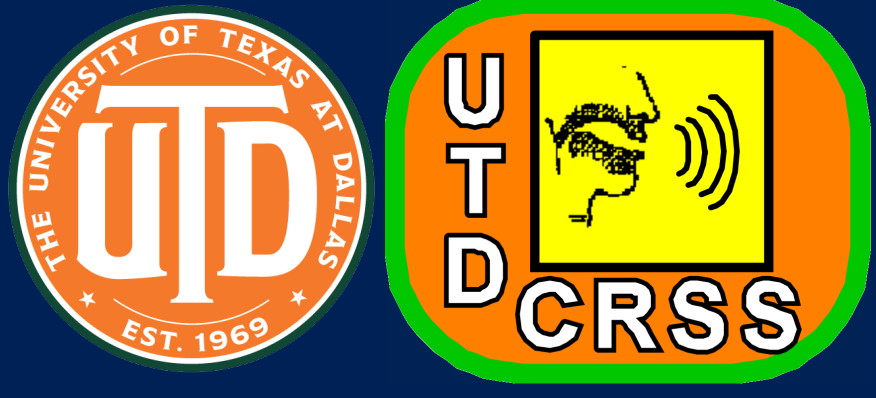




AALBORG UNIVERSITY  
DENMARK



# Improved Vocal Effort Transfer Vector Estimation for Vocal Effort-Robust Speaker Verification

Iván López-Espejo<sup>1,2</sup>, Santi Prieto<sup>3</sup>, Alfonso Ortega<sup>4</sup> and Eduardo Lleida<sup>4</sup>

<sup>1</sup>Department of Electronic Systems, Aalborg University, Denmark

<sup>2</sup>Center for Robust Speech Systems (CRSS), The University of Texas at Dallas, USA

<sup>3</sup>VeriDas | das-Nano, Navarre, Spain

<sup>4</sup>ViVoLab, Aragón Institute for Engineering Research (I3A), University of Zaragoza, Spain

