

Voice Conversion Using Duration-Embedded Bi-HMMs for Expressive Speech Synthesis

Source: TRANSACTIONS ON AUDIO, SPEECH, AND LANGUAGE PROCESSING

Author : Chung-Hsien Wu, Chi-Chun Hsia, Jhing-Fa Wang

Professor : 陳嘉平

Reporter : 楊治鏞

Introduction

- Voice conversion methods have previously been proposed to convert the speech signals uttered by a speaker to the other speaker with limited speech data.
- With a limited size of speech database, this study proposes a spectral and a prosodic voice conversion model as the post-process of the TTS system for expressive speech synthesis.

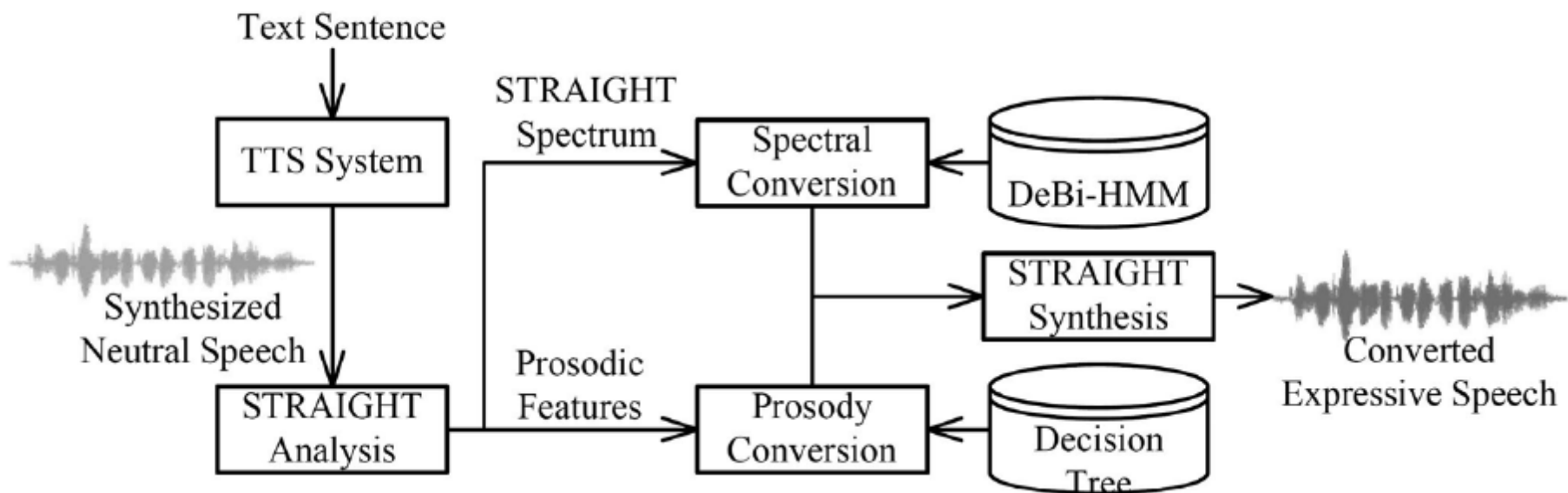


Fig. 1. Architecture of the proposed expressive speech synthesis system.

PROPOSED CONVERSION METHOD

- The Bi-HMM is adopted to model the spectral envelope evolution of source and target speech, separately but synchronously.
- The conversion function for each state is estimated with MMSE criterion under the conditional normal assumption.

