# MODULAR GLOBAL VARIANCE ENHANCEMENT FOR VOICE CONVERSION SYSTEMS

Hadas Benisty, David Malah, and Koby Crammer

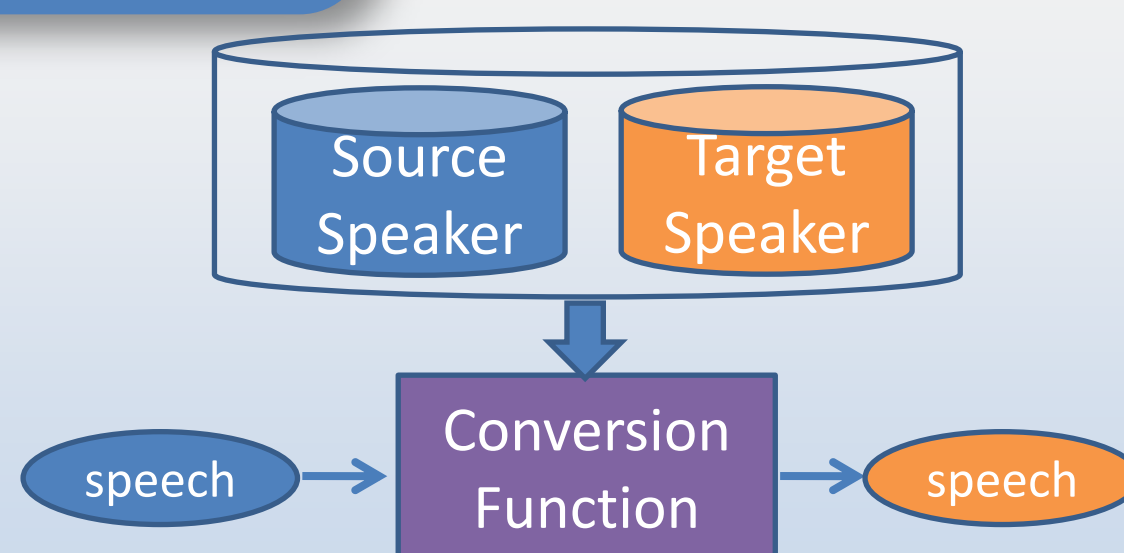EE Department, Technion – Israel Institute of Technology, Haifa, Israel

