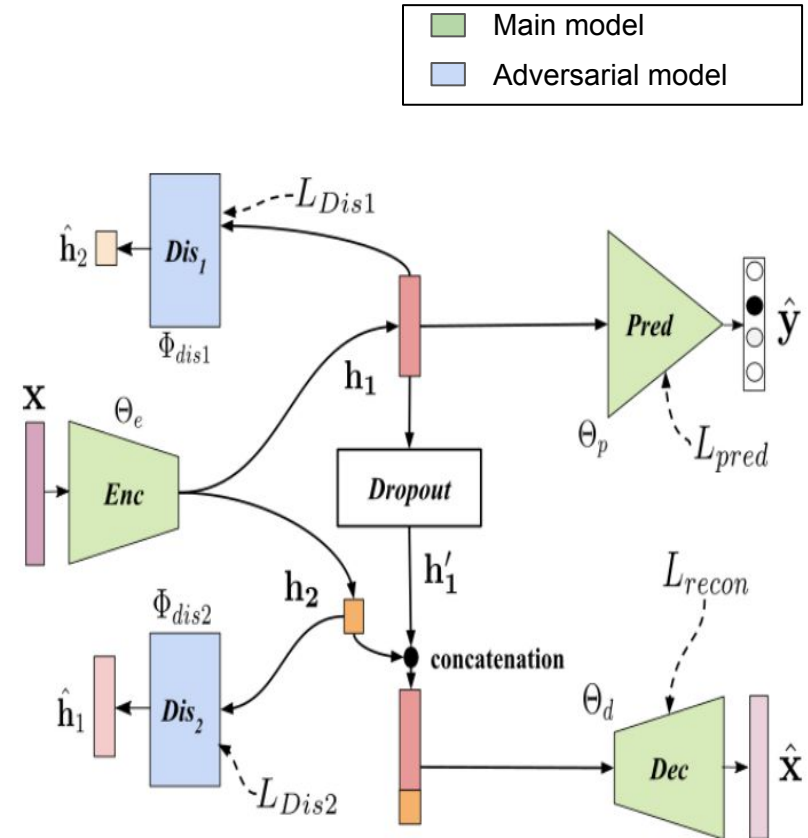- Decoder (*Dec*): Reconstruct input

### Adversarial model

- Disentanglers (*Dis₁* and *Dis₂*): Make $h_1$ and $h_2$ poor predictors of each other

**Legend:**
- Main model
- Adversarial model



### Adversarial Training[*]

$$L_{main} = \alpha L_{pred}\left(\mathbf{y}, \hat{\mathbf{y}}\right) + \beta L_{recon}\left(\mathbf{x}, \hat{\mathbf{x}}\right)$$

$$L_{adv} = L_{Dis1}\left(\mathbf{h}_2, \hat{\mathbf{h}}_2\right) + L_{Dis2}\left(\mathbf{h}_1, \hat{\mathbf{h}}_1\right)$$

$$\min_{\Theta_e, \Theta_d, \Theta_p} \max_{\Phi_{dis1}, \Phi_{dis2}} L_{main} + \gamma L_{adv}$$

16

*Jaiswal, A., Wu, R.Y., Abd-Almageed, W. and Natarajan, P., 2018. Unsupervised adversarial invariance. In Advances in Neural Information Processing Systems (pp. 5092-5102).

**USC** Viterbi
School of Engineering

SAiL

## Input

- Speech representations (MFCC, x-vectors etc)
- Speaker labels

## Main model

- Predictor (*Pred*): Predicts speakers

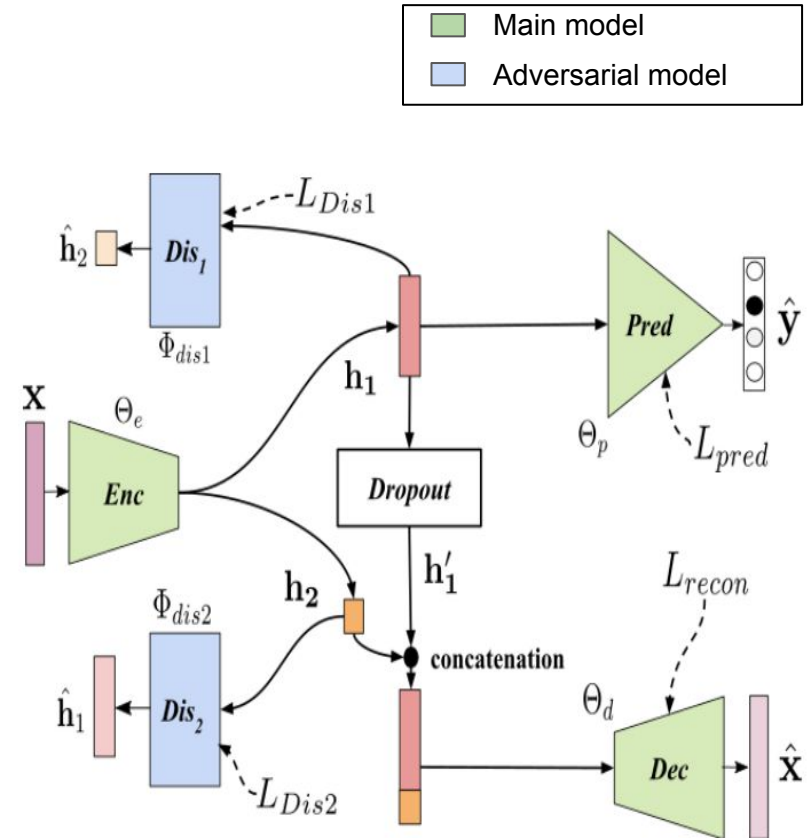- Decoder (*Dec*): Reconstruct input

## Adversarial model

- Disentanglers (*Dis$_1$* and *Dis$_2$*): Make h$_1$ and h$_2$ poor predictors of each other

**Legend:**
- Main model
- Adversarial model



## Adversarial Training*

$$L_{main} = \alpha L_{pred}(\mathbf{y}, \hat{\mathbf{y}}) + \beta L_{recon}(\mathbf{x}, \hat{\mathbf{x}})$$

$$L_{adv} = L_{Dis1}(\mathbf{h}_2, \hat{\mathbf{h}}_2) + L_{Dis2}(\mathbf{h}_1, \hat{\mathbf{h}}_1)$$

$$\min_{\Theta_e, \Theta_d, \Theta_p} \max_{\Phi_{dis1}, \Phi_{dis2}} L_{main} + \gamma L_{adv}$$