Bi-HMM Spectral Conversion Model

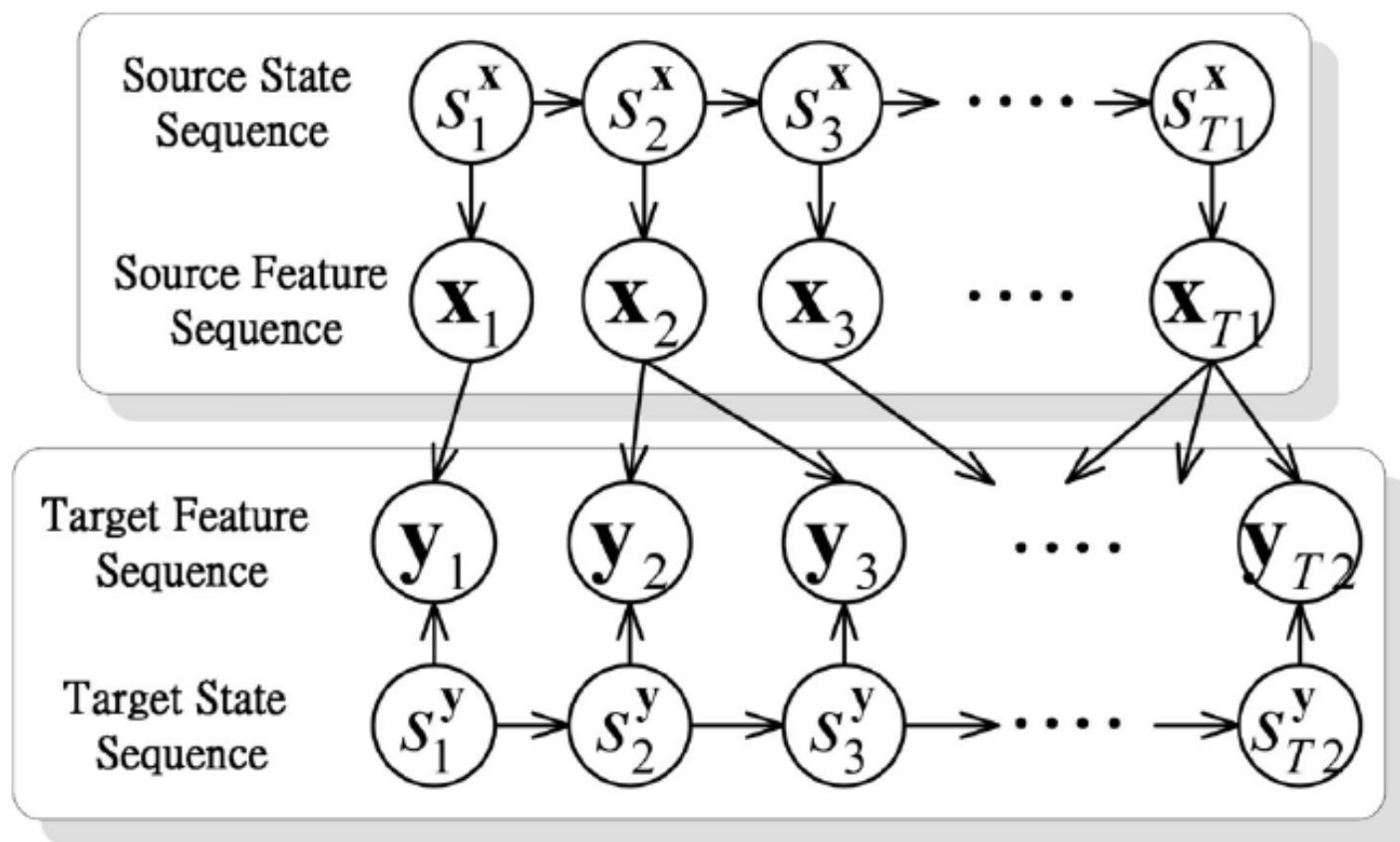


Fig. 2. State alignment between the source and target HMMs which compose the Bi-HMM.

Bi-HMM Spectral Conversion Model

- The upper portion depicts the relationship between the source feature vector sequence $X = \{x_1, x_2, \dots, x_{T_1}\}$ with length T_1 and its corresponding state sequence $S_X = \{s_{x_1}, s_{x_2}, \dots, s_{x_{T_1}}\}$.
- The lower portion depicts the relationship between the target feature vector sequence $Y = \{y_1, y_2, \dots, y_{T_2}\}$ and its corresponding state sequence $S_Y = \{s_{y_1}, s_{y_2}, \dots, s_{y_{T_2}}\}$.

