



Sequential Voice Conversion Using Grid-Based Approximation

H. Benisty, D. Malah and K. Crammer

Electrical Engineering Department
Technion – Israel Institute of Technology

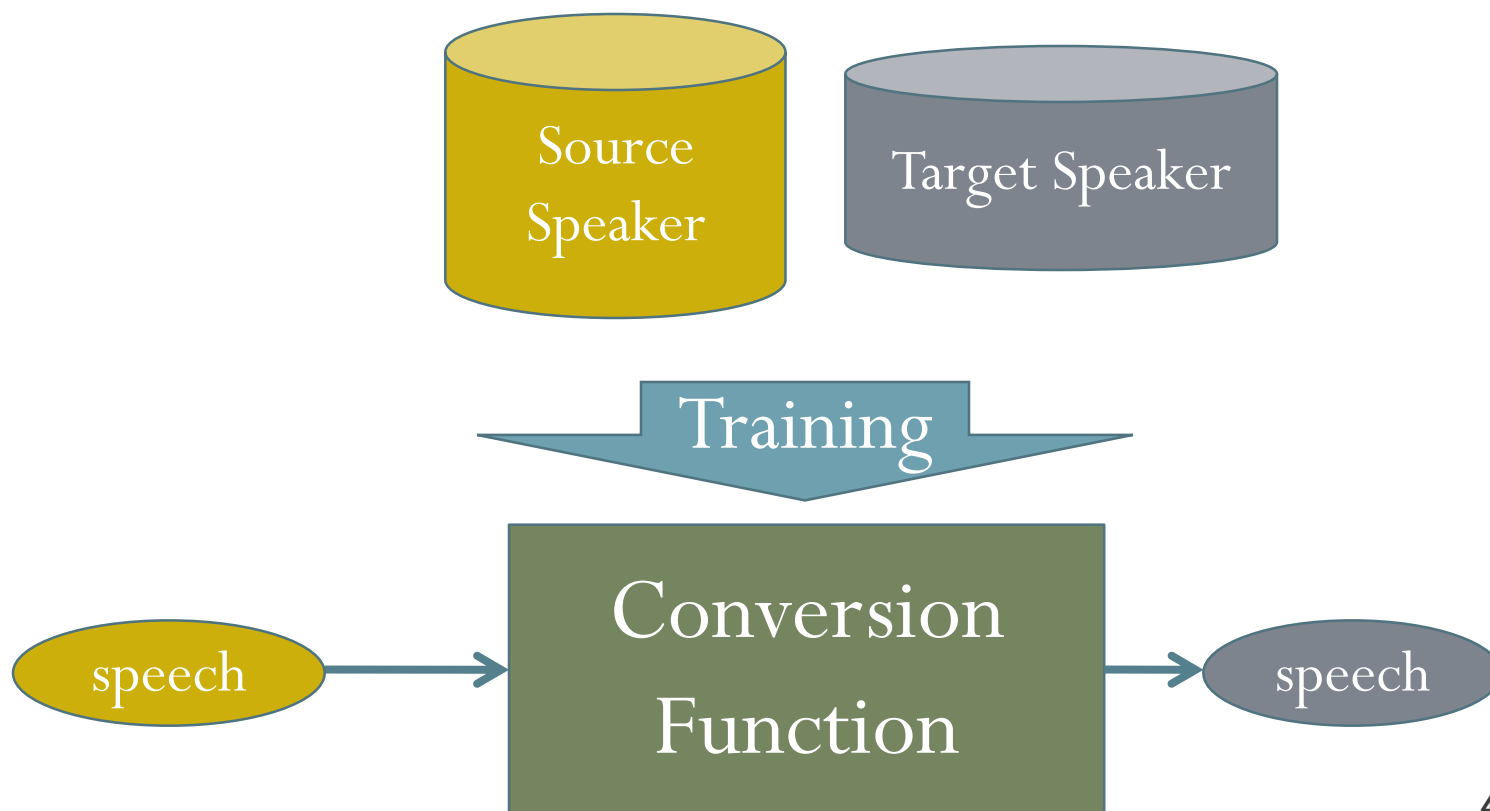
IEEEI 2014 – Eilat, Israel

Agenda

- Introduction
 - Motivation
 - General VC Scheme
- A Common Method - GMM-based conversion
- Grid-Based (GB) Conversion
 - Formulation
- Experimental Results
- Conclusion