- h$_1$ : speaker discriminative information

- h$_2$ : nuisance information

17

*Jaiswal, A., Wu, R.Y., Abd-Almageed, W. and Natarajan, P., 2018. Unsupervised adversarial invariance. In Advances in Neural Information Processing Systems (pp. 5092-5102).
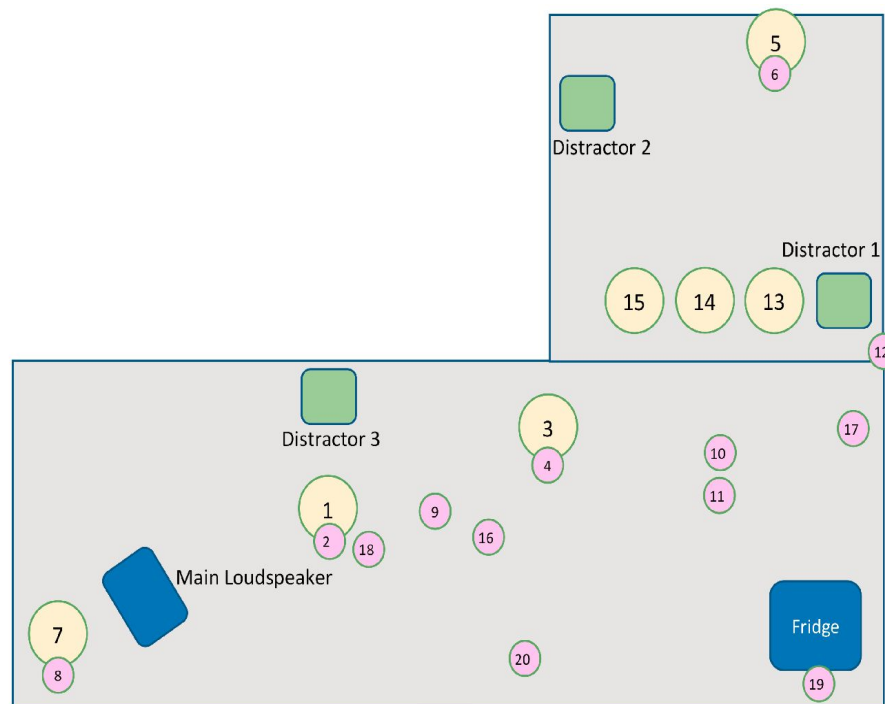
## Training data (VoxCeleb[1])

- Training set of VoxCeleb
  - Vox 1 (Dev)
  - Vox2 (Dev and test)
- No artificial augmentation
- 1.2M data samples
- 7323 unique speakers

## Input features

- x-vectors using pre-trained model[2]

## Training data (VoxCeleb[1])

- Training set of VoxCeleb
  - Vox 1 (Dev)
  - Vox2 (Dev and test)
- No artificial augmentation
- 1.2M data samples
- 7323 unique speakers

## Input features

- x-vectors using pre-trained model[2]

## Evaluation data (VOiCES[3])
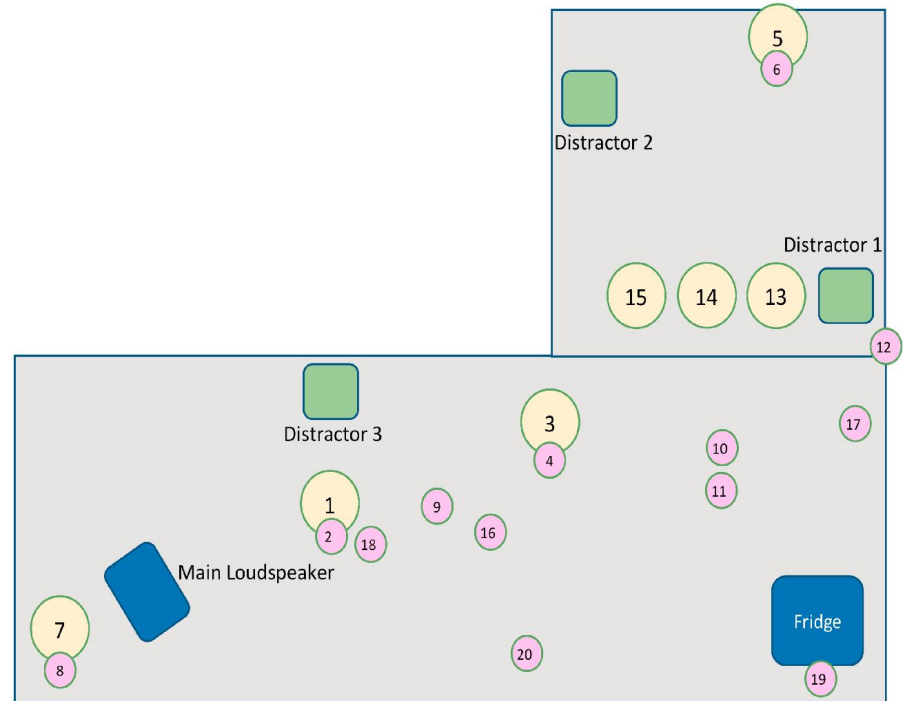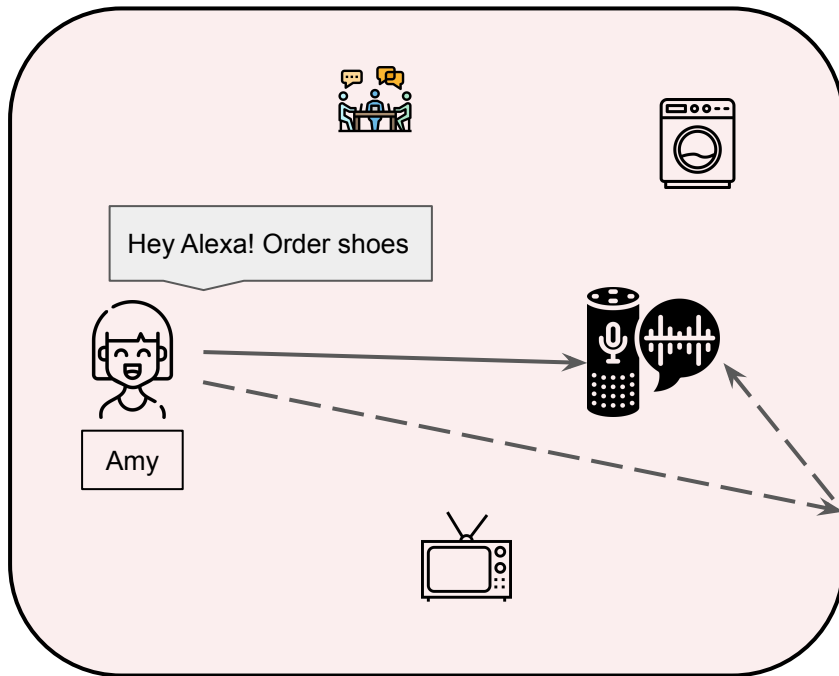
- 4 rooms
- 12-18 microphones
- 200 unique speakers
- 2 subsets: Voices-dev, Voices-eval

