

Probabilistic Feature Mapping Based on Trajectory HMMs

Heiga Zen Yoshihiko Nankaku Keiichi Tokuda (Nagoya Institute of Technology)

1. Introduction

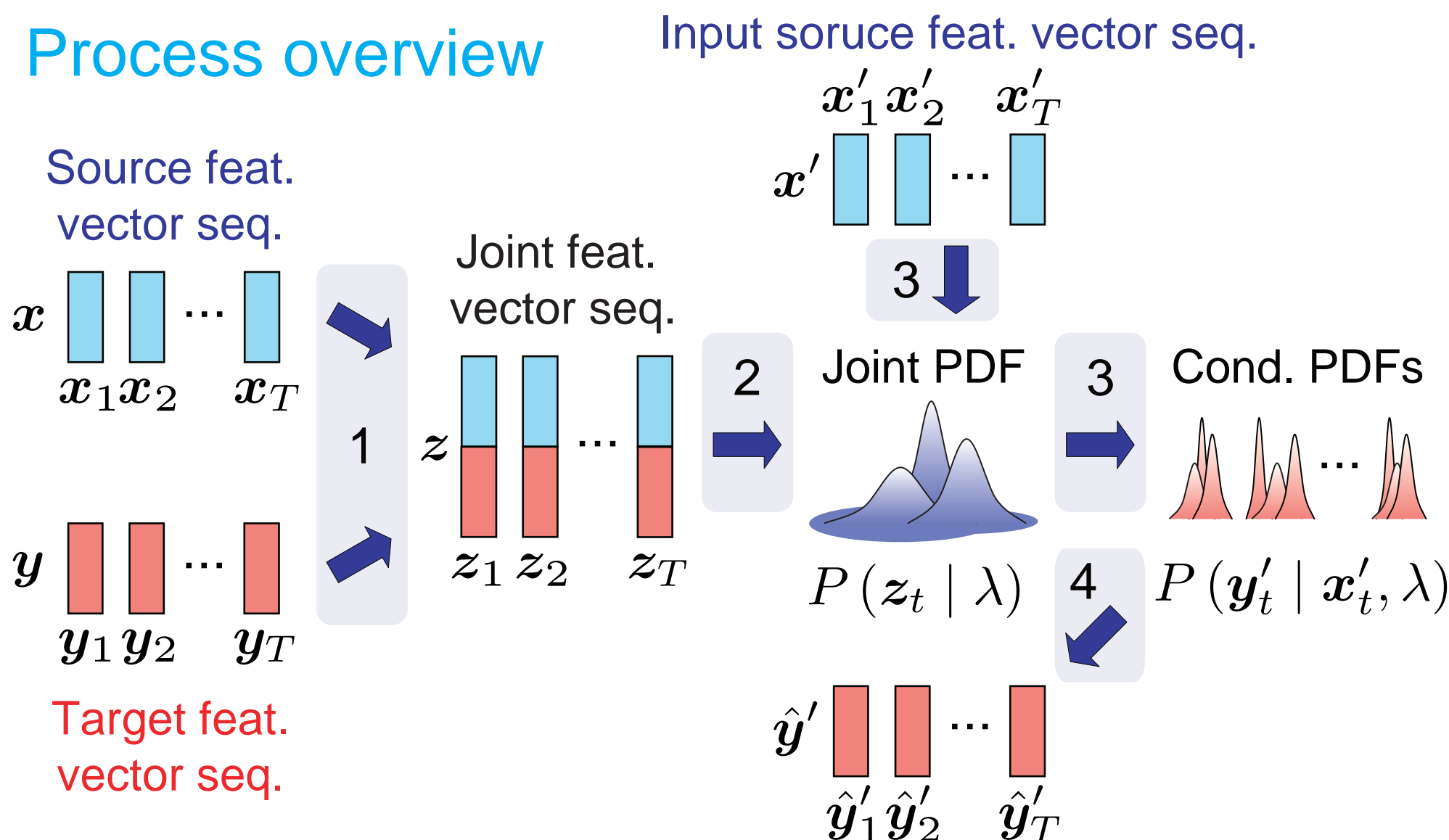
GMM / HMM-based feature mapping

- Continuously transform src. features into tgt. ones
 - Model joint PDF betw. src. & tgt. by GMM or HMMs
 - Map src. feats. into tgt. ones via conditional PDFs
- Applications
 - Voice conversion [Stylianou;'98, Kain;'98]
 - Acoustic-articulatory inv. mapping [Shiga;'04, Toda;'08]
 - Noise compensation [Droppo;'02, Cui;'07]
- Problem
 - Frame-by-frame mapping
 - ⇒ Inappropriate dynamic characteristics
 - Mapping with dynamic feature constraints [Toda;'07]
 - ⇒ Improper in the sense of statistical modeling

Trajectory GMM / HMM-based feature mapping

- Trajectory HMM [Zen;'07]
 - Impose explicit relationships betw. static & dyn. feats.
 - ⇒ HMM is reformulated as a trajectory model
 - Avoid limitations of HMMs w/o additional parameters
 - * Frame-wise conditional independence assumption
 - * Piece-wise constant statistics within an HMM state
- Feature mapping based on trajectory HMMs / GMMs
 - Entire utterance-level mapping
 - ⇒ Appropriate static & dynamic characteristics
 - Using dyn. feat. constraints in both training & mapping
 - ⇒ Make training & mapping consistent