Bi-HMM Spectral Conversion Model

- For a given state $s_X = i$ and an input source feature vector \mathbf{x} , the converted target vector $\tilde{\mathbf{y}}$ is predicted using the following conversion function:

$$\tilde{\mathbf{y}} = F_{s_X=i}(\mathbf{x}) = \sum_{j=1}^J \left\{ p_i(j | \mathbf{x}, \Lambda) \times \left[\mu_{i,j}^Y + \Sigma_{i,j}^{YX} \left(\Sigma_{i,j}^{XX} \right)^{-1} \left(\mathbf{x} - \mu_{i,j}^X \right) \right] \right\}$$

- where J denotes the number of mixture components;
 $\mu_{i,j}^X$ and $\mu_{i,j}^Y$ denote the mean vectors of source and target feature vector sequences in mixture j of state i .

Bi-HMM Spectral Conversion Model

- $p_i(j | \mathbf{x}, \Lambda)$ represents the conditional probability where belongs to the mixture j in state i and is estimated as

$$p_i(j | \mathbf{x}, \Lambda) = \frac{w_{i,j} N(\mathbf{x}; \mu_{i,j}^{\mathbf{x}}; \mu_{i,j}^{\mathbf{xx}})}{\sum_{k=1}^J w_{i,k} N(\mathbf{x}; \mu_{i,k}^{\mathbf{x}}; \mu_{i,k}^{\mathbf{xx}})}$$

- where $w_{i,j}$ denotes the weight of mixture j in state i .

Bi-HMM Spectral Conversion Model

- Given the Bi-HMM parameter set $\Lambda = \{\Lambda_x, \Lambda_y\}$ and the source and target feature vector sequences X and Y , the state sequences S_X and S_Y with maximum joint probability $p(Y, S_Y, X, S_X | \Lambda)$ is obtained as follows:

$$\begin{aligned} (S_Y^*, S_X^*) &= \arg \max_{S_Y, S_X} p(Y, S_Y, X, S_X | \Lambda) \\ &= \arg \max_{S_Y, S_X} p(Y, S_Y | X, S_X, \Lambda_Y) p(X, S_X | \Lambda_X) \\ &= \arg \max_{S_Y, S_X} p(Y | X, S_X, S_Y, \Lambda_Y) \\ &\quad \times p(S_Y | \Lambda_Y) p(X | S_X, \Lambda_X) p(S_X | \Lambda_X) \end{aligned}$$

Bi-HMM Spectral Conversion Model

- Given the input X , a candidate target \tilde{Y} is predicted for a specific state sequence S_X . The optimal state sequence pair $(S_X^*, S_{\tilde{Y}}^*)$ in the conversion phase is obtained by

$$\begin{aligned} (S_X^*, S_{\tilde{Y}}^*) = \arg \max_{S_Y, S_X} & p(\tilde{Y} | S_{\tilde{Y}}, \Lambda_Y) p(S_{\tilde{Y}} | \Lambda_Y) \\ & \times p(X | S_X, \Lambda_X) p(S_X | \Lambda_X) \end{aligned}$$

Bi-HMM Spectral Conversion Model

- The target vector sequence with maximum probability in Bi-HMM is obtained by the following conversion function:

$$\tilde{Y} = F_{S_X^*}(X) = \left\{ F_{S_{x_1}^*}(x_1), F_{S_{x_2}^*}(x_2), \dots, F_{S_{x_{T1}}^*}(x_{T1}) \right\}$$

Embedded Duration Model

- When the speech signal stays in state i , $i = 1, \dots, I$ with self-transition probability a_{ii} , for τ frames, the implicit duration probability density is a geometric distribution $d_i(\tau) = a_{ii}^{\tau-1} (1 - a_{ii})$.
- This exponential state duration density is inappropriate for most speech signals.