## Goal

- Many voice conversion methods produce muffled synthesized outputs due to over-smoothing of the converted spectra
- GV enhancement – used for muffling reduction and commonly applied as an integrated part of the conversion system
- **We propose a new modular method for GV enhancement, applied as a post-processing block**
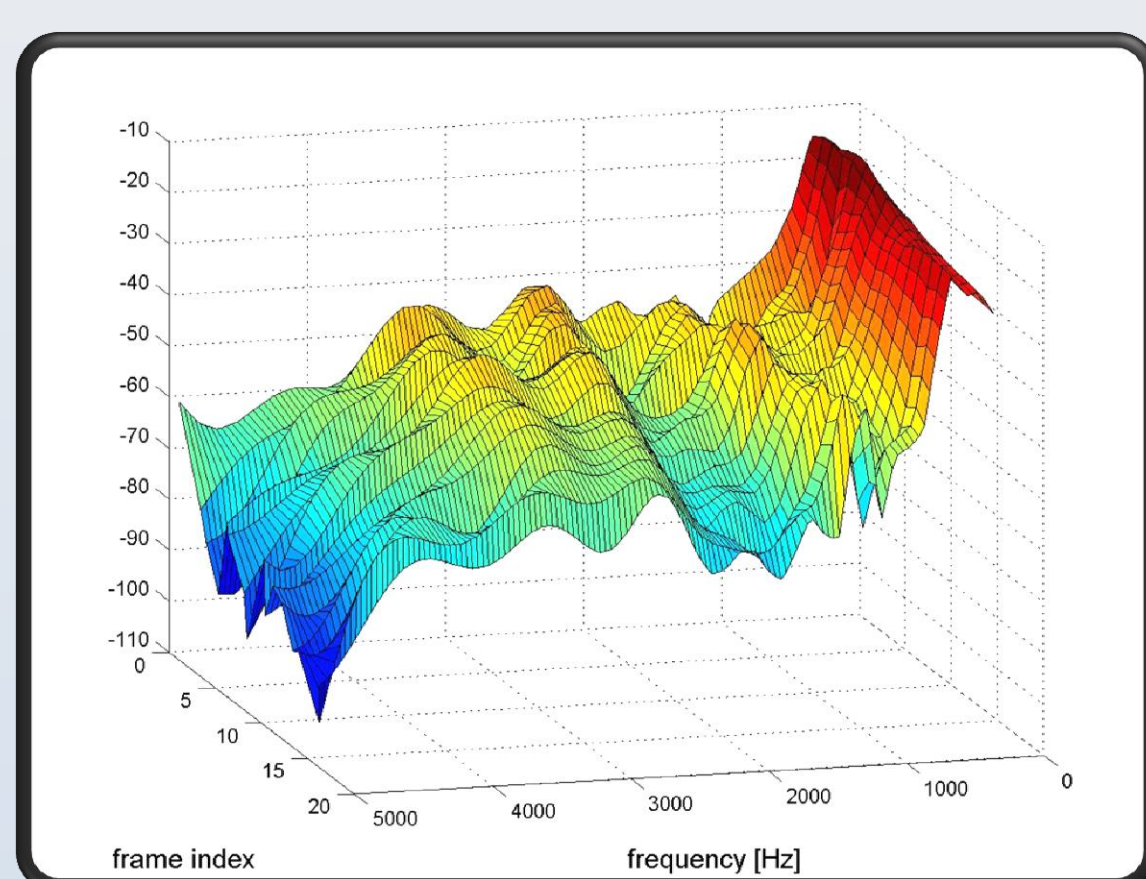
## Voice Conversion

- Transform a sentence said by a source speaker, to sound as if a target speaker had said it, based on pre-recorded training set