2. GMM-based mapping [Kain;'98]

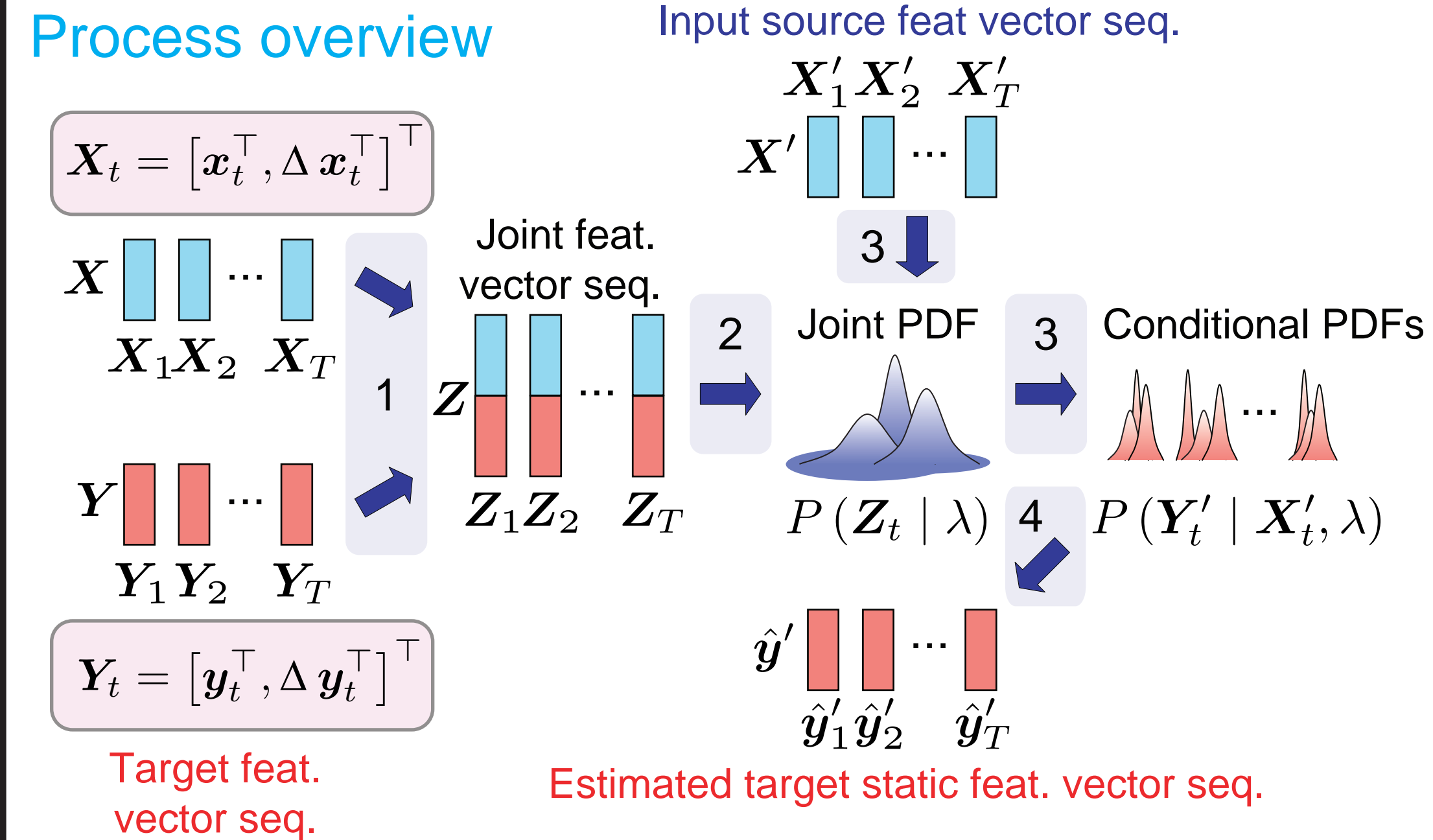


- Make joint feat. vector z_t from src. & tgt. vectors x_t & y_t
- Model **frame-level** joint PDF $P(z_t | \lambda)$ by GMM
- Convert joint PDF $P(z_t | \lambda)$ to cond. PDF $P(y'_t | x'_t, \lambda)$
- Estimate \hat{y}'_t from conditional PDF by MMSE

$$\hat{y}'_t = \sum_{i=1}^N \gamma_i \left[\mu_i^{(y)} + \Sigma_i^{(yx)} \Sigma_i^{(xx)^{-1}} (x'_t - \mu_i^{(x)}) \right]$$

Each frame is transformed independently
⇒ Mapped features are sometimes discontinuous

3. GMM-based mapping with dyn. feats. [Toda;'07]



- 1~3 are the same as those of mapping w/o dynamic features
- Estimate \hat{y}' from conditional PDF by ML **under the constraints between static & dynamic features ($Y' = W y'$)**