1. Chung, J.S., Nagrani, A. and Zisserman, A., 2018. Voxceleb2: Deep speaker recognition. *arXiv preprint arXiv:1806.05622*.
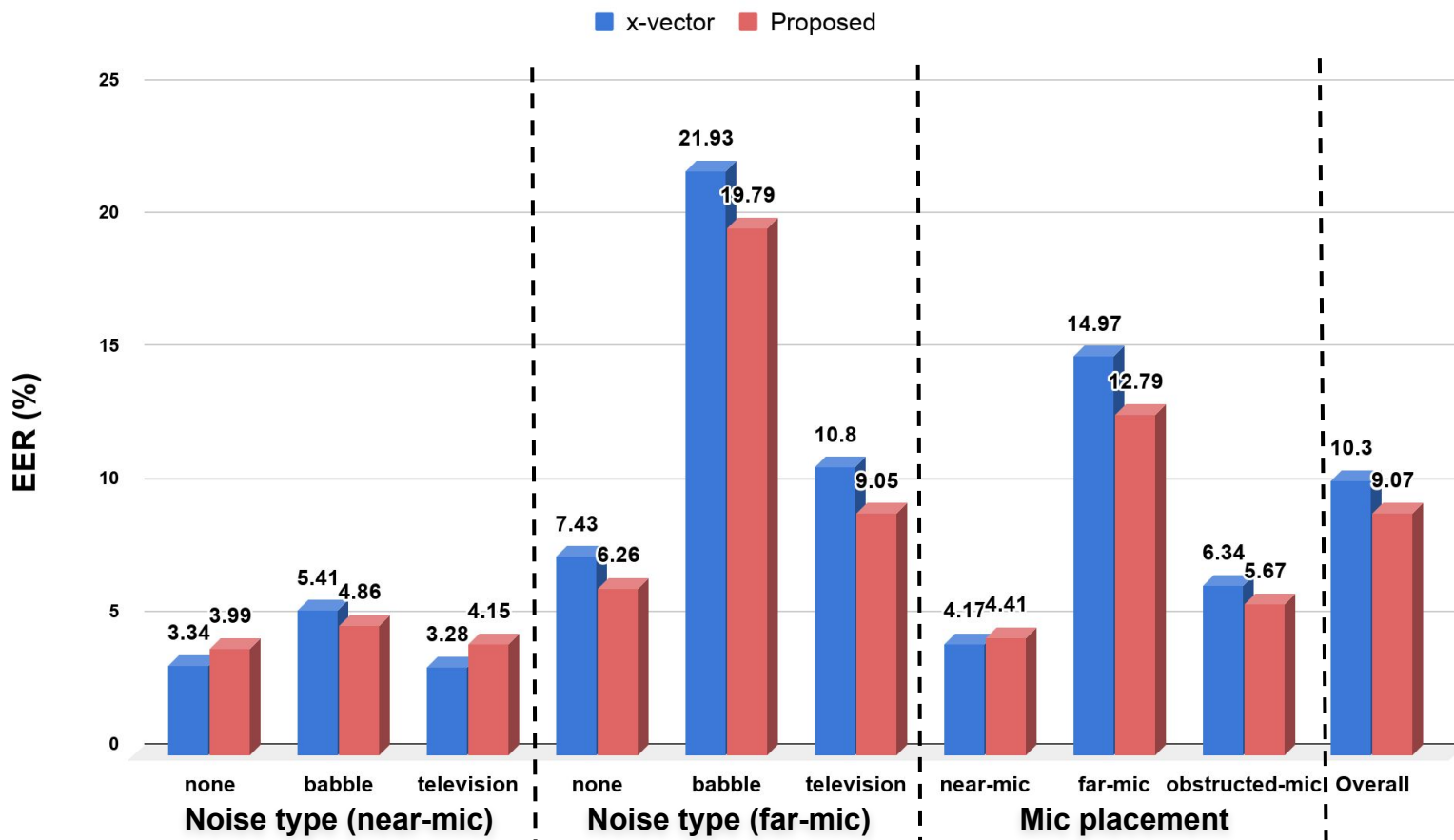2. https://kaldi-asr.org/models/m7
3. Richey, Colleen, Maria A. Barrios, Zeb Armstrong, Chris Bartels, Horacio Franco, Martin Graciarena, Aaron Lawson et al. "Voices obscured in complex environmental settings (voices) corpus." *arXiv preprint arXiv:1804.05053* (2018).
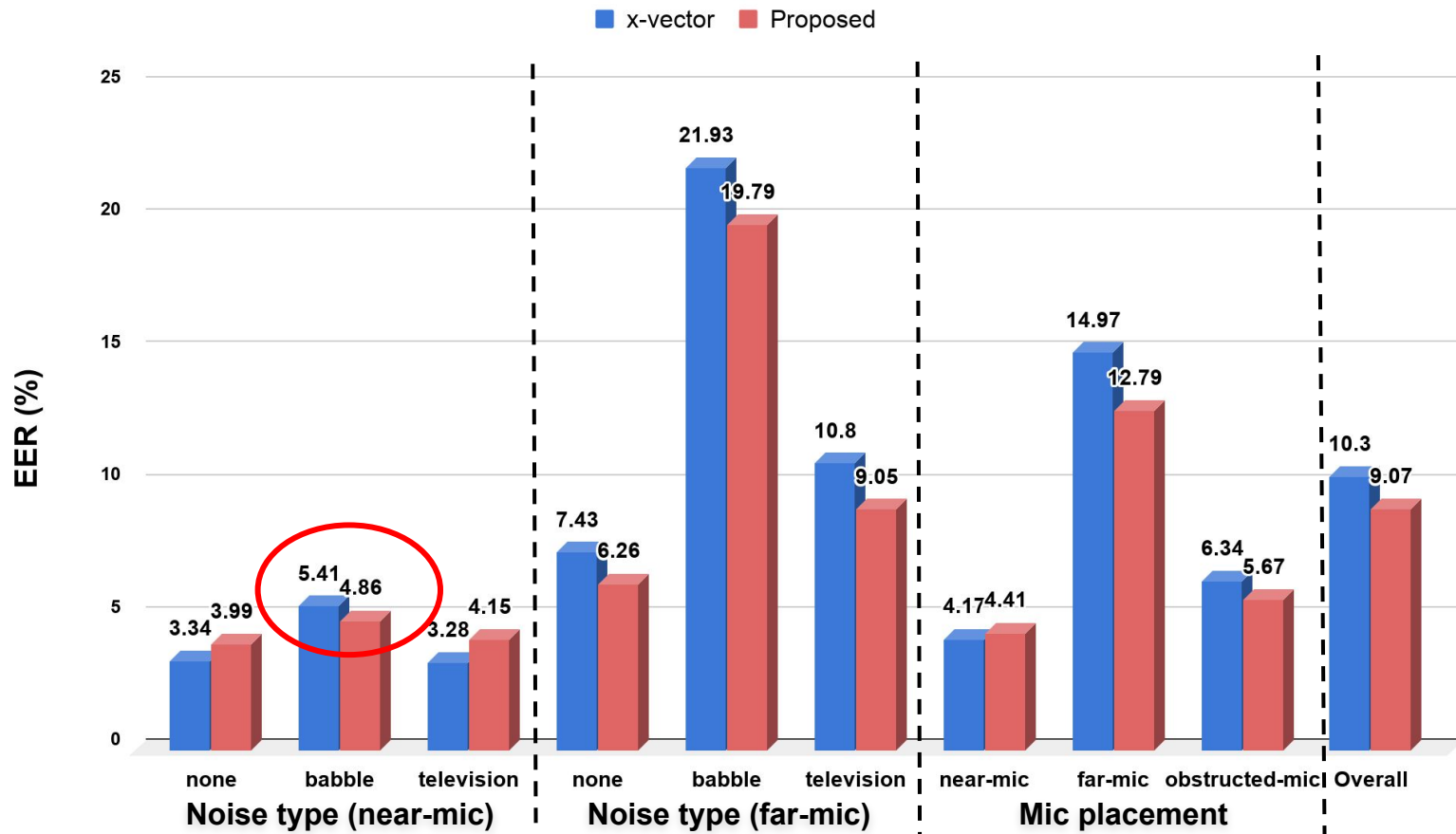
- Dimensionality reduction: Linear Discriminant Analysis (LDA)
  - x-vector - dimension 150, Proposed - dimension 96
- Verification scoring: Probabilistic LDA (PLDA)