ivl@es.aau.dk, sprieto@veridas.com, {ortega,lleida}@unizar.es



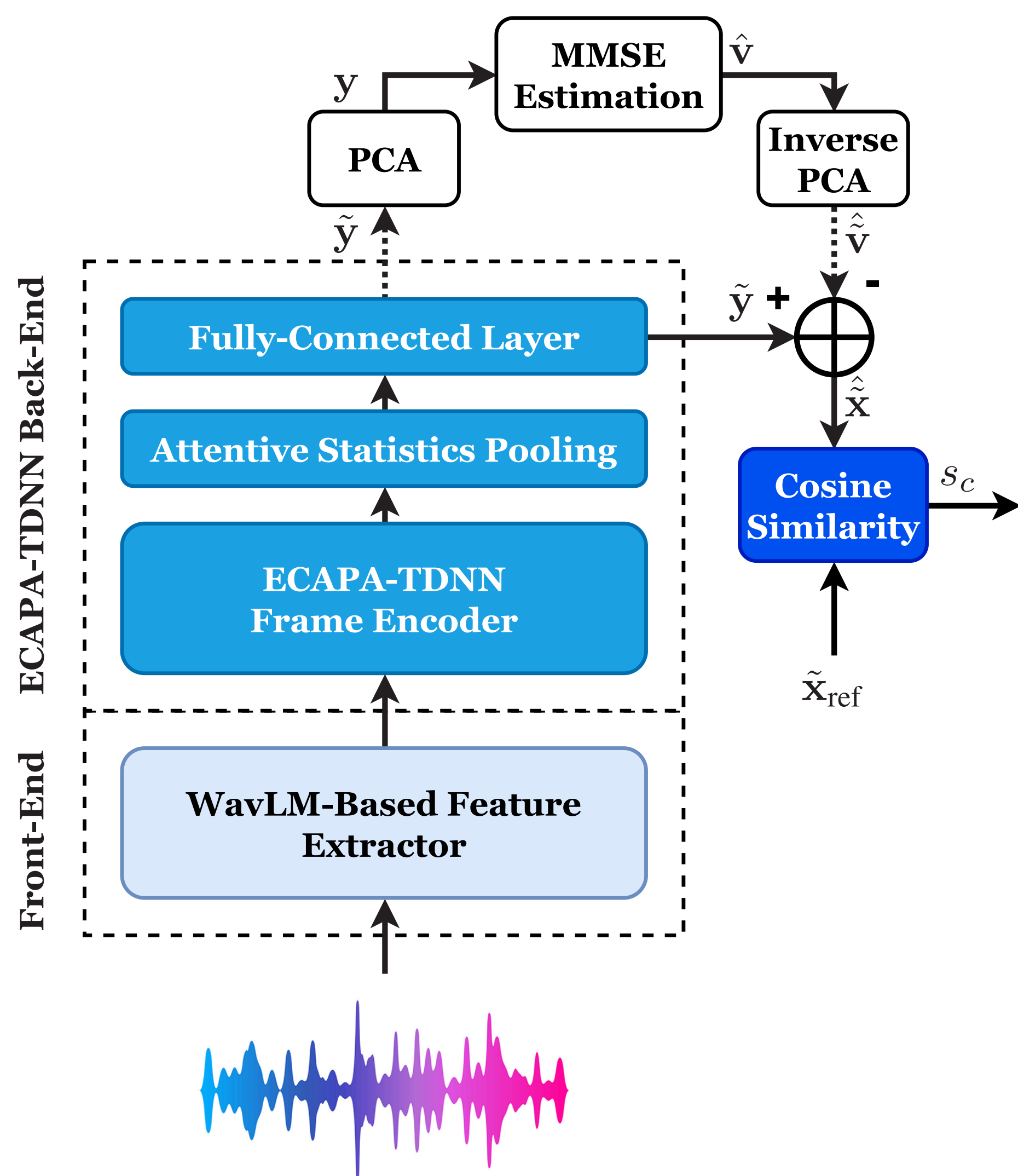
Funded by  
the European Union

## Introduction

- Speaker verification performance tends to dramatically drop in the presence of non-neutrally-phonated (e.g., *shouted and whispered*) speech
- Previous work explored a series of *minimum mean square error (MMSE)* techniques estimating **normal speaker embeddings from non-neutrally-phonated ones**
- **MEMLIN** (Multi-Environment Model-based Linear Normalization) provided the best performance in terms of equal error rate (EER) when dealing with both shouted and whispered speech
- In this work **we tackle a MEMLIN's shortcoming**, which is explained in the next box 🙋

## System Overview

- Speaker embedding compensation is applied *only* in case that the embedding comes from non-neutrally-phonated speech
- The ECAPA-TDNN back-end is trained on an augmented version of the **VoxCeleb2** dataset



- **Shouted-normal speech corpus:** Paired utterances in Finnish from 22 speakers
- **Whispered-normal speech corpus:** Paired utterances in English from 36 speakers (*CHAINS*)
- Due to speech data scarcity, experiments are performed by following a **leave-one-speaker-out cross-validation** strategy

## Normal Speaker Embedding Estimation

$\tilde{\mathbf{x}} \in \mathbb{R}^D$ : Normal embedding |  $\tilde{\mathbf{y}} \in \mathbb{R}^D$ : Non-neutrally-phonated embedding |  $\tilde{\mathbf{v}} \in \mathbb{R}^D$ : Vocal effort transfer vector

### MMSE Estimation

$$\tilde{\mathbf{y}} = \tilde{\mathbf{x}} + \tilde{\mathbf{v}} \rightarrow \text{Assuming } \tilde{\mathbf{y}} \text{ is modeled by a } K\text{-component GMM} \rightarrow \hat{\mathbf{x}} = \tilde{\mathbf{y}} - \underbrace{\sum_{k=1}^K P(k|\tilde{\mathbf{y}})\hat{\mathbf{v}}^{\{k\}}}_{\hat{\mathbf{v}}}$$