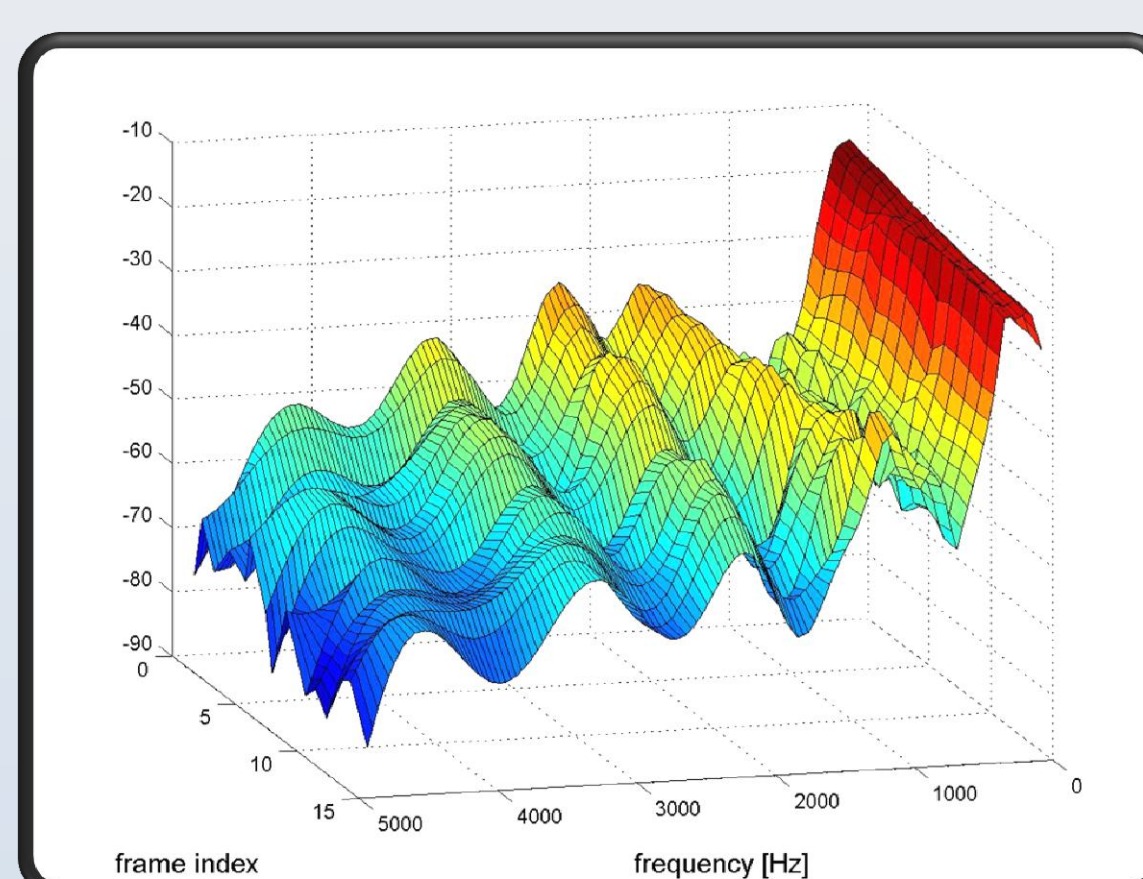Source Speaker — Target Speaker

speech → Conversion Function → speech

## Voice Conversion Using GMM

- Linear Conversion based on a Gaussian Mixture Model (GMM) [Stylianou, 1998], [Kain & Macon, 1998]
- A common approach for spectral conversion
- Minimizes the mean Log Spectral distortion (LSD) between converted feature vectors and target vectors
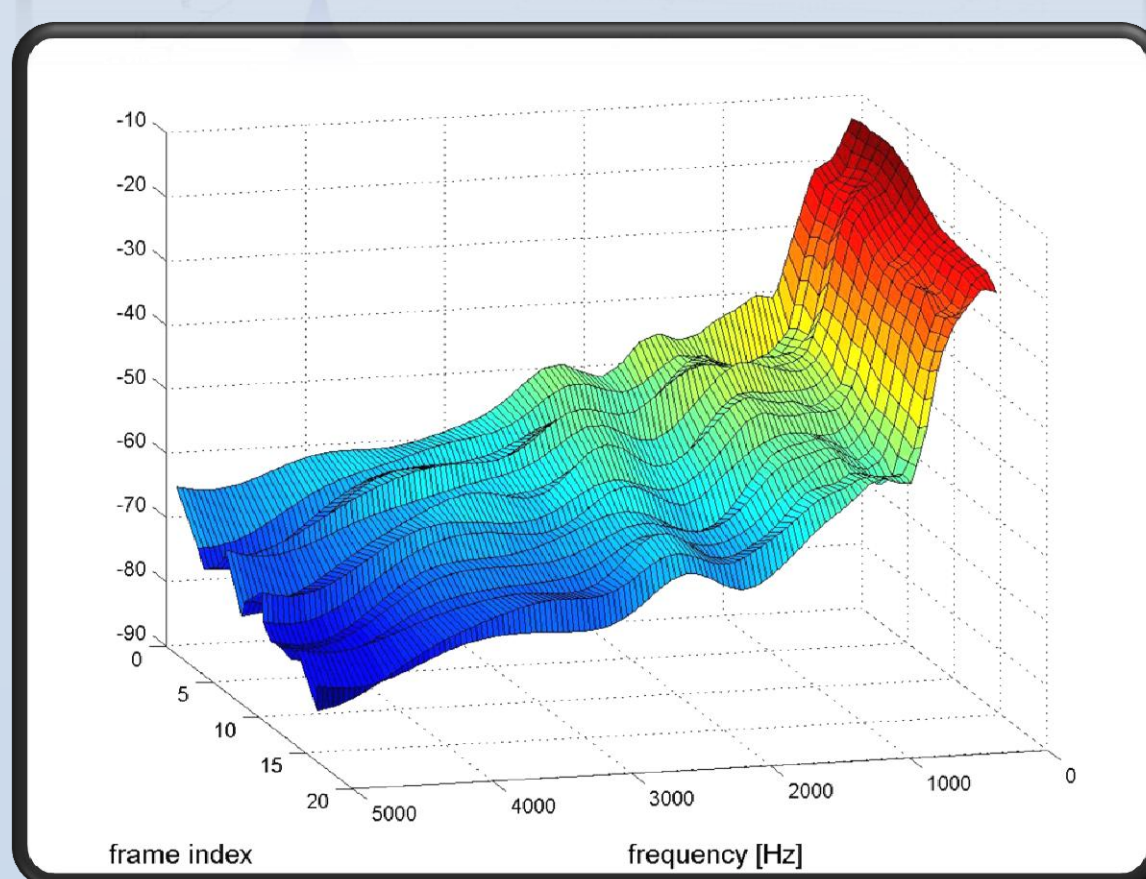- Characterized by smoothed spectral envelopes causing a **muffling effect**:
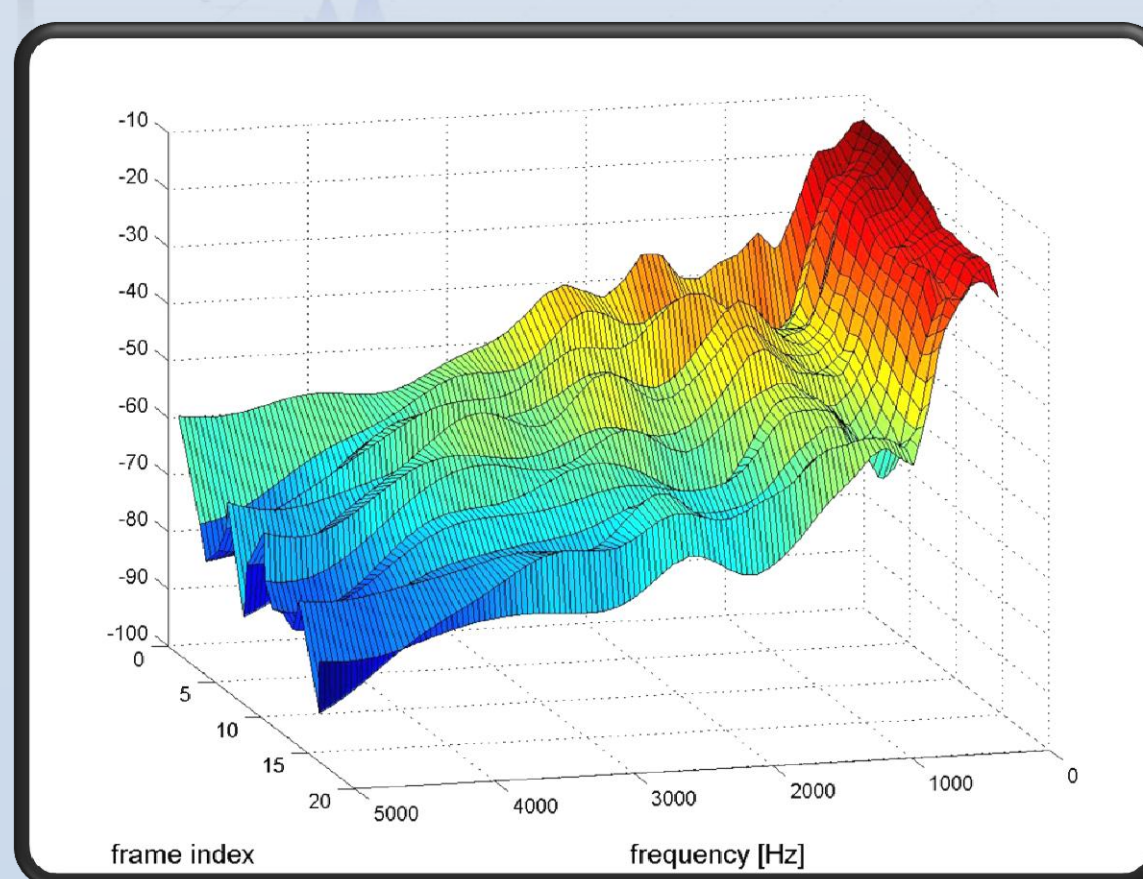
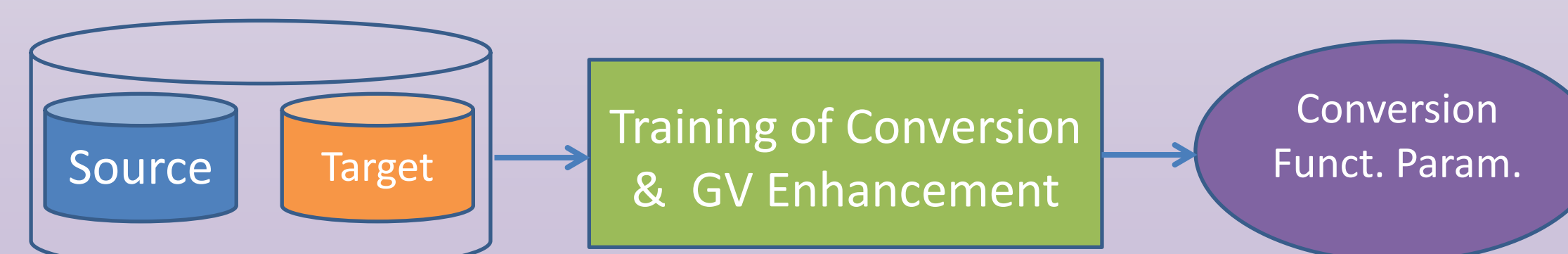### Spectral-Envelope Evolution in Time


Source Speaker


Target Speaker


Converted signal [Stylianou, 1998]


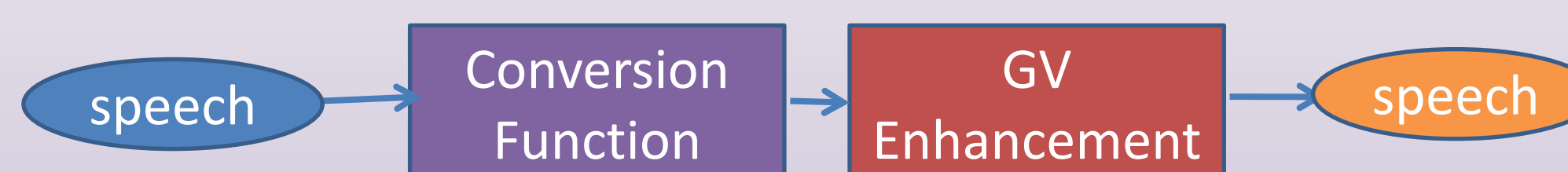LS-GMM followed by GV enhancement **(our work)**

## GV Enhancement

- GV enhancement methods have been proposed to overcome the muffling effect:
  - ML estimation [Toda et. al., 2007]
  - Constrained GMM (CGMM) [Benisty and Malah, 2011]
- These enhancement methods are integrated into the training process of the conversion

Source — Target → Training of Conversion & GV Enhancement → Conversion Funct. Param.

## Proposed Modular GV Enhancement

- **GV Enhancement Using an LSD Constraint**
  - Designed **independently** of any specific conversion scheme and applied as a **post-processing** block

  speech → Conversion Function → GV Enhancement → speech

  - The extent of GV enhancement is **controlled** by the allowed **spectral distance** the enhanced and the originally converted output, as specified by the user

## Experimental Results

- **Evaluated Methods**
  - GMM-based Conversion (LS-GMM) [Stylianou, 1998]
  - LS-GMM followed by **our** GV enhancement
  - CGMM [Benisty and Malah, 2011]

- **Objectively**
  - For a given mean LSD, CGMM leads to higher GV than our method
- **Subjectively**
  - Our method was selected by the majority of listeners as **better** than CGMM, both in terms of **quality** and **similarity** to the target

### Objective Evaluations

| Conversion Method | Mean LSD [dB] | Mean Norm. GV |
|---|---|---|
| LS-GMM | 6.2 | 0.1 |
| Enhanced $\theta_{LSD}=1dB$ | 6.4 | 0.2 |
| Enhanced $\theta_{LSD}=2dB$ | 6.7 | 0.3 |
| Enhanced $\theta_{LSD}=4dB$ | 7.3 | 0.4 |
| CGMM | 7.3 | 0.9 |

## GV Enhancement Using an LSD Constraint

- **Input**
  - A sequence of converted feature vectors $\tilde{\mathbf{Y}}_{1:T} \triangleq \left(\tilde{\mathbf{y}}_1, \quad \tilde{\mathbf{y}}_2, \quad ..., \quad \tilde{\mathbf{y}}_T\right)^T$
- **Output**
  - A sequence of **enhanced** feature vectors $\tilde{\mathbf{Z}}_{1:T} \triangleq \left(\tilde{\mathbf{z}}_1, \quad \tilde{\mathbf{z}}_2, \quad ..., \quad \tilde{\mathbf{z}}_T\right)^T$
- The enhanced sequence is the solution of:

$$\tilde{\mathbf{Z}}_{1:T} = \arg\max_{\mathbf{Z}_{1:T}} \mathrm{NGV}\{\mathbf{Z}_{1:T}\}$$
$$\text{s.t} \quad \overline{\mathrm{LSD}}\left(\mathbf{Z}_{1:T}, \tilde{\mathbf{Y}}_{1:T}\right) \le \theta_{LSD}$$

  - $\mathrm{NGV}\{\mathbf{Z}_{1:T}\}$ - the normalized GV of the sequence $\mathbf{Z}_{1:T}$, evaluated by:

$$\mathrm{NGV}\{\mathbf{Z}_{1:T}\} \triangleq \frac{1}{P}\sum_{p=1}^{P}\frac{\mathrm{Var}\{\tilde{\mathbf{Z}}_{1:T}(p)\}}{\mathrm{Var}\{\mathbf{Z}_{1:T}(p)\}}$$

  - $\overline{\mathrm{LSD}}\left(\mathbf{Z}_{1:T}, \tilde{\mathbf{Y}}_{1:T}\right)$ - mean Log spectral Distortion between the converted and enhanced sequences
  - $\theta_{LSD}$ - pre-set threshold value for the mean LSD in dB
- The solution is obtained with explicit terms for mean LSD and NGV

$$\overline{\mathrm{LSD}}\left(\tilde{\mathbf{Z}}_{1:T}, \tilde{\mathbf{Y}}_{1:T}\right) \approx \frac{\kappa}{T}\|\tilde{\mathbf{Z}}_{1:T} - \tilde{\mathbf{Y}}_{1:T}\|_{2,1} \qquad \kappa \triangleq 10\sqrt{2}/\ln 10$$

$$\mathrm{NGV}\{\tilde{\mathbf{Y}}_{1:T}\} = \frac{1}{P}\left\|\Delta_T \cdot \tilde{\mathbf{Y}}_{1:T} \cdot \mathbf{C}^{-\frac{1}{2}}\right\|_2^2 \qquad \Delta_T \triangleq \frac{1}{\sqrt{T}}\left(\mathbf{I}_{T\times T} - \frac{1}{T}\mathrm{ones}(T,T)\right)$$

$$\mathbf{C} \triangleq diag\left(\mathrm{Var}\{\mathbf{Y}(1)\}, \quad ..., \quad \mathrm{Var}\{\mathbf{Y}(P)\}\right) \quad \mathrm{Var}\{\mathbf{Y}(p)\} \text{ - GV of spectral features related to the target speaker}$$

### Subjective Evaluations


Quality - AB


Identity - XAB