General Conversion Setup

- **The goal:** modify a source speaker's speech to sound as if spoken by a target speaker



Applications

- Personalize the output of Text-to-Speech (TTS) systems
- Dialog systems
- Vocal pathology
- Entertainment
- Toys & games

Speech Characteristics

- The identity of a speaker is associated with:
 - Prosody attributes - pitch, duration and energy
 - Spectral envelope
- Pitch - usually modified using a simple statistical mean and variance scaling
- Most VC methods deal with **spectral envelope conversion**

Training Data

- Textual Content

- Parallel data set – the two speakers say the same text
- Non-parallel data set – source-target correspondence is learned during training
- Cross-lingual data set

GMM-Based Conversion

[Stylianou et. al., 1998; Kain & Macon, 1998]

- Given a **parallel** and **aligned** source and target training vectors $\{\mathbf{x}^k, \mathbf{y}^k\}_1^N \in \mathbb{R}^P$ (represented by Mel Frequency Cepstrum Coefficients - MFCCs)

- A GMM is trained using the source vectors:

$$p(\mathbf{x}) = \sum_{m=1}^M \alpha_m N(\mathbf{x}; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$$

- The conversion function is formulated, using the GMM parameters, as a weighted sum of linear Bayesian estimators of the target spectra:

$$\mathcal{F}(\mathbf{x}) = \sum_{m=1}^M \alpha_m (A_m \mathbf{x} + b_m)$$

GMM-Based Conversion - Cont.

- Exhaustive training - Expectation Maximization
 - Often leading to ill-conditioning if the dataset used is too small
- Linear conversion → over-smoothed spectral envelopes → muffled synthesizes speech
 - Proposed remedy – global variance enhancement
[Toda et al, 2007; Benisty and Malah, 2011; Benisty et al., 2012]

Grid-Based (GB) Conversion

- Based on sequential **Bayesian tracking**
- Expressed as a sequential estimation problem of tracking the **target spectrum** based on the observed **source spectrum**
- The converted MFCC vectors are sequentially evaluated using a **weighted sum** of the target training set used as **grid-points**

Grid-based conversion

Bayesian Tracking

- The target spectrum \mathbf{y}_t , is assumed to follow a first order Markov dynamics: $\mathbf{y}_t = f_t(\mathbf{y}_{t-1}, \mathbf{u}_t)$
 \mathbf{u}_t - an i.i.d. noise sequence.

- The source spectrum \mathbf{x}_t , depends on the target spectrum through: $\mathbf{x}_t = h_t(\mathbf{y}_t, \mathbf{v}_t)$
 \mathbf{v}_t - an i.i.d. measurement noise

- The Bayesian optimal estimation for the target spectrum is:

$$\hat{\mathbf{y}}_t = E[\mathbf{y}_t | \mathbf{x}_{1:T}] = \int p(\mathbf{y}_t | \mathbf{x}_{1:T}) \mathbf{y}_t d\mathbf{y}_t$$

- In practice an analytical derivation is impossible since we do not have a closed form for $p(\mathbf{y}_t | \mathbf{x}_{1:T})$, so we use **Grid-Based approximation**

Grid-based conversion

Discrete Approximation

- Given: parallel unaligned training sets — $\{\mathbf{x}^k\}_1^{N_x} \in \mathfrak{R}^P, \{\mathbf{y}^k\}_1^{N_y} \in \mathfrak{R}^P$
and a sequence of test source vectors - $\mathbf{x}_{1:T}$
- We evaluate the posterior probability as a discrete sum:

$$p(\mathbf{y}_t | \mathbf{x}_{1:T}) = \sum_{k=1}^{N_y} w_{t|t}^k \delta(\mathbf{y}_t = \mathbf{y}^k)$$

where the posterior weights are: $w_{t|t}^k = p(\mathbf{y}_t = \mathbf{y}^k | \mathbf{x}_{1:T})$