$$\hat{y}' = \arg \max_{Y'} P(Y' | X', \lambda) = \arg \max_{y'} P(W y' | X', \lambda) \approx (W^T D_q^{(y)^{-1}} W)^{-1} W^T D_q^{(y)^{-1}} E_q^{(y)} \quad (\text{Viterbi approx.})$$

$$E_q^{(y)} = [E_{tq_1}^{(y)^T}, E_{tq_2}^{(y)^T}, \dots, E_{tq_T}^{(y)^T}]^T \quad q = \{q_1, q_2, \dots, q_T\}$$

$$D_q^{(y)} = \text{diag} [D_{q_1}^{(y)}, D_{q_2}^{(y)}, \dots, D_{q_T}^{(y)}]$$

$$E_{ti}^{(y)} = \mu_i^{(y)} + \Sigma_i^{(yx)} \Sigma_i^{(xx)^{-1}} (X'_t - \mu_i^{(x)})$$

$$D_i^{(y)} = \Sigma_i^{(yy)} - \Sigma_i^{(yx)} \Sigma_i^{(xx)^{-1}} \Sigma_i^{(xy)}$$

- Entire utterance-level mapping
 - ⇒ Mapped feats. has proper static & dyn. characteristics
- Dynamic feature constraints are used only on mapping
 - ⇒ Inconsistency between training & mapping

4. Trajectory HMM / Trajectory GMM

Trajectory HMM [Zen;'07]

- Derived from HMMs with explicit dyn. feat. constraints
- Underlying generative model of HMM-based speech synth.

