



# STATISTICAL TEXT-TO-SPEECH SYNTHESIS WITH IMPROVED DYNAMICS

Stas Tiomkin<sup>†‡</sup> and David Malah<sup>†</sup>

<sup>†</sup>Department of Electrical Engineering, Technion - I.I.T., <sup>‡</sup>Speech Technologies Group, IBM, Israel.



## Main Text-To-Speech (TTS) Synthesis Methods

- Concatenative TTS - CTTS
  - Concatenation of natural samples from an inventory.
  - Advantages
    - \* Natural quality speech.
    - \* Low computational complexity
  - Disadvantages
    - \* Essential quality degradation with a footprint reduction.
    - \* Speaker dependent.
- Statistical TTS - STTS
  - Speech parameters generation from statistical models.
  - Advantages
    - \* Low footprint.
    - \* Voice modification.
  - Disadvantages
    - \* Insufficient speech features dynamics.
    - \* Muffled and buzzy speech.

## Conventional STTS

- Speech feature vector over an entire utterance
  - $\mathbf{c} = [\mathbf{c}_1, \mathbf{c}_2, \dots, \mathbf{c}_T]$
  - \*  $\mathbf{c}_i = [c_{i1}, c_{i2}, \dots, c_{id}]$
- Augmented space vector
  - $\mathbf{o} = [\mathbf{o}_1, \mathbf{o}_2, \dots, \mathbf{o}_T]$
  - \*  $\mathbf{o}_i = [\mathbf{c}_i, \Delta^1 \mathbf{c}_i, \Delta^2 \mathbf{c}_i]$
  - \*  $\Delta^{1,2} \mathbf{c}_i = \sum_{j=-L}^{L-1} \omega_j^{1,2} \mathbf{c}_{i+j}$
  - \*  $\mathbf{o} = \mathbf{W}\mathbf{c}$

- Statistical model
  - $P(\mathbf{o}) \sim \mathcal{N}(\mathbf{m}, \mathbf{U})$
- Optimal Solution

$$\mathbf{c}^{opt} = \underset{\mathbf{c}}{\operatorname{argmin}} \ln(P(\mathbf{o}))|_{\mathbf{o}=\mathbf{W}\mathbf{c}} = (\mathbf{W}^T \mathbf{U}^{-1} \mathbf{W})^{-1} \mathbf{W}^T \mathbf{U}^{-1} \mathbf{m}. \quad (1)$$