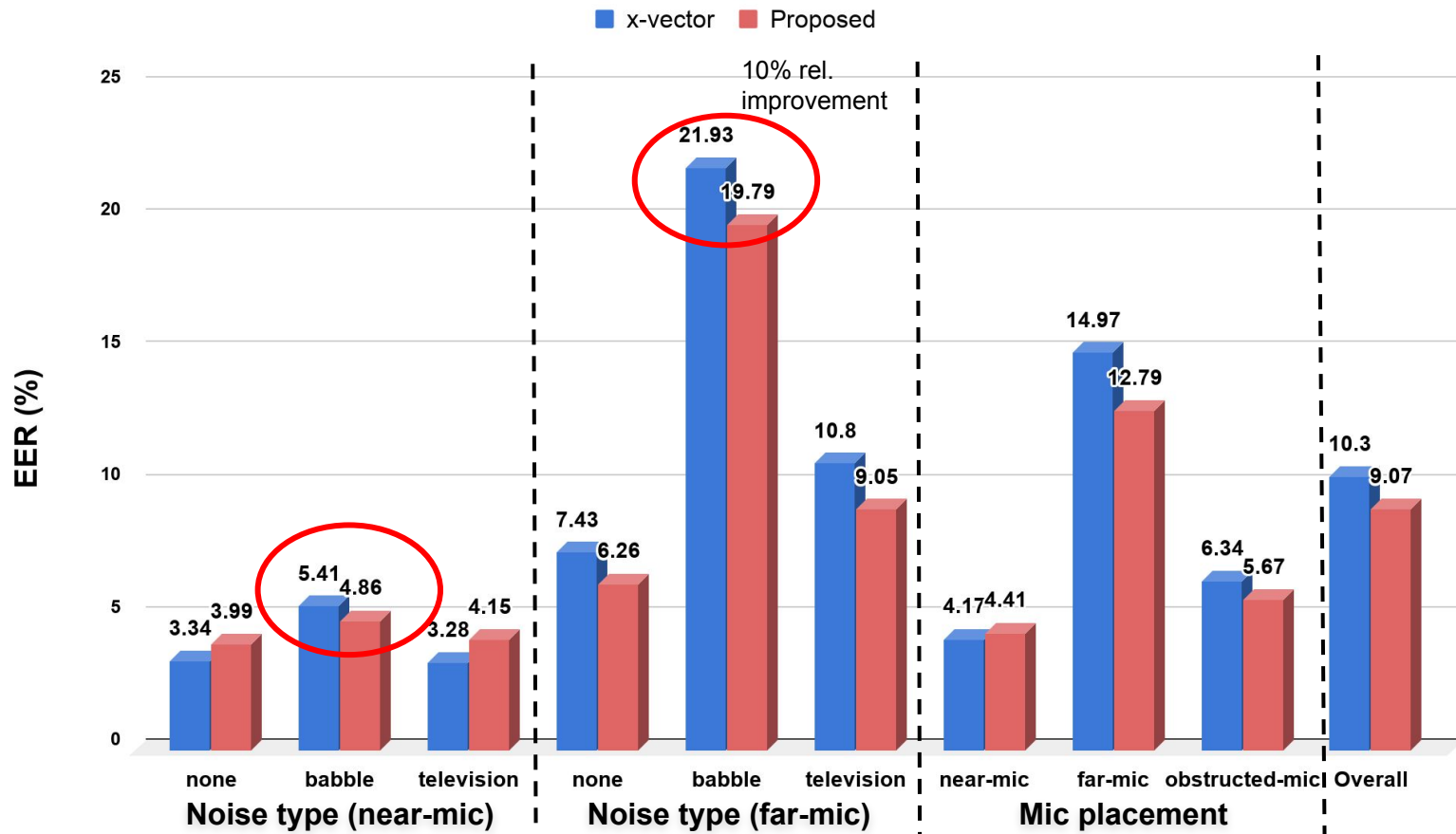Speaker verification performance on VOiCES-eval

21

- Dimensionality reduction: Linear Discriminant Analysis (LDA)
  - x-vector - dimension 150, Proposed - dimension 96
- Verification scoring: Probabilistic LDA (PLDA)