$$P(z | \lambda) = \sum_{\mathbf{q}} P(\mathbf{q} | \lambda) P(z | \mathbf{q}, \lambda) \quad \mu_{\mathbf{q}} = [\mu_{q_1}^T, \dots, \mu_{q_T}^T]^T$$

$$P(z | \mathbf{q}, \lambda) = \mathcal{N}(z; \bar{z}_{\mathbf{q}}, P_{\mathbf{q}}) \quad \mu_i = [\mu_i^{(x)^T}, \mu_i^{(y)^T}]^T$$

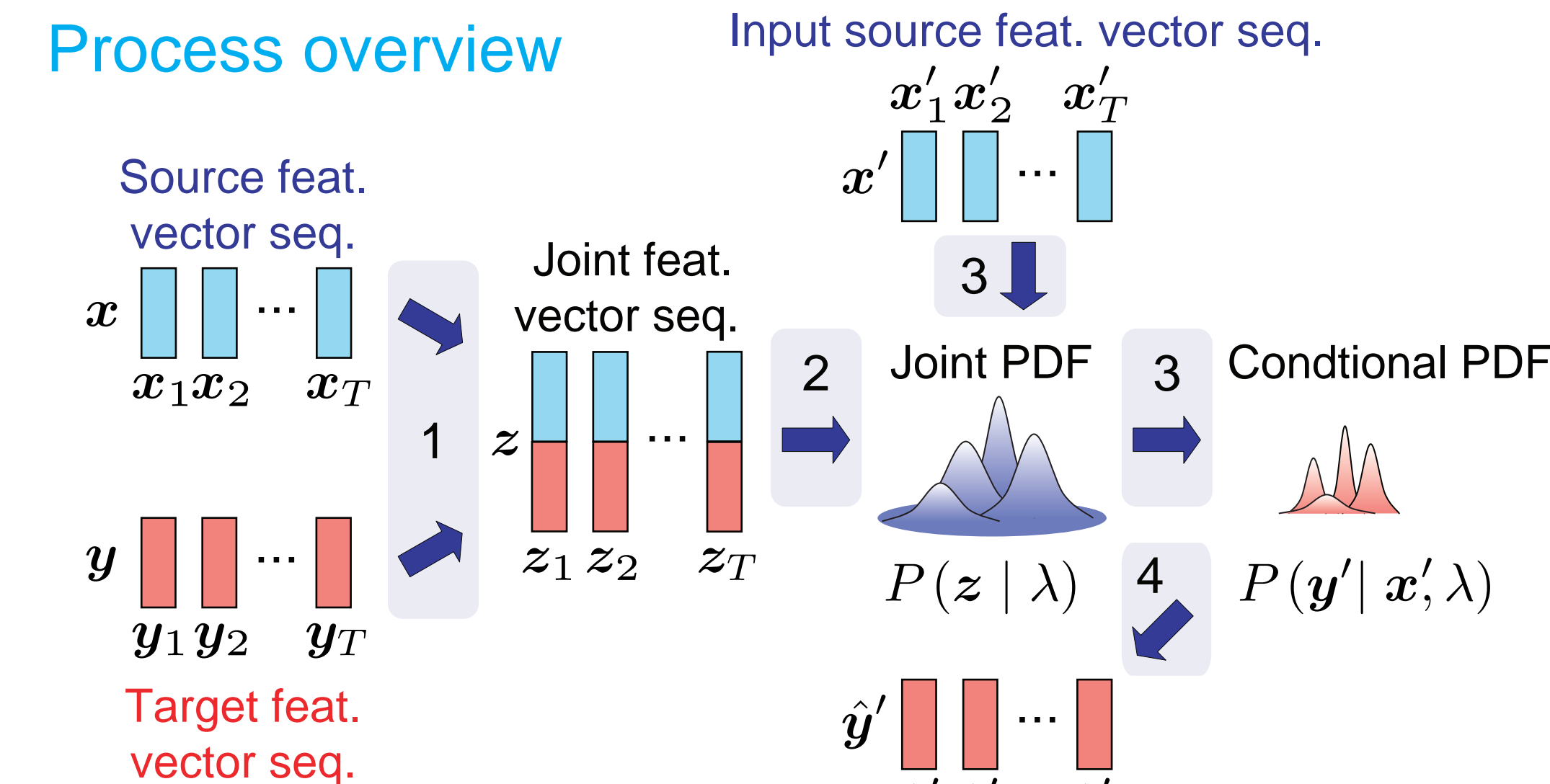
$$R_{\mathbf{q}} \bar{z}_{\mathbf{q}} = r_{\mathbf{q}} \quad \Sigma_{\mathbf{q}}^{-1} = \text{diag} [\Sigma_{q_1}^{-1}, \dots, \Sigma_{q_T}^{-1}]$$