## Generated Speech Features

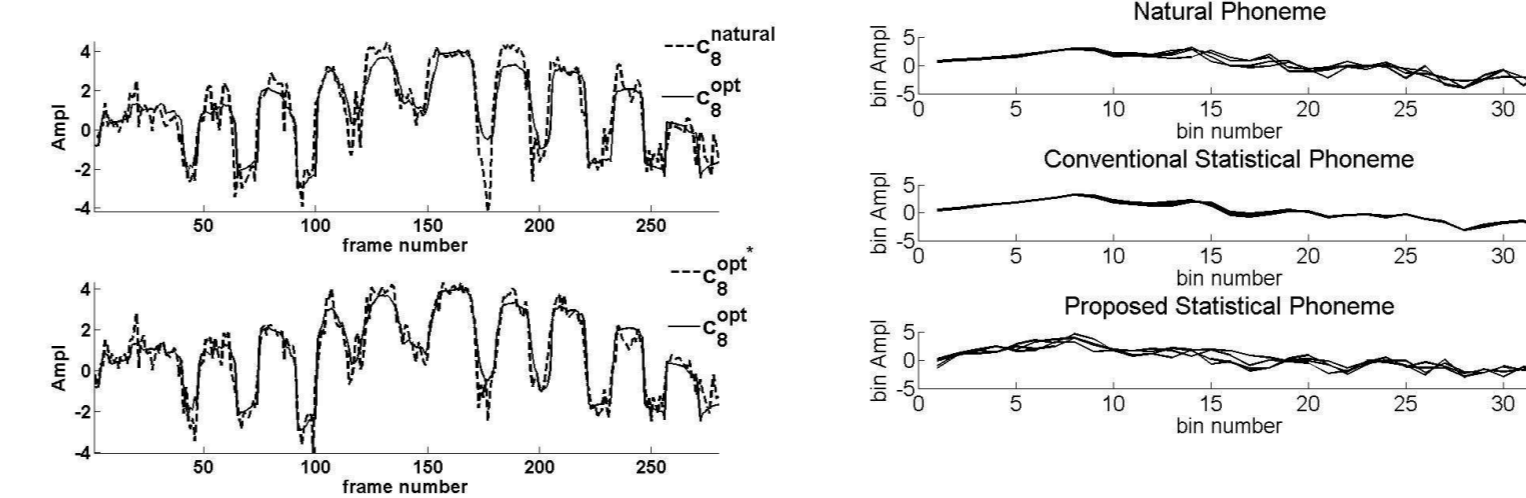


FIGURE 1: Evolution of a particular frequency component. The upper plot shows natural features (dashed line) and conventionally generated features (solid line). The bottom plot shows features generated by the proposed method (dashed line), described below, and conventionally generated features (solid line).

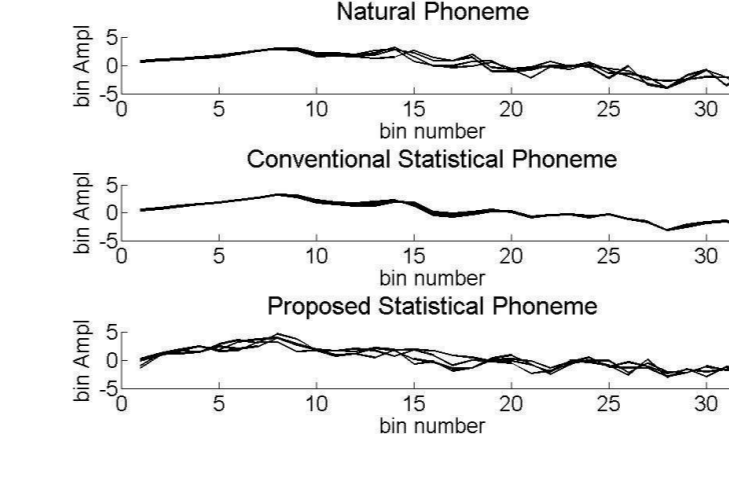


FIGURE 2: Evolution of frames within a particular phoneme in natural phoneme (top plot), in conventionally generated phoneme (middle plot), and in a phoneme generated by the proposed method (bottom plot). The conventionally generated phoneme lacks inter-frames dynamics, appeared in natural phoneme.

## Drawbacks of Conventional STTS

- Statistically generated components are over-smoothed.
- Over-smoothing causes muffled and buzzy speech.
- $\Delta^{1,2} \mathbf{c}$  do not model speech dynamics sufficiently.

## Previously Proposed Solution to Over-Smoothing

- Global Variance (GV) approach has been proposed<sup>a</sup>.
  - Applying a penalty for global variance reduction.
  - GV improves statistically generated speech quality.
  - GV is computationally complex.

<sup>a</sup>Tomoki Toda, Keiichi Tokuda, Speech parameter generation algorithm considering global variance for HMM-based speech synthesis, In INTERSPEECH-2005, 2801-2804

## Arranging Data as Quasi-Periodic Sequences

A phoneme of  $T_i$  frames is represented as a one-dimensional coefficients sequence of length  $dT_i$  with a period  $d$ .