## Speaker verification performance on VOiCES-eval

■ x-vector ■ Proposed



EER (%)

| | x-vector | Proposed |
|---|---|---|
| none (near-mic) | 3.34 | 3.99 |
| babble (near-mic) | 5.41 | 4.86 |
| television (near-mic) | 3.28 | 4.15 |
| none (far-mic) | 7.43 | 6.26 |
| babble (far-mic) | 21.93 | 19.79 |
| television (far-mic) | 10.8 | 9.05 |
| near-mic | 4.17 | 4.41 |
| far-mic | 14.97 | 12.79 |
| obstructed-mic | 6.34 | 5.67 |
| Overall | 10.3 | 9.07 |

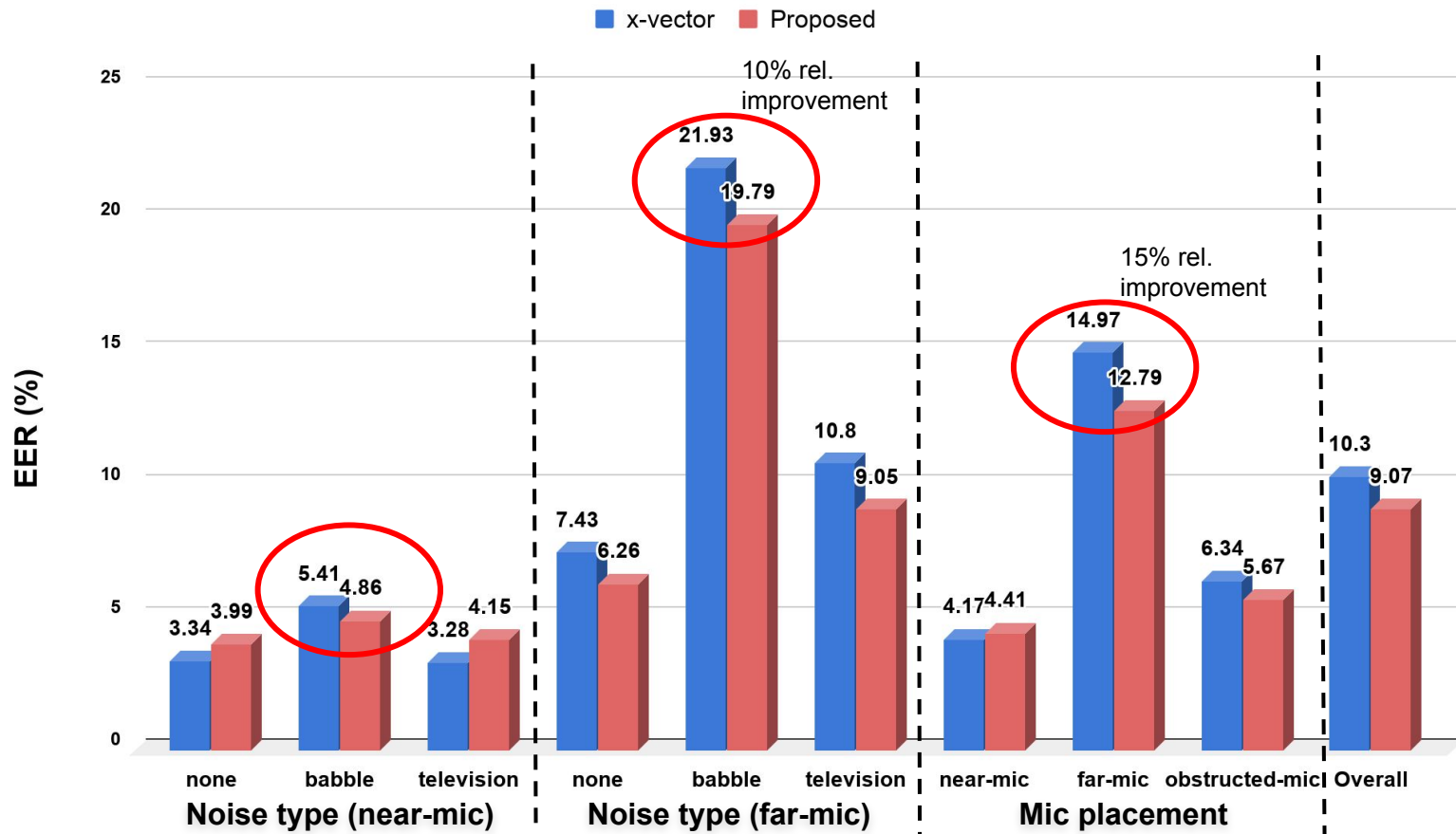**Noise type (near-mic)** | **Noise type (far-mic)** | **Mic placement**
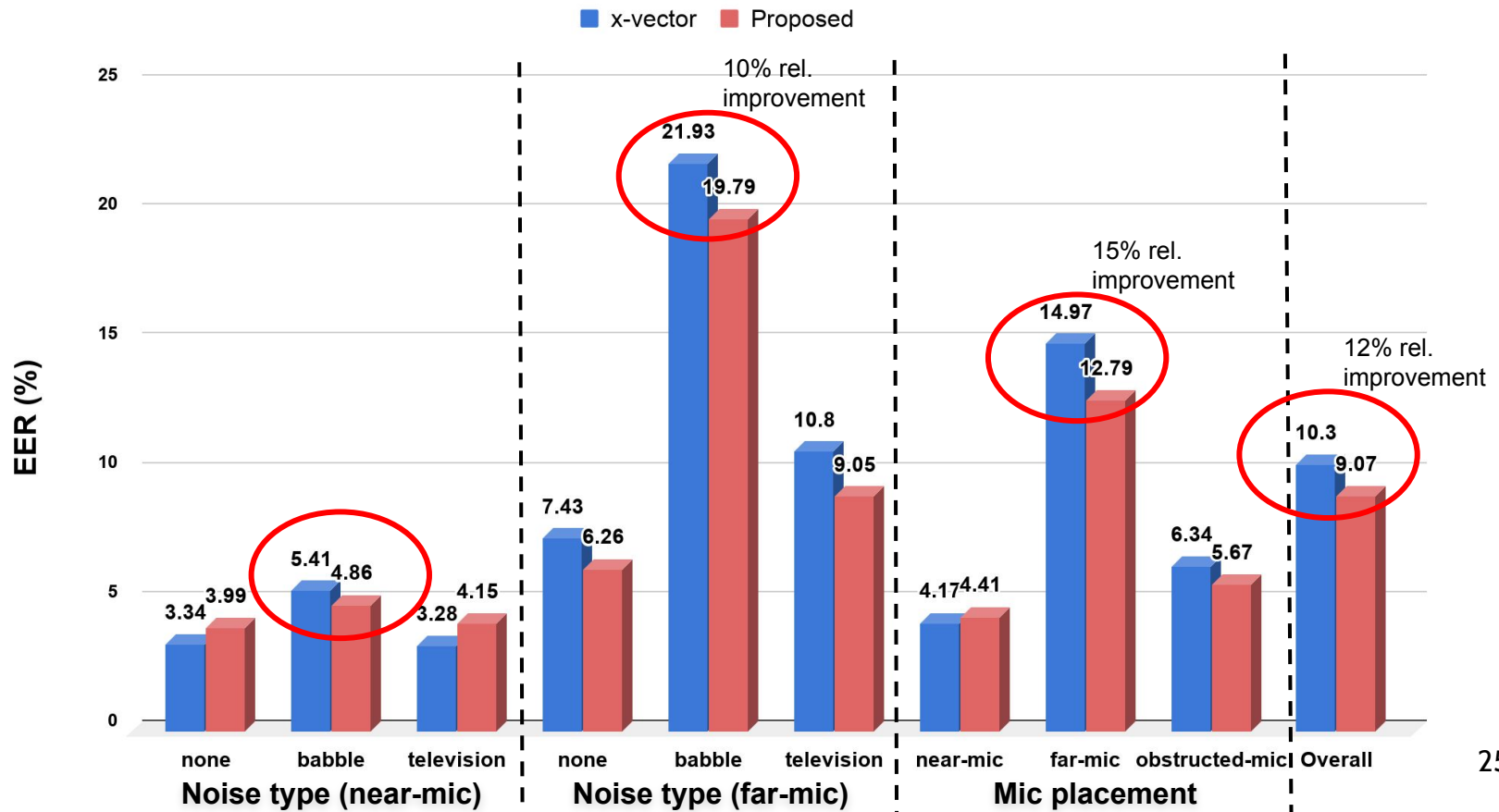
22

