# Non-parallel Voice Conversion Using i-Vector PLDA:
## Towards Unifying Speaker Verification and Transformation

**Tomi Kinnunen, Lauri Juvela, Paavo Alku, Junichi Yamagishi**

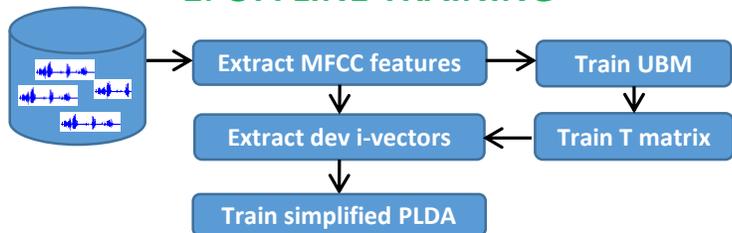UNIVERSITY OF EASTERN FINLAND

A! Aalto University

NII

Fully non-parallel system: no frame-level alignment, parallel reference speaker or parallel source/target data required.
Source-to-target **conversion function trained on utterance-level speaker features** instead of low-level frame features.

## 1. OFFLINE TRAINING

Extract MFCC features → Train UBM

Extract dev i-vectors ← Train T matrix
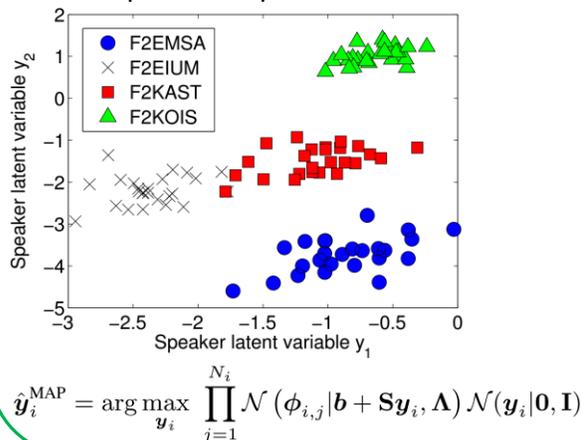
Train simplified PLDA

An **i-vector** is a single low-dimensional 'feature vector' of a speech utterance. We model i-vector distributions with **simplified probabilistic linear discriminant analysis** (PLDA) model:
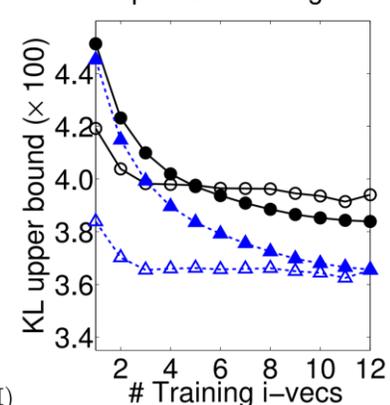
j:th i-vector of i:th speaker

Speaker factor matrix

Speaker variable. **Convert this**

Residual: within-speaker variation related to content, microphone, etc.

bias

$$\boldsymbol{\phi}_{i,j} = \boldsymbol{b} + \mathbf{S}\boldsymbol{y}_i + \boldsymbol{\varepsilon}_{i,j}$$

Datapoint = a speech utterance



Nonparallel training data

Parallel training data

F2EMSA, F2EIUM, F2KAST, F2KOIS

Q=100 (test), Q=100 (training), Q=800 (test), Q=800 (training)

$$\hat{\boldsymbol{y}}_i^{\mathrm{MAP}} = \arg\max_{\boldsymbol{y}_i} \prod_{j=1}^{N_i} \mathcal{N}\left(\boldsymbol{\phi}_{i,j} | \boldsymbol{b} + \mathbf{S}\boldsymbol{y}_i, \boldsymbol{\Lambda}\right) \mathcal{N}\left(\boldsymbol{y}_i | \boldsymbol{0}, \mathbf{I}\right)$$

## 2. TRAINING: extract latent speaker features with PLDA

$\boldsymbol{\Phi}_{\mathrm{src}} = \{\boldsymbol{\phi}_n\}$ ➔ MAP estimate of $\mathbf{y}_{\mathrm{src}}$

$\boldsymbol{\Phi}_{\mathrm{tar}} = \{\boldsymbol{\phi}_m\}$ ➔ MAP estimate of $\mathbf{y}_{\mathrm{tar}}$

No 'pairing': independent extraction of the latent speaker features for each speaker. Number of utterances can also vary.

## 3. VOICE CONVERSION: predict target speaker i-vector (i.e. GMM)

PLDA decomposition of source speaker test utterance:

$$\boldsymbol{\phi}_{\mathrm{src}} = \boldsymbol{b} + \mathbf{S}\hat{\boldsymbol{y}}_{\mathrm{src}} + \boldsymbol{e}_{\mathrm{src}}$$
$$\hat{\boldsymbol{\phi}}_{\mathrm{tar}} = \boldsymbol{b} + \mathbf{S}\hat{\boldsymbol{y}}_{\mathrm{tar}} + \boldsymbol{e}_{\mathrm{src}} = \boldsymbol{\phi}_{\mathrm{src}} + \mathbf{S}(\hat{\boldsymbol{y}}_{\mathrm{tar}} - \hat{\boldsymbol{y}}_{\mathrm{src}})$$

Replace source latent identity variable with that of the target speaker.
The predicted target speaker i-vector then defines target GMM means and frame-level conversion:

$$\boldsymbol{\mu}_c^{\mathrm{tar}} = \boldsymbol{m}_c + \mathbf{T}_c \hat{\boldsymbol{\phi}}_{\mathrm{tar}} \quad ➔ \quad \hat{\boldsymbol{y}}_t = \boldsymbol{x}_t + \sum_{c=1}^{C} P(c|\boldsymbol{x}_t)\left(\boldsymbol{\mu}_c^{\mathrm{tar}} - \boldsymbol{\mu}_c^{\mathrm{src}}\right)$$