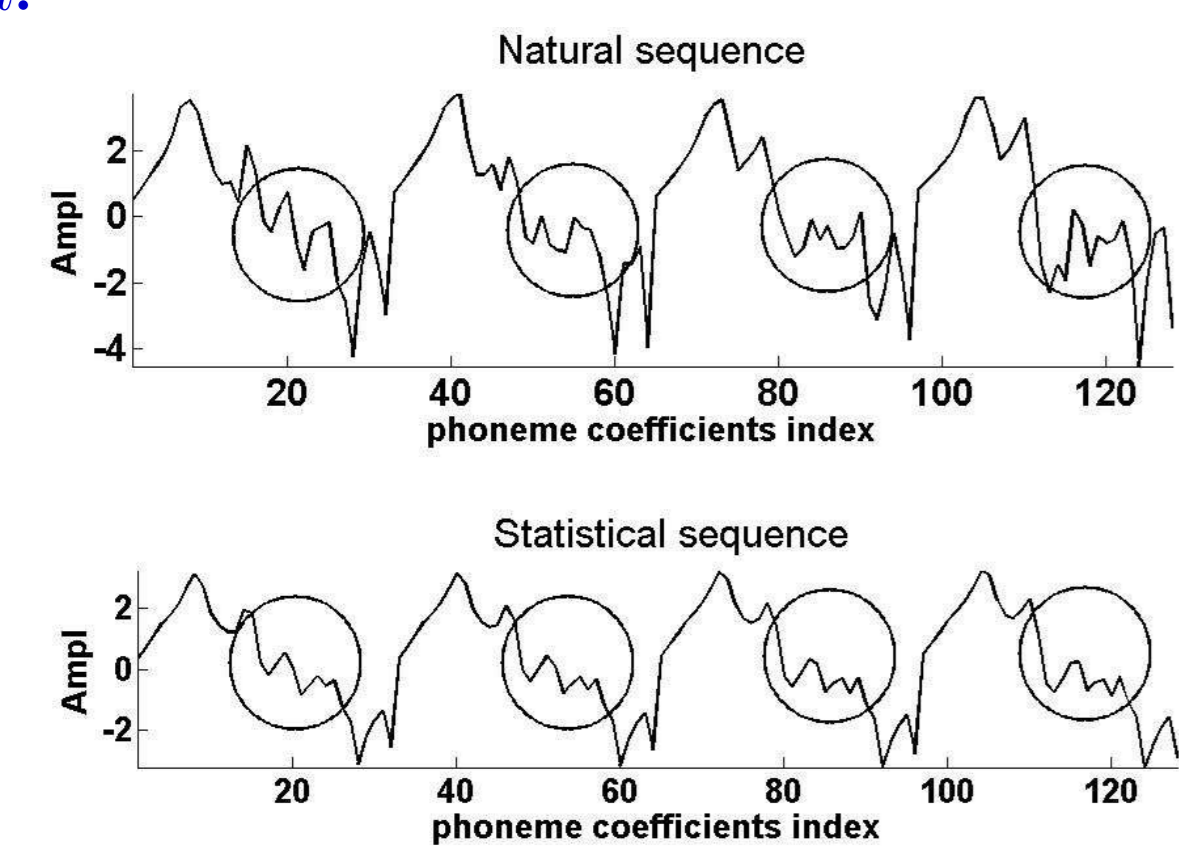
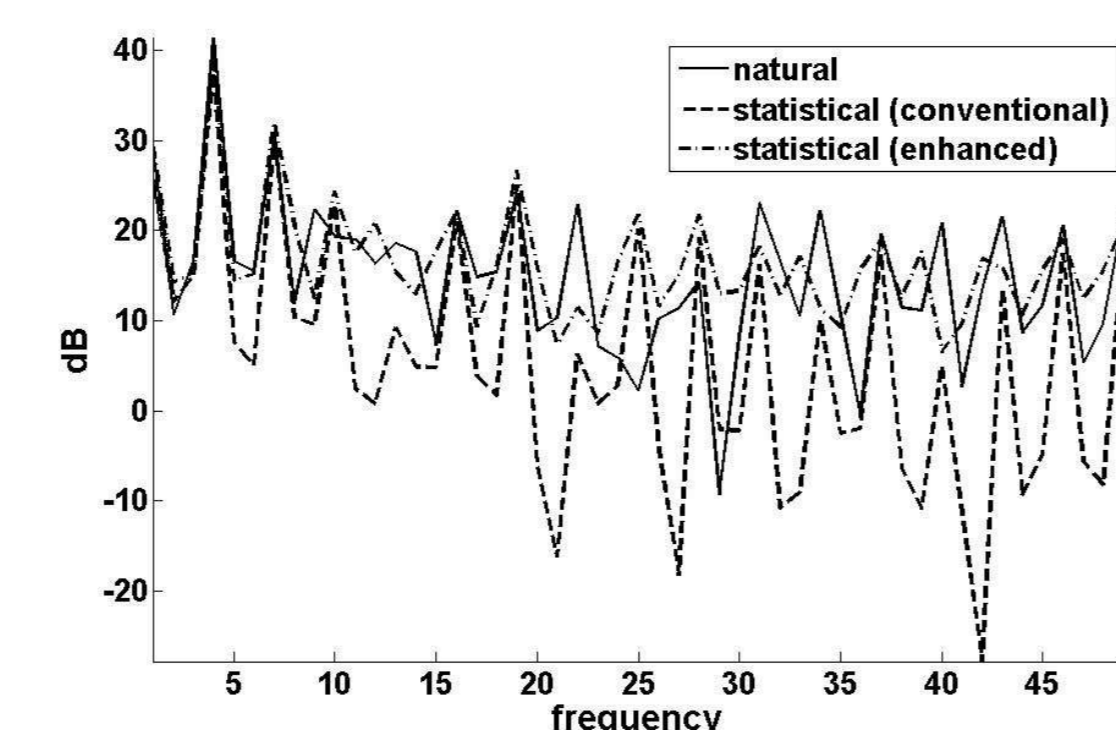