$$R_{\mathbf{q}} = W^T \Sigma_{\mathbf{q}}^{-1} W = P_{\mathbf{q}}^{-1} \quad \Sigma_i^{-1} = \begin{bmatrix} \Omega_i^{(xx)} & \Omega_i^{(xy)} \\ \Omega_i^{(yx)} & \Omega_i^{(yy)} \end{bmatrix}$$

$$r_{\mathbf{q}} = W^T \Sigma_{\mathbf{q}}^{-1} \mu_{\mathbf{q}}$$

- $\bar{z}_{\mathbf{q}}$ is determined as a *smooth trajectory*
 - ⇒ Speech dynamics can be modeled explicitly
- Intra & inter-frame covariance matrix $P_{\mathbf{q}}$ is *full*
 - ⇒ Capture both intra & inter-frame correlations

5. Trajectory GMM (HMM)-based mapping



- Make joint feat. vec. seq. z from src. & tgt. vec. seqs. x & y

$$z = [z_1^T, z_2^T, \dots, z_T^T]^T \quad z_t = [x_t^T, y_t^T]^T$$

- Model **utterance-level** joint PDF $P(z | \lambda)$ by trajectory GMM

$$P(z | \lambda) = \sum_{\mathbf{q}} P(\mathbf{q} | \lambda) P(z | \mathbf{q}, \lambda)$$

$$P(z | \mathbf{q}, \lambda) = \mathcal{N}(z; \bar{z}_{\mathbf{q}}, P_{\mathbf{q}})$$

- Convert joint PDF $P(z | \lambda)$ to conditional PDF $P(y' | x', \lambda)$

$$P(z | \mathbf{q}, \lambda) = \mathcal{N} \left(\begin{bmatrix} x \\ y \end{bmatrix}; \begin{bmatrix} \bar{x}_{\mathbf{q}} \\ \bar{y}_{\mathbf{q}} \end{bmatrix}, \begin{bmatrix} P_{\mathbf{q}}^{(xx)} & P_{\mathbf{q}}^{(xy)} \\ P_{\mathbf{q}}^{(yx)} & P_{\mathbf{q}}^{(yy)} \end{bmatrix} \right)$$

$$P(y' | x', \lambda) = \sum_{\mathbf{q}} \gamma_{\mathbf{q}} \cdot P(y' | x', \mathbf{q}, \lambda)$$

$$P(y' | x', \mathbf{q}, \lambda) = \mathcal{N}(y'; \tilde{y}_{\mathbf{q}}, \tilde{P}_{\mathbf{q}}^{(yy)})$$

$$\tilde{y}_{\mathbf{q}} = \bar{y}_{\mathbf{q}} + P_{\mathbf{q}}^{(yx)} C_{\mathbf{q}}^{(xx)} (x' - \bar{x}_{\mathbf{q}})$$

$$\tilde{P}_{\mathbf{q}}^{(yy)} = P_{\mathbf{q}}^{(yy)} - P_{\mathbf{q}}^{(yx)} C_{\mathbf{q}}^{(xx)} P_{\mathbf{q}}^{(xy)}$$