Embedded Duration Model

- In this study, a Gamma duration model is adopted as follows:

$$d_i(\tau | \eta_i, \nu_i) = \frac{\eta_i^{\nu_i}}{\Gamma(\nu_i)} \tau^{\nu_i-1} \exp(-\eta_i \tau)$$

- where $\Gamma(\cdot)$ is gamma function; $\eta_i > 0$ and $\nu_i > 0$ are parameters of the gamma distribution.

Embedded Duration Model

- The E-step calculates the expectation of the log likelihood of the complete data given new estimates, $\hat{\eta}_i$ and $\hat{\nu}_i$, after having the following current estimates, given by η_i and ν_i

$$Q(\hat{\eta}_i, \hat{\nu}_i | \eta_i, \nu_i) = \sum_{t=1}^T \left\{ \xi_{t \in t_s}(i) \cdot \left[-\log \Gamma(\hat{\nu}_i) + \hat{\nu}_i \log \hat{\eta}_i + (\hat{\nu}_i - 1) \log \tau_t - \hat{\eta}_i \tau_t \right] \right\}$$

- where t_s is the starting frame of a state, T represents the total number of frames for a training utterance

Embedded Duration Model

$$\xi_{t \in t_s}(i) = \delta(s_t - i) \delta(t - t_s) = P(s_{t \in t_s} = i | X, \Lambda)$$

where $\delta(\cdot)$ is the Kronecker delta function.

- Unfortunately, no closed-form solution to new estimate was derived.
- Since gamma distribution $d_i(\tau | \eta_i, \nu_i)$ contains the mean ν_i / η_i and variance ν_i / η_i^2 , the parameters are then derived empirically from the sample mean and variance.

Prosody Conversion by Decision Tree

- The prosody conversion decision tree (PDT) is derived from the prosodic features of a syllable.
- This work adopted the ratio of syllable duration, pitch mean, pitch dynamic range, energy mean, and energy dynamic range for prosody conversion.
- Table 1 shows the linguistic question set used to train the decision tree, which includes the linguistic features in word and sentence levels.



LINGUISTIC QUESTION SET USED TO TRAIN THE DECISION TREE FOR PROSODY CONVERSION

Word level features	Tone type of the syllable
	Word length represented in syllable number
	Part of speech of the word
	Punctuation marks ahead and behind the word
Sentence level features	Word position in the phrase (front, middle, end)
	Number of phrases in the sentence
	Word position in the sentence



EXPERIMENTS

- A total number of distinct tonal syllables of 1410 were used as the basic units to synthesize neutral speech.
- Six expressive styles were applied, namely happiness, sadness, anger, confusion, apology, and question.