- **Limitation of MEMLIN:** The set of partial estimates  $\{\hat{\mathbf{v}}^{\{k\}}; k = 1, \dots, K\}$  is pre-computed (during an offline training stage) and fixed

- **To overcome MEMLIN's shortcoming, we propose MMSE<sub>v</sub>:**

1. We jointly model  $\tilde{\mathbf{v}}$  and  $\tilde{\mathbf{y}}$  by a  $K$ -component GMM  $p(\tilde{\mathbf{z}} = (\tilde{\mathbf{v}}, \tilde{\mathbf{y}}))$
2. Estimation is carried out in a **principal component analysis (PCA)** domain to face data scarcity

Let  $\mathbf{W}_L$  be a  $D \times L$  PCA transform matrix, where  $L \ll D = 256$

$$\mathbf{v} = \mathbf{W}_L^T \tilde{\mathbf{v}}, \quad \mathbf{y} = \mathbf{W}_L^T \tilde{\mathbf{y}}$$

$$p(\mathbf{z} = (\mathbf{v}, \mathbf{y}) \in \mathbb{R}^{2L}) = \sum_{k=1}^K P(k) \mathcal{N}(\mathbf{z} | \mu_z^{\{k\}}, \Sigma_z^{\{k\}}), \quad \mu_z^{\{k\}} = \begin{pmatrix} \mu_v^{\{k\}} \\ \mu_y^{\{k\}} \end{pmatrix}, \quad \Sigma_z^{\{k\}} = \begin{pmatrix} \Sigma_{vv}^{\{k\}} & \Sigma_{vy}^{\{k\}} \\ \Sigma_{yv}^{\{k\}} & \Sigma_{yy}^{\{k\}} \end{pmatrix}$$

### MMSE<sub>v</sub> Compensation

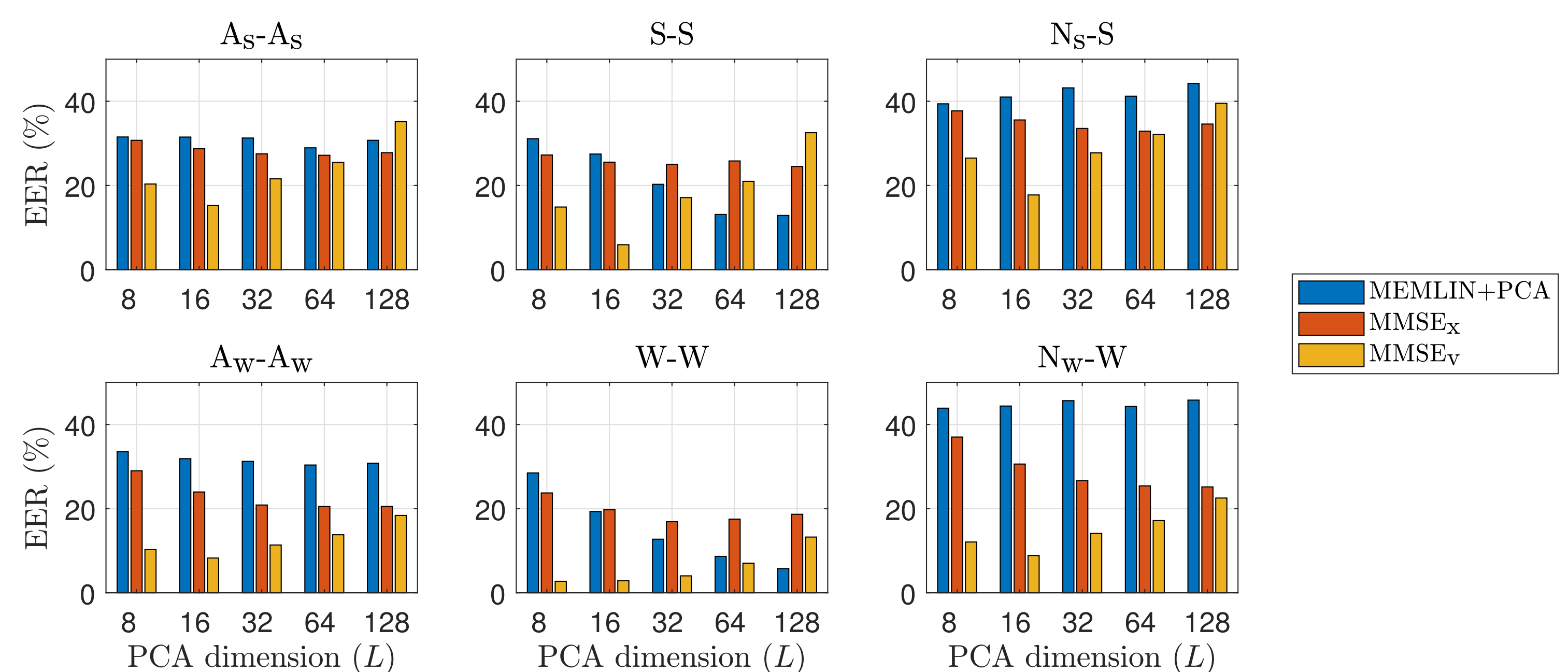
$$\hat{\mathbf{v}} = \mathbb{E}(\mathbf{v}|\mathbf{y}) = \sum_{k=1}^K P(k|\mathbf{y}) \underbrace{\mathbb{E}(\mathbf{v}|\mathbf{y}, k)}_{\hat{\mathbf{v}}^{\{k\}}} \rightarrow \hat{\mathbf{x}} = \tilde{\mathbf{y}} - \underbrace{\mathbf{W}_L \hat{\mathbf{v}}}_{\hat{\mathbf{v}}}$$