- The optimal Bayesian estimation for the target spectrum is evaluated as a discrete sum of the target training vectors:

$$\Rightarrow \hat{\mathbf{y}}_t = \sum_{k=1}^{N_y} w_{t|t}^k \mathbf{y}^k$$

Grid-based conversion

Sequential Estimation

- The posterior weights are sequentially evaluated using two stages:

1. Prediction: $w_{t|t-1}^k \approx \sum_{l=1}^{N_y} w_{t-1|t-1}^l p(\mathbf{y}_t = \mathbf{y}^k | \mathbf{y}_{t-1} = \mathbf{y}^l)$

2. Update: $w_{t|t}^k \approx \frac{w_{t|t-1}^k p(\mathbf{x}_t | \mathbf{y}_t = \mathbf{y}^k)}{\sum_{l=1}^{N_y} w_{t|t-1}^l p(\mathbf{x}_t | \mathbf{y}_t = \mathbf{y}^l)}$

- $p(\mathbf{y}_t = \mathbf{y}^k | \mathbf{y}_{t-1} = \mathbf{y}^l)$ - evidence probability
- $p(\mathbf{x}_t | \mathbf{y}_t = \mathbf{y}^k)$ - likelihood probability

Grid-based conversion

Evidence Modeling

- The evidence probability $p(\mathbf{y}_t = \mathbf{y}^k | \mathbf{y}_{t-1} = \mathbf{y}^l)$ expresses the transition probability from state \mathbf{y}^l to \mathbf{y}^k
- In natural speech, spectral feature vectors related to consecutive time frames are typically similar, but not identical.
- We model the transition probability as:

$$p(\mathbf{y}_t = \mathbf{y}^k | \mathbf{y}_{t-1} = \mathbf{y}^l) \propto \exp \left\{ -\frac{1}{2} \max \left(\frac{\text{MCD}(\mathbf{y}^k, \mathbf{y}^l)}{R_y}, 1 \right)^2 \right\}$$

- R_y – a parameter
- MCD – Mel Cepstral Distortion:

$$\text{MCD}(\mathbf{y}^k, \mathbf{y}^l) \triangleq \frac{10\sqrt{2}}{\ln 10} \|\mathbf{y}^k - \mathbf{y}^l\|_2$$

Grid-based conversion

Likelihood Modeling

- The likelihood probability $p(\mathbf{x}_t | \mathbf{y}_t = \mathbf{y}^k)$ expresses the probability of a source vector given target training vectors
- We model the likelihood probability as:

$$p(\mathbf{x}_t | \mathbf{y}_t = \mathbf{y}^k) \propto \sum_{m=1}^{N_x} p(\mathbf{x}^m | \mathbf{y}_t = \mathbf{y}^k) \exp \left\{ -\frac{\text{MCD}^2(\mathbf{x}_t, \mathbf{x}^m)}{2R_x^2} \right\}$$

- R_x – a parameter
- $p(\mathbf{x}^m | \mathbf{y}_t = \mathbf{y}^k)$ - the discrete likelihood

Grid-based conversion

Discrete Likelihood Modeling

- The discrete likelihood $p(\mathbf{x}^m | \mathbf{y}_t = \mathbf{y}^k)$ expresses the correspondence between the source and target training vectors
- We use a parallel (but need not be aligned) and phonetically labeled data set to model the discrete likelihood as:

$$p(\mathbf{x}^m | \mathbf{y}_t = \mathbf{y}^k) \propto \begin{cases} 1 & \mathbf{x}^m \text{ and } \mathbf{y}^k \text{ belong to the} \\ & \text{same mid-utterance} \\ 0 & \text{otherwise} \end{cases}$$