STATISTICS OF EXPRESSIVE SPEECH DATABASES

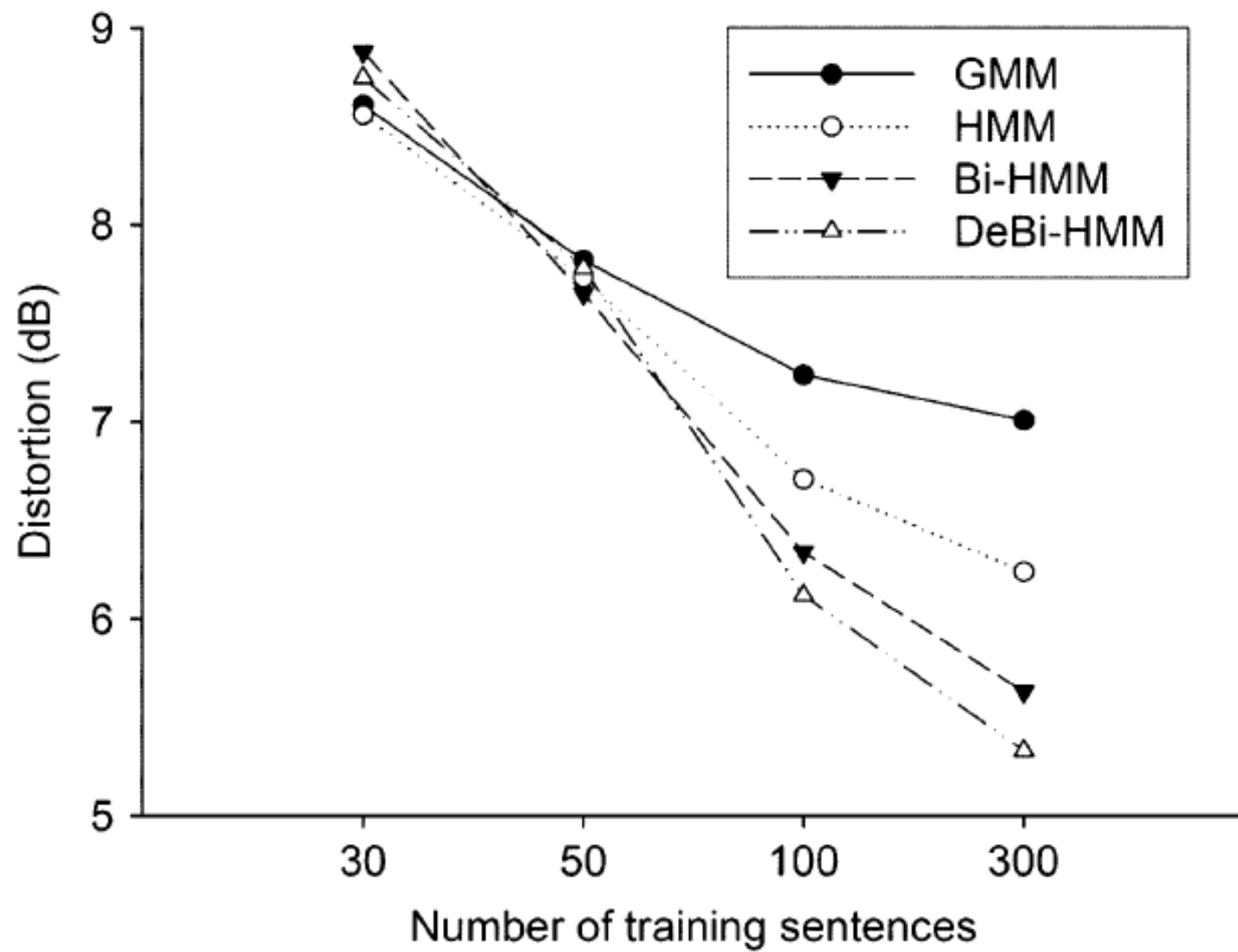
Expressive Style	Happiness				Sadness				Anger			
Number of sentences	30	50	100	300	30	50	100	300	30	50	100	300
Avg. appearance number of context indep. sub-syllables	14.24	21.81	47.00	124.62	16.24	24.10	41.78	136.33	12.28	22.52	46.28	119.53
Number of syllables	459	680	1420	3930	501	760	1310	4140	369	680	1408	3637
Number of words	270	395	855	2246	277	425	728	2262	216	386	776	1951
Length (min:sec)	1:31	2:28	4:56	15:15	1:46	3:02	5:51	17:41	1:11	1:58	4:02	11:54
Expressive Style	Confusion				Question				Apology			
Number of sentences	30	50	100	300	30	50	100	300	30	50	100	300
Avg. appearance number of context indep. sub-syllables	13.07	21.19	42.03	113.21	13.95	22.33	43.66	123.50	11.50	20.14	35.48	131.86
Number of syllables	393	630	1290	3420	426	685	1312	3750	342	602	1072	3963
Number of words	244	366	705	1839	241	374	737	1963	212	381	656	2359
Length (min:sec)	1:25	2:23	4:46	14:06	1:09	1:56	3:46	12:31	1:15	2:10	4:08	13:12