- Estimate \hat{y}' from conditional PDF by MMSE or MAP

– MMSE

$$\hat{y}'_{\text{MMSE}} = \sum_{\mathbf{q}} \gamma_{\mathbf{q}} \left[\tilde{y}_{\mathbf{q}} + P_{\mathbf{q}}^{(yx)} C_{\mathbf{q}}^{(xx)} (x' - \bar{x}_{\mathbf{q}}) \right]$$

– MAP

$$\hat{y}'_{\text{MAP}} = \left(\sum_{\mathbf{q}} \gamma'_{\mathbf{q}} \cdot \tilde{P}_{\mathbf{q}}^{(yy)^{-1}} \right)^{-1} \sum_{\mathbf{q}} \gamma'_{\mathbf{q}} \cdot \tilde{P}_{\mathbf{q}}^{(yy)^{-1}} \left[\tilde{y}_{\mathbf{q}} + P_{\mathbf{q}}^{(yx)} C_{\mathbf{q}}^{(xx)} (x' - \bar{x}_{\mathbf{q}}) \right]$$

- Entire utterance-level transform
 - ⇒ Mapped feats. have proper static & dyn. characteristics
- Dyn. feat. constraints are used in both training & mapping
 - ⇒ Make training & mapping consistent

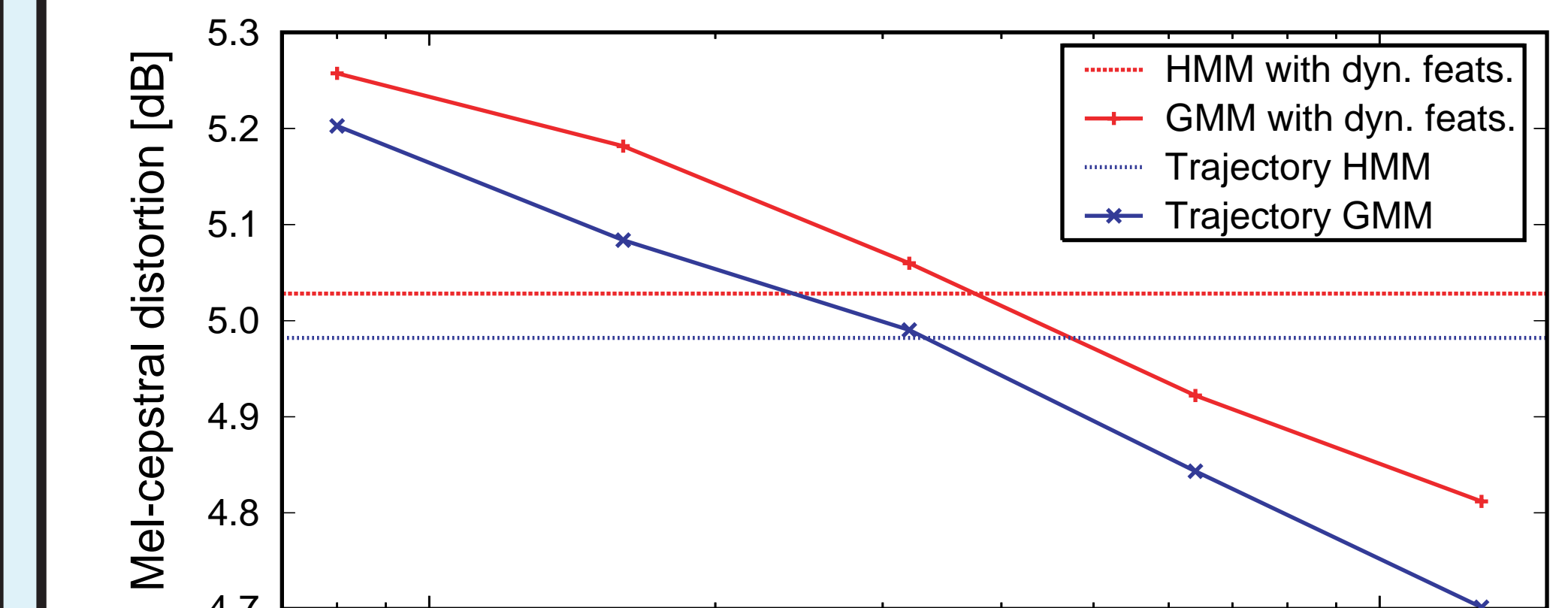
6. Voice conversion experiment

Conditions

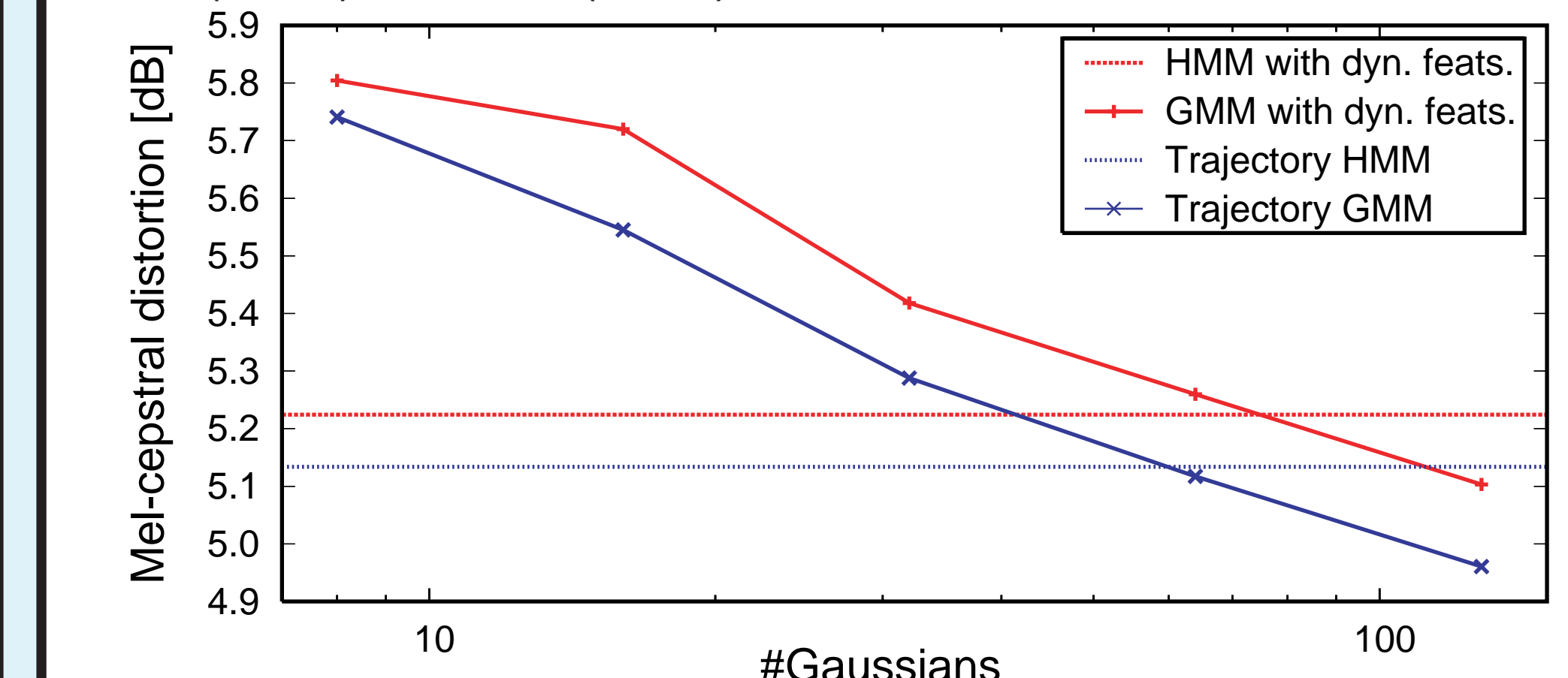
Database	CMU ARCTIC speech database
Training data	Speakers BDL, CLB, RMS, & SLT, first 593 utts. (Mapping: CLB→SLT, BDL→RMS)
Test data	Last 40 utterances
Sampling freq.	16 kHz
Analysis win.	25-ms Blackman window / 5-ms shift
Feature vec.	0~24 order Mel-cepstral coefficients, Δ & $\Delta\Delta$
Topology	(Trajectory) HMM: 3-state, left-to-right no-skip, monophone 1-mix (121 states) (Trajectory) GMM: 128-mix
Mapping	Spectrum: GMM (HMM) with dyn. or proposed F0: Linearly transformed in the log domain