Grid-Based (GB) Conversion Algorithm Summary

- **Input:** a sequence of feature vectors
- **Initialization:** set the initial weights, $\left\{w_{t|t-1}^0\right\}_{k=1}^{N_y}$
- **Main Iteration:** for $t = 1, \dots, T$, perform the following steps:

1. Evaluate the prior weights: $\left\{w_{t|t-1}^k\right\}_{k=1}^{N_y}$
2. Evaluate the posterior weights: $\left\{w_{t|t}^k\right\}_{k=1}^{N_y}$
3. Obtain the converted spectra:

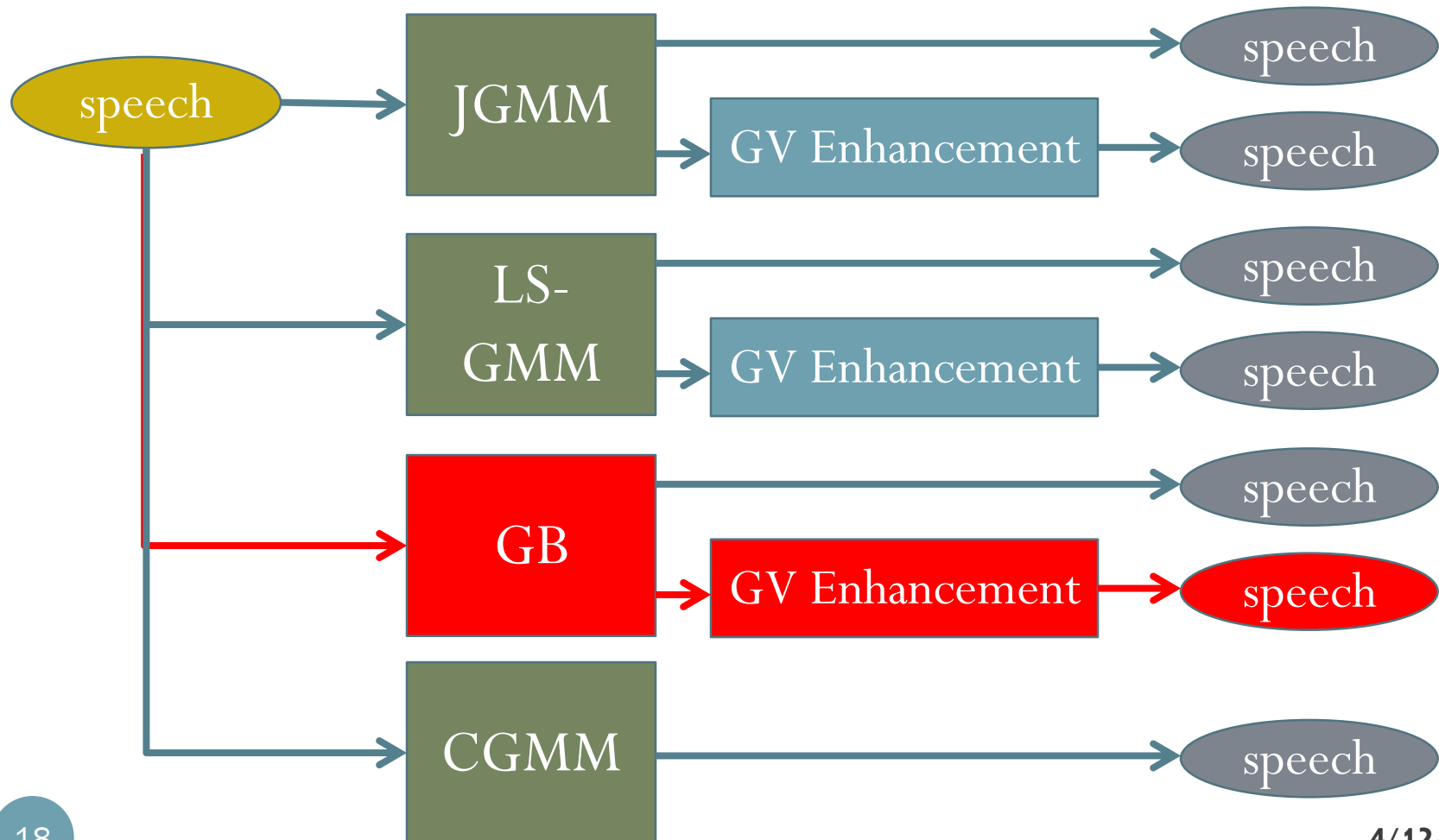
$$\hat{\mathbf{y}}_t = \sum_{k=1}^{N_y} w_{t|t}^k \mathbf{y}^k$$