1. **Combination weights:**  $P(k|\mathbf{y}) = \frac{p(\mathbf{y}|k)P(k)}{\sum_{k'=1}^K p(\mathbf{y}|k')P(k')}$
2. **Partial estimates:**  $\mathbb{E}(\mathbf{v}|\mathbf{y}, k) = \mu_v^{\{k\}} + \Sigma_{vy}^{\{k\}} (\Sigma_{yy}^{\{k\}})^{-1} (\mathbf{y} - \mu_y^{\{k\}})$

- Both  $\mathbf{W}_L$  and  $p(\mathbf{z})$  are calculated from paired normal and non-neutrally-phonated embeddings

## Experimental Results and Discussion

- **EER (%)** is the chosen speaker verification metric | Use of  $K = 8$ -component GMMs
- Embedding compensation experiments are carried out by employing *E-T+WavLM as the baseline system*
- **MMSE<sub>x</sub>:** MMSE estimator equivalent to MMSE<sub>v</sub> that directly estimates  $\tilde{\mathbf{x}}$  from  $\mathbb{E}(\mathbf{x}|\mathbf{y})$



### Shouted and normal speech:

Condition	E-T+MFCC	E-T+WavLM	MEMLIN	MEMLIN+PCA	MMSE <sub>x</sub>	MMSE <sub>v</sub>
$A_S-A_S$	19.96	17.11	15.62	31.50	28.72	<b>15.22</b>
$N_S-N_S$	9.73	<b>7.25</b>	<b>7.25</b>	<b>7.25</b>	<b>7.25</b>	<b>7.25</b>
$S-S$	11.58	9.94	10.44	27.46	25.53	<b>5.91</b>
$N_S-S$	25.28	21.76	20.74	41.00	35.56	<b>17.74</b>

### Whispered and normal speech:

Condition	E-T+MFCC	E-T+WavLM	MEMLIN	MEMLIN+PCA	MMSE <sub>x</sub>	MMSE <sub>v</sub>
$A_W-A_W$	16.54	11.24	<b>8.25</b>	31.87	23.95	8.27
$N_W-N_W$	1.21	<b>0.62</b>	<b>0.62</b>	<b>0.62</b>	<b>0.62</b>	<b>0.62</b>
$W-W$	4.38	5.26	4.00	19.31	19.77	<b>2.87</b>
$N_W-W$	12.81	9.81	11.47	44.38	30.59	<b>8.86</b>