- Dimensionality reduction: Linear Discriminant Analysis (LDA)
  - x-vector - dimension 150, Proposed - dimension 96
- Verification scoring: Probabilistic LDA (PLDA)

**Speaker verification performance on VOiCES-eval**



■ x-vector  ■ Proposed

10% rel. improvement

EER (%)

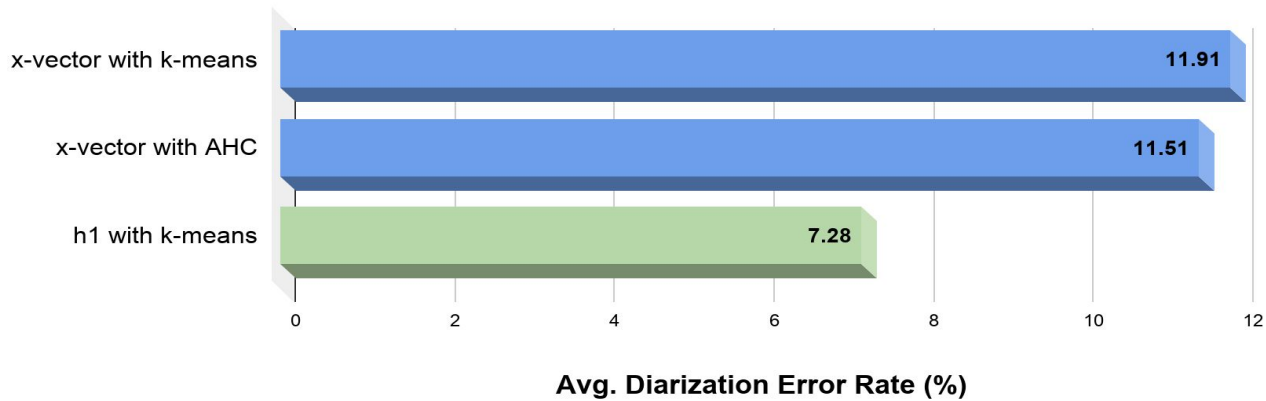| Noise type (near-mic) | | | Noise type (far-mic) | | | Mic placement | | | |
|---|---|---|---|---|---|---|---|---|---|
| none | babble | television | none | babble | television | near-mic | far-mic | obstructed-mic | Overall |
| 3.34 / 3.99 | 5.41 / 4.86 | 3.28 / 4.15 | 7.43 / 6.26 | 21.93 / 19.79 | 10.8 / 9.05 | 4.17 / 4.41 | 14.97 / 12.79 | 6.34 / 5.67 | 10.3 / 9.07 |

23

- Dimensionality reduction: Linear Discriminant Analysis (LDA)
  - x-vector - dimension 150, Proposed - dimension 96
- Verification scoring: Probabilistic LDA (PLDA)



Speaker verification performance on VOiCES-eval

24

- Dimensionality reduction: Linear Discriminant Analysis (LDA)
  - x-vector - dimension 150, Proposed - dimension 96
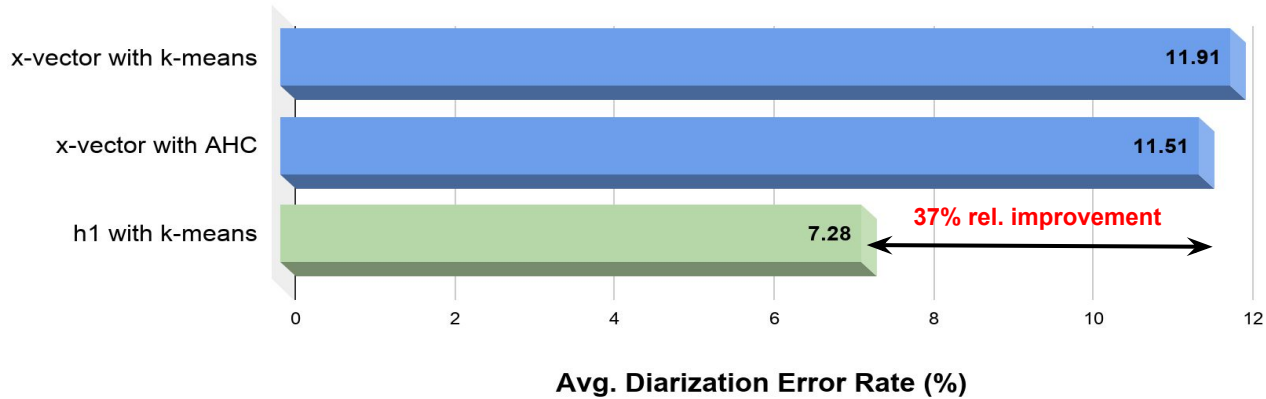- Verification scoring: Probabilistic LDA (PLDA)