Objective evaluation (mel-cepstral distortions)

CLB (female) → SLT (female)

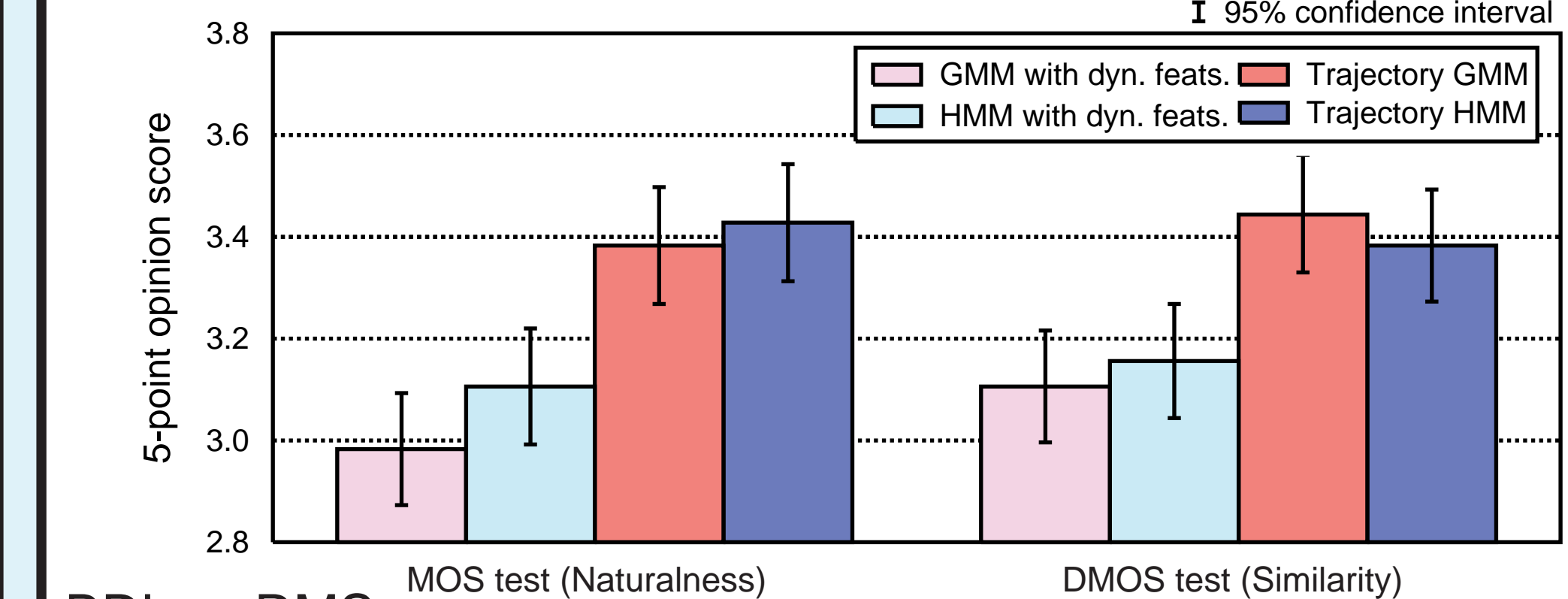


BDL (male) → RMS (male)



Subjective evaluation (MOS & DMOS)

CLB → SLT



BDL → RMS