- **Output:** a sequence of converted vectors $\boxed{\hat{\mathbf{y}}_{1:T}}$

Examined Methods

- Conversion Methods
 - Joint GMM (JGMM) [Kain & Macon, 1998]
 - Least squares GMM (LS-GMM) [Stylianou et. al., 1998]
 - Proposed GB
- GV-Enhancement Methods
 - Integrated during training –
Constrained GMM (CGMM) [Benisty and Malah, 2011]
 - A Post processing block [Benisty et. al., 2012]

Examined Setups

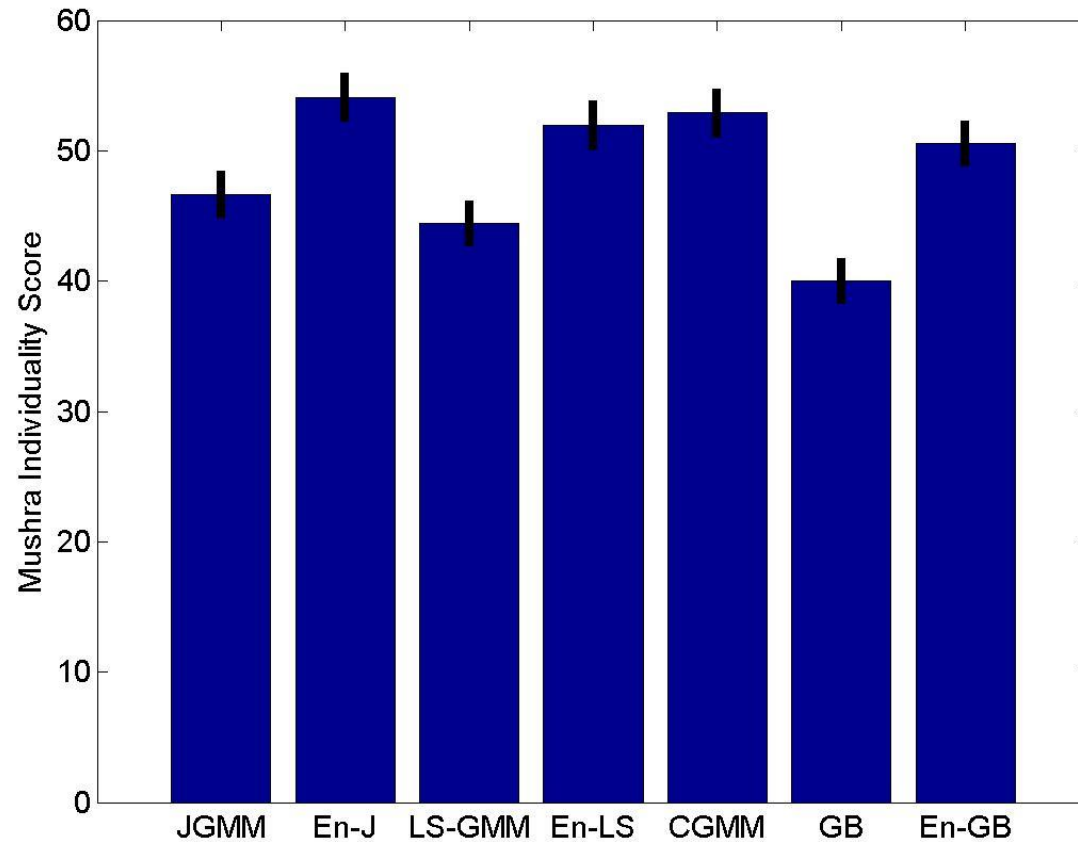


Training Period

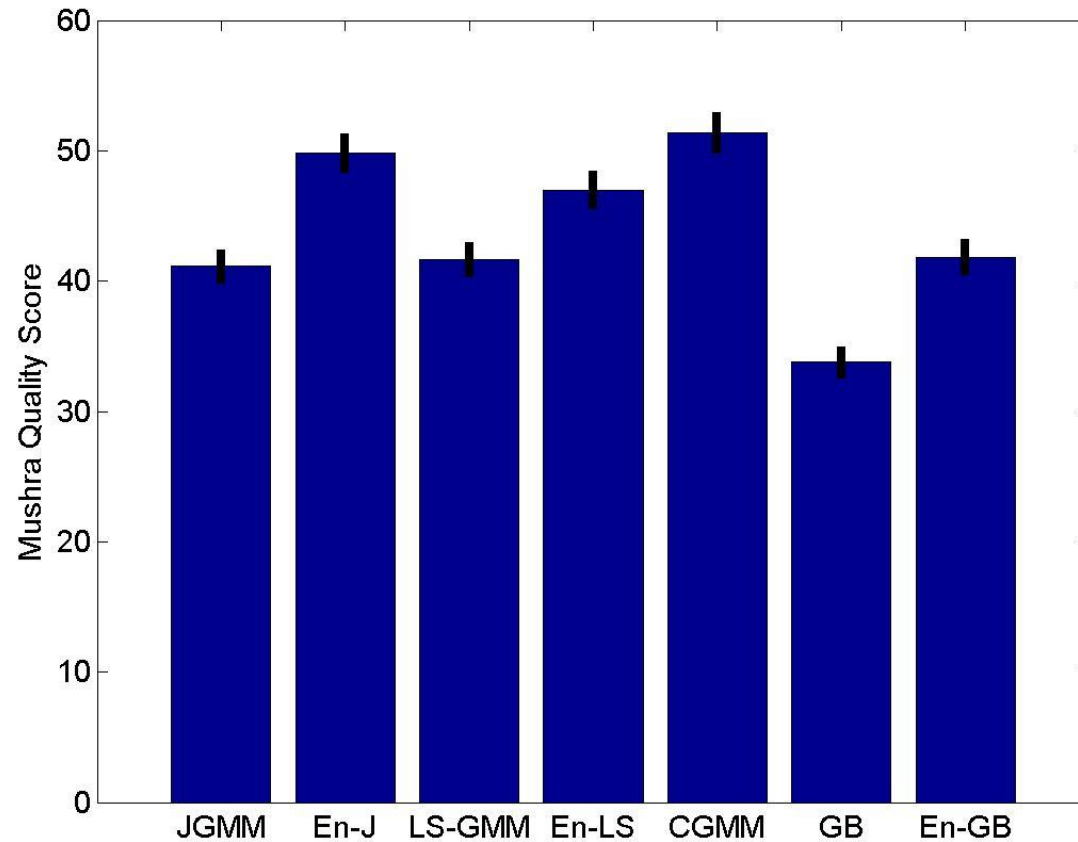
Method	Training time	Conversion time per frame
JGMM	7 h	11 msec
LS-GMM	8.5 h	11 msec
CGMM	11 h	11 msec
GB	10 sec	10 msec
GB enhancement	None in training	23 msec

Subjective Evaluations

Individuality Tests



Subjective Evaluations Quality Tests



Conclusion

- Easily trained using either small or large scale datasets
- Does require a parallel dataset but **without** time alignment
- Subjectively –
 - Comparable quality to the classical GMM-based methods (without enhancement)
 - Comparable individuality enhanced GMM-based methods
- Sound samples are available at:
<http://sipl.technion.ac.il/Info/hadas/sound-samples.htm>

Thank You