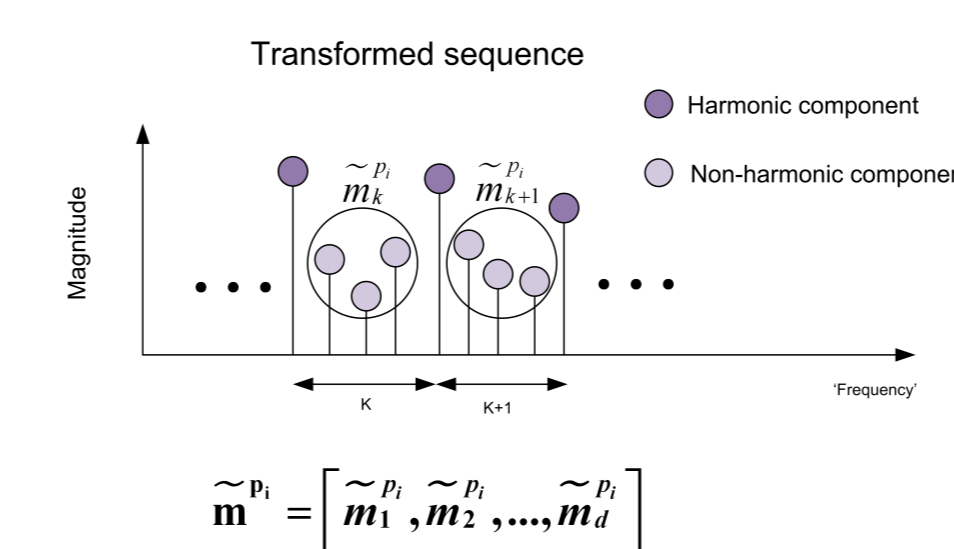
FIGURE 3: Features frames of a natural sequence (4 frames on the same plot) at the top; statistical sequence at the bottom. The frame-to-frame variations between circled regions demonstrate the low dynamics in the statistical sequence as compared to the natural sequence.

## Modeling Speech Features Dynamics in the Transform Domain

- Apply a DFT of length  $dT_i$  to the quasi-periodic sequence:
  - Harmonic frequencies:  $k = lT_i, l = 1, 2, \dots, d$
  - Non-harmonic frequencies:  $k = k + 1, k + 2, \dots, k + T_i - 1$ , where  $k = 1, T_i, 2T_i, \dots, (d-1)T_i$
- Non-harmonic content (NHC) in statistically generated phonemes is much lower, compared to NHC in natural phonemes.