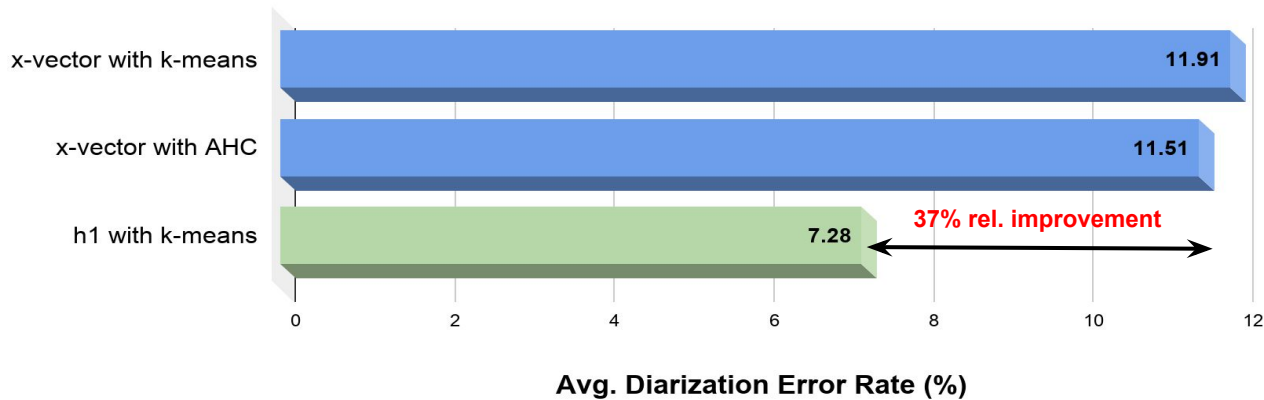## Speaker verification performance on VOiCES-eval



25

**Speaker Diarization performance on AMI meeting corpus compared to two competitive baselines (oracle SAD, known num. speakers)**



**Avg. Diarization Error Rate (%)**

**Speaker Diarization performance on AMI meeting corpus compared to two competitive baselines (oracle SAD, known num. speakers)**



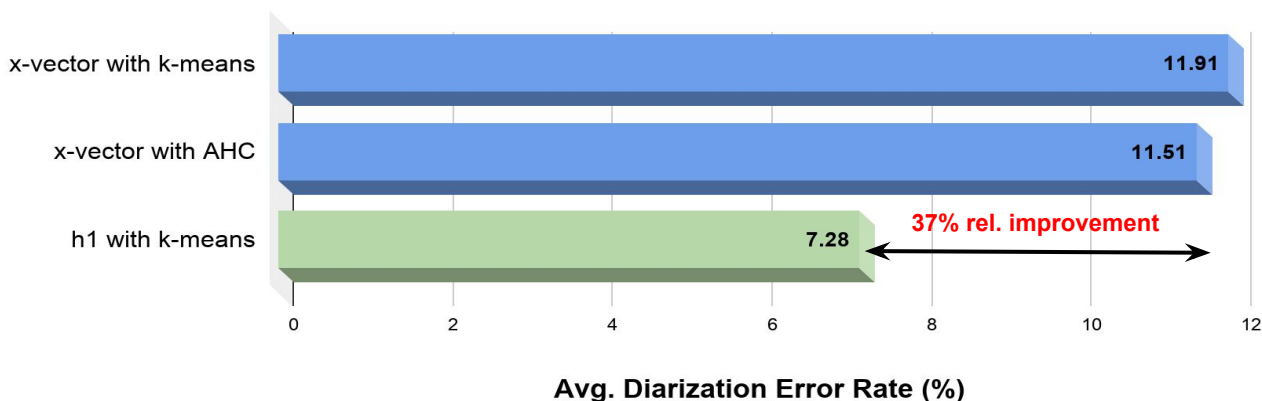| | Avg. Diarization Error Rate (%) |
|---|---|
| x-vector with k-means | 11.91 |
| x-vector with AHC | 11.51 |
| h1 with k-means | 7.28 — 37% rel. improvement |

**Avg. Diarization Error Rate (%)**

27

**Speaker Diarization performance on AMI meeting corpus compared to two competitive baselines (oracle SAD, known num. speakers)**



Avg. Diarization Error Rate (%)

## Conclusions

- Proposed novel speaker embeddings

    - Disentangled speaker and nuisance factors from speaker embeddings

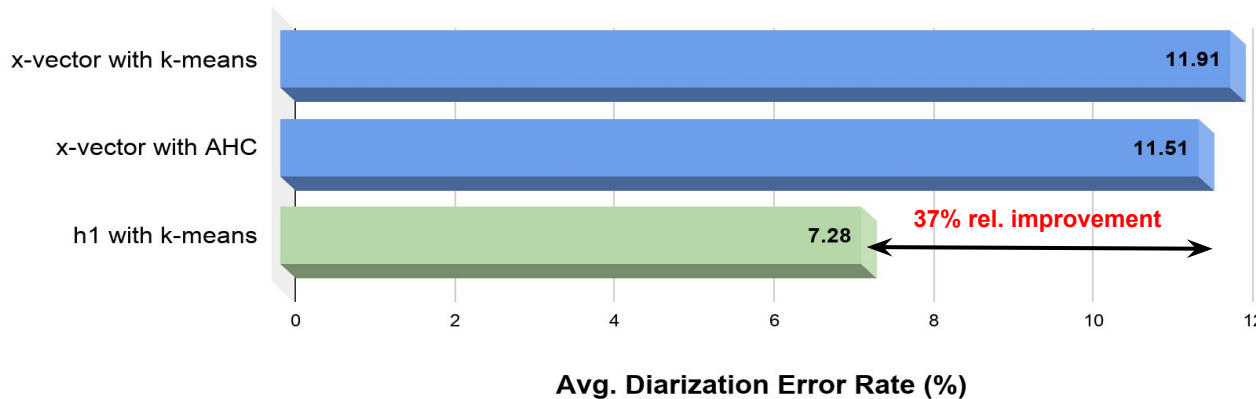    - No prior knowledge of specific nuisance factors during training

**Speaker Diarization performance on AMI meeting corpus compared to two competitive baselines (oracle SAD, known num. speakers)**