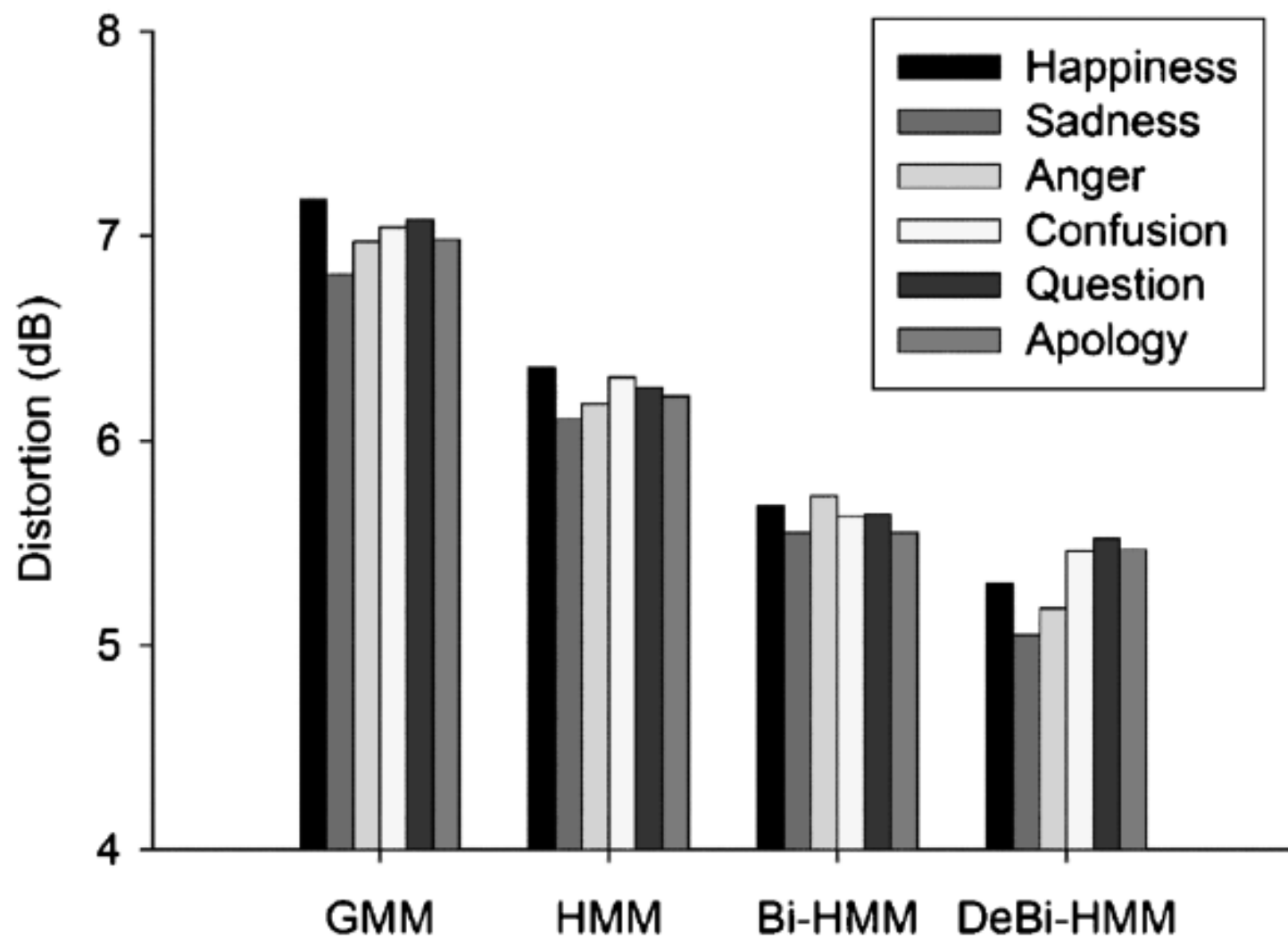
Objective test

- Distortion between the converted (from neutral) and the target (happy, sad and angry) speech is calculated using the following equation:

$$Avg_Dis = \frac{20}{\ln 10} \left(\frac{1}{T} \sum_{