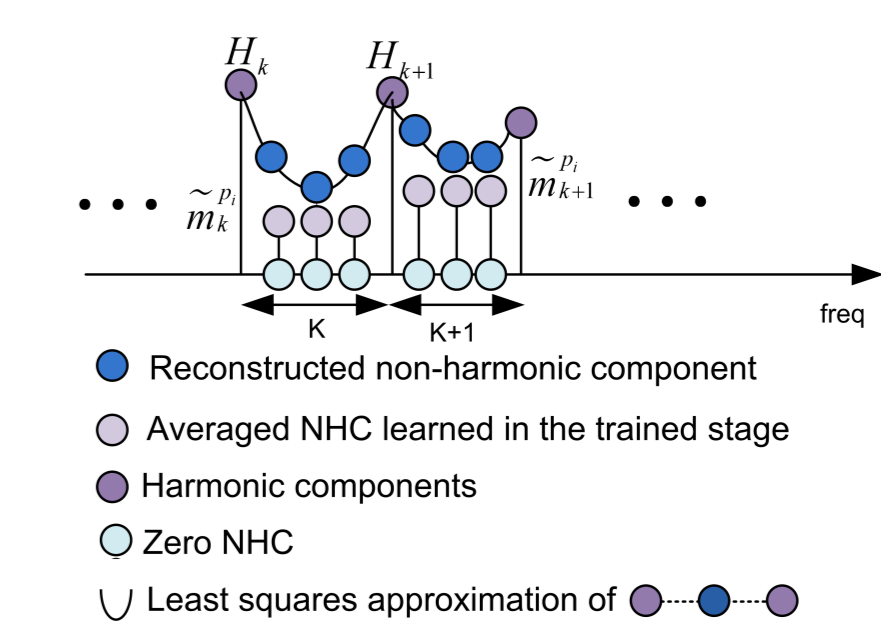


## Improving Speech Features Dynamics - Learning Non-Harmonic Components Statistics



1. Apply DFT to quasi-periodic sequences of natural segments from the phoneme  $p_i$ .
2. Extract NHC between every two harmonic components.
3. Compute  $\tilde{m}_k^{p_i}$  as an average value of NHC, pertaining to the  $k$ -th non-harmonic interval.
4. Compose  $\tilde{\mathbf{m}}^{p_i}$  for the phoneme  $p_i$ .

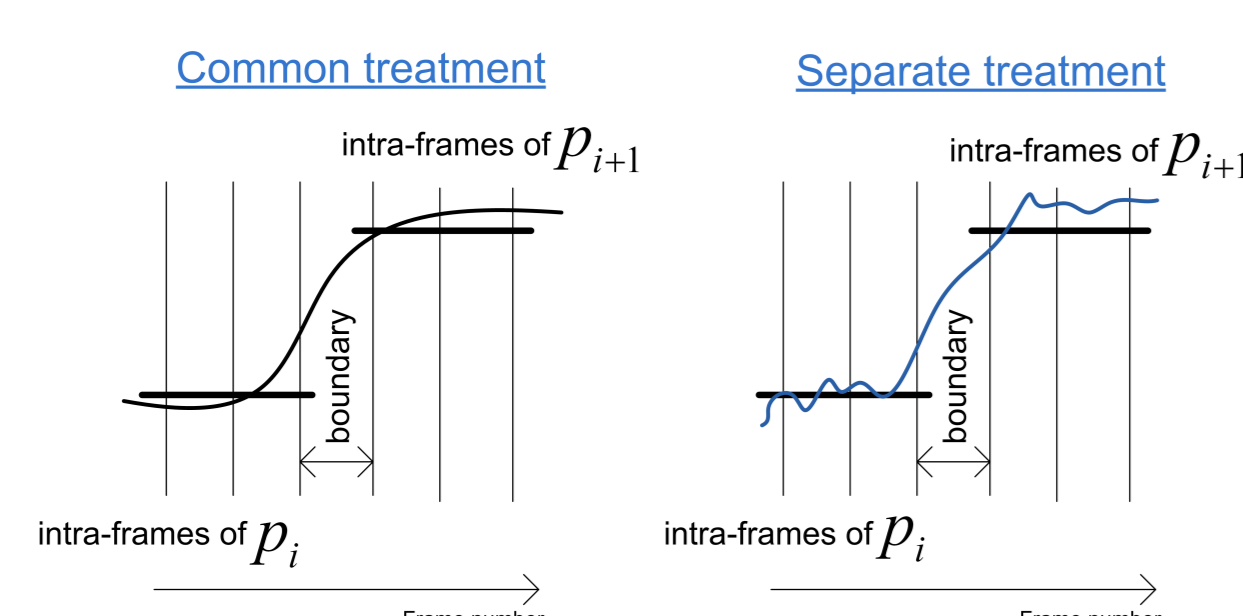
## Improving Speech Features Dynamics - Enhancing Non-Harmonic Component