**Avg. Diarization Error Rate (%)**

Chart data:
- x-vector with k-means: 11.91
- x-vector with AHC: 11.51
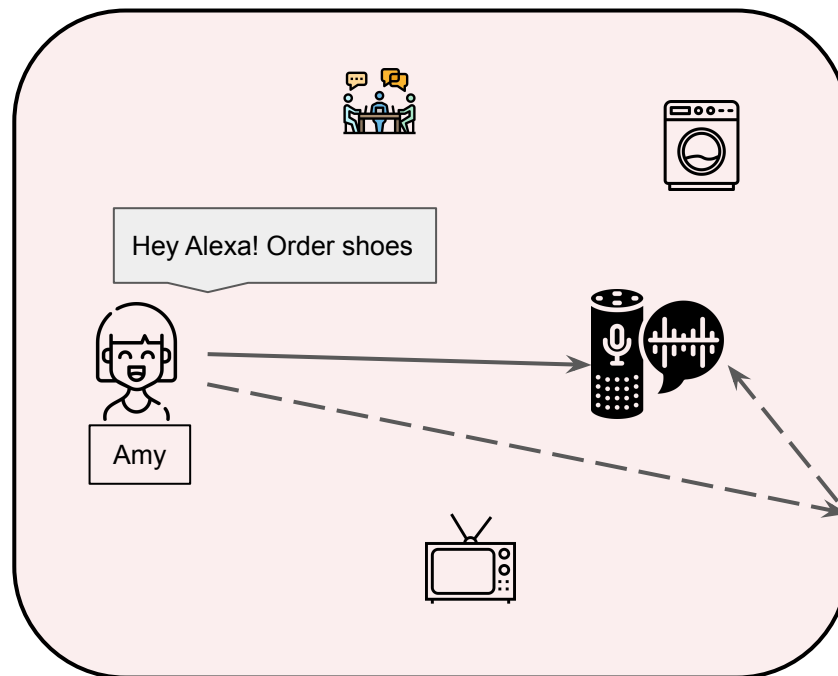- h1 with k-means: 7.28 — **37% rel. improvement**

## Conclusions

- Proposed novel speaker embeddings

  - Disentangled speaker and nuisance factors from speaker embeddings

  - No prior knowledge of specific nuisance factors during training

- Improves speaker verification performance in challenging conditions

  - Particularly babble noise (10% EER) and far-field recording conditions (15% EER)
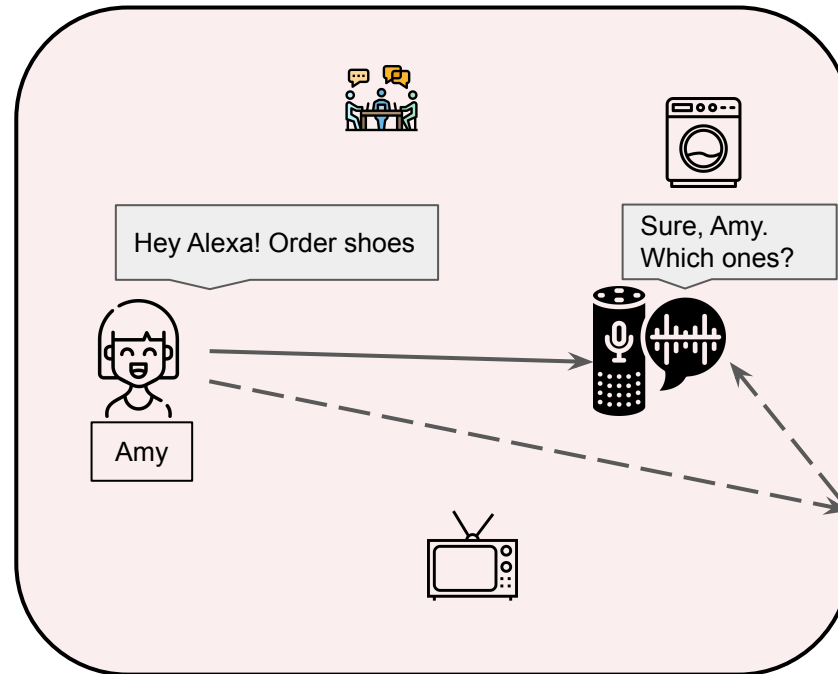
**Speaker Diarization performance on AMI meeting corpus compared to two competitive baselines (oracle SAD, known num. speakers)**

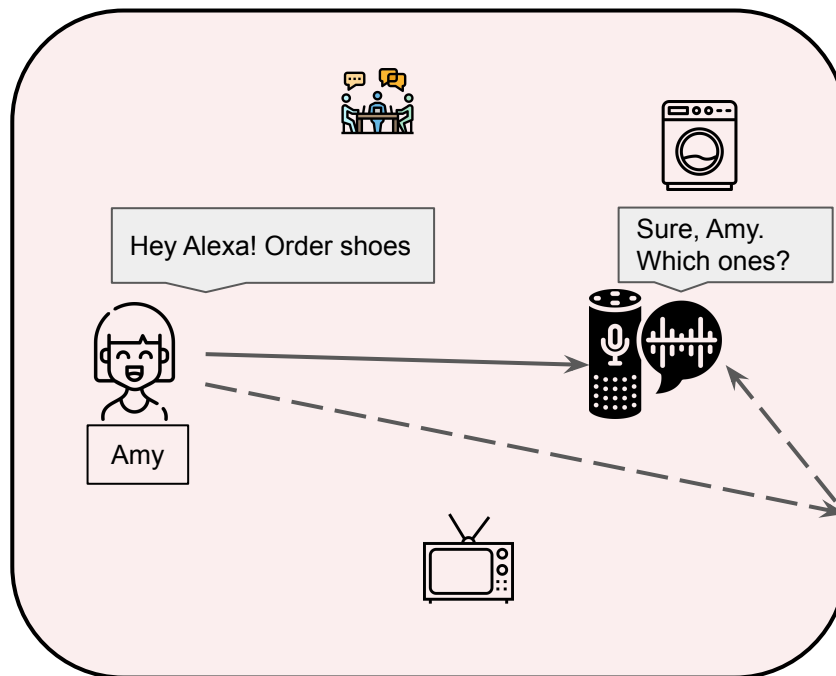| Method | Avg. Diarization Error Rate (%) |
|---|---|
| x-vector with k-means | 11.91 |
| x-vector with AHC | 11.51 |
| h1 with k-means | 7.28 — 37% rel. improvement |

**Avg. Diarization Error Rate (%)**

## Conclusions

- Proposed novel speaker embeddings

  - Disentangled speaker and nuisance factors from speaker embeddings

  - No prior knowledge of specific nuisance factors during training

- Improves speaker verification performance in challenging conditions

  - Particularly babble noise (10% EER) and far-field recording conditions (15% EER)

- Improves speaker diarization performance on AMI meeting corpus (37% DER)

## Future work

- Improve performance in babble noise scenario

- Evaluate disentangled speaker embeddings in presence of other nuisance factors, such as affective state, lexical content

- Train with more basic speech representations, which contain more variability useful for disentangling

33

# We gratefully acknowledge the support of USG for this work

**Contact Info**
Name: Raghuveer Peri
Email: rperi@usc.edu