1. Replicate the model mean  $\mathbf{m}_i$   $T_i$ -times.
2. Apply a DFT of length  $T_i$  to the periodic sequence.
  - NHC in the transformed sequence are exactly zero.
3. Find a least squares approximation by a second order polynomial of the points:  $H_k, \dots, \tilde{m}_k^{p_i}, \dots, H_{k+1}$ .

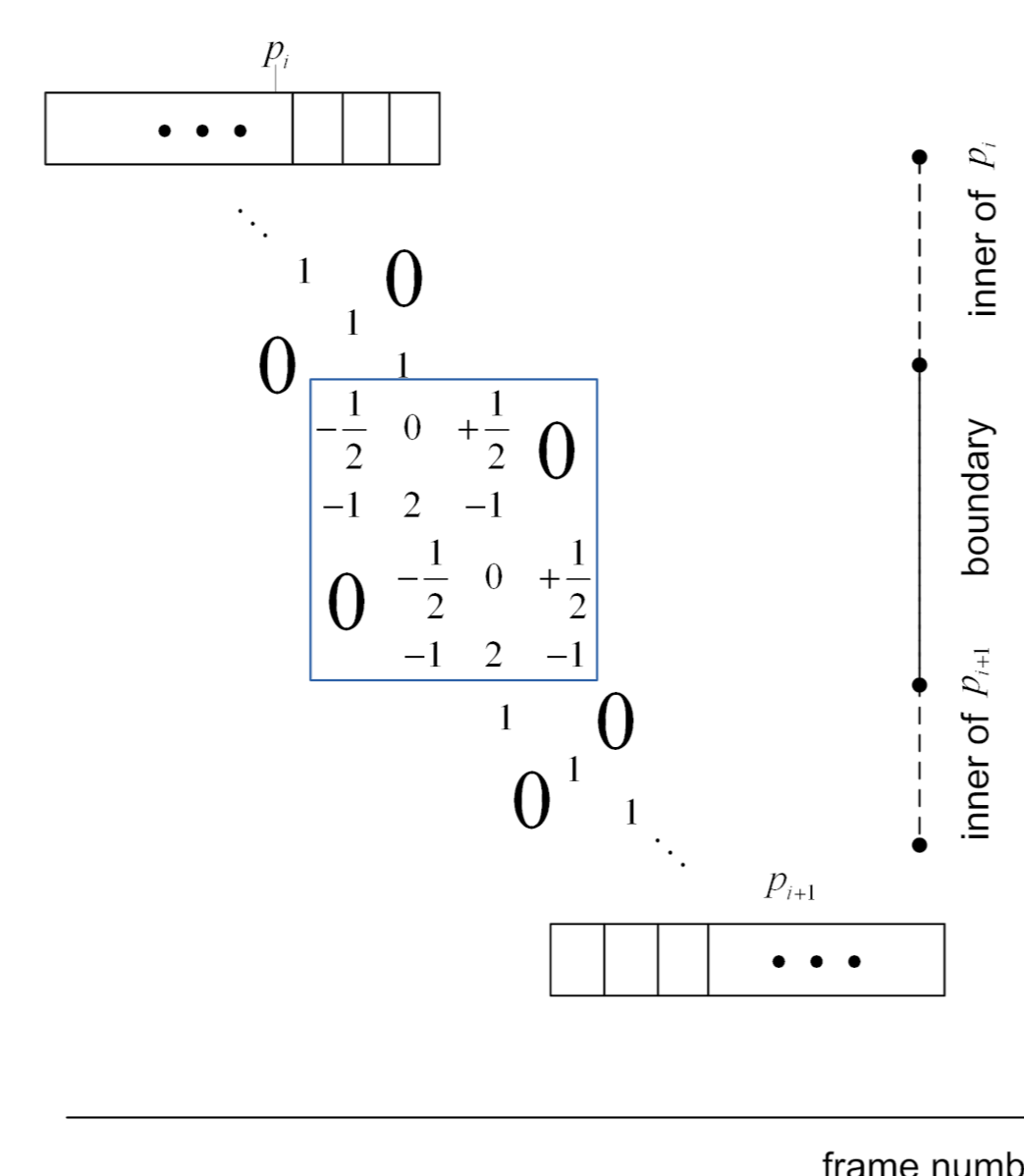
## Utterance-Level Synthesis

- Problem Setting
  - The inter phoneme boundaries and intra phoneme frames are treated equally.
  - Frames of both types are over-smoothed.
    - \* Smooth inter phoneme transitions are desired.
    - \* Intra phoneme smoothing reduces features dynamics.
  - We propose a separate treatment for inter phoneme boundaries and intra phoneme frames.



## Modified Linear Transformation

- Connect smoothly adjacent phonemes, using  $\Delta^{1,2} \mathbf{c}$  at boundaries.
- Enable enhanced features dynamics within phonemes.
- $\mathbf{W}$  in (1) is replaced by  $\tilde{\mathbf{W}}$ :



## Synthesis scheme

- Given an utterance, composed of  $K$  phonemes,  $p_1, p_2, \dots, p_K$ , having lengths  $T_1, T_2, \dots, T_K$
- Enhance the non-harmonic components in each phoneme
  - Construct model mean and model covariance matrix:
    - $\tilde{\mathbf{m}} = \left[ \underbrace{m_{p_1}^{stt}}_{\times T_1}, \underbrace{m_{p_1}^{\Delta^1}}_{\times T_1}, \underbrace{m_{p_1}^{\Delta^2}}_{\times T_1}, \dots, \underbrace{m_{p_K}^{\Delta^1}}_{\times T_K}, \underbrace{m_{p_K}^{\Delta^2}}_{\times T_K}, \underbrace{m_{p_K}^{stt}}_{\times T_K} \right]$
    - \*  $m_{p_i}^{stt}$  is the static features mean of  $p_i$
    - \*  $m_{p_i}^{\Delta^{1,2}}$  are the dynamic features mean of  $p_i$
    - $\tilde{\mathbf{U}} = \operatorname{diag} \left[ \underbrace{U_{p_1}^{stt}}_{\times T_1}, \underbrace{U_{p_1}^{\Delta^1}}_{\times T_1}, \underbrace{U_{p_1}^{\Delta^2}}_{\times T_1}, \dots, \underbrace{U_{p_K}^{\Delta^1}}_{\times T_K}, \underbrace{U_{p_K}^{\Delta^2}}_{\times T_K}, \underbrace{U_{p_K}^{stt}}_{\times T_K} \right]$
    - \*  $U_{p_i}^{stt}$  is the static features covariance of  $p_i$
    - \*  $U_{p_i}^{\Delta^{1,2}}$  are the dynamic features covariances of  $p_i$
  - Find the optimal solution, using  $\tilde{\mathbf{W}}, \tilde{\mathbf{m}}$  and  $\tilde{\mathbf{U}}$  in (1):
    - $\tilde{\mathbf{c}}^{opt} = (\tilde{\mathbf{W}}^T \tilde{\mathbf{U}}^{-1} \tilde{\mathbf{W}})^{-1} \tilde{\mathbf{W}}^T \tilde{\mathbf{U}}^{-1} \tilde{\mathbf{m}}$ .

## Results & Conclusions

- Subjective evaluation
  - Comparison test 'A vs B'
    - \* 'A' - conventional approach
    - \* 'B' - proposed approach
  - Ten arbitrary utterances were generated, using each approach.
  - Twelve listeners, graduate and undergraduate students, participated.
  - 81.7% of listeners preferred 'B' utterances.
- The proposed method:
  - improves speech features dynamics.
  - is not computationally complex.
  - is a data-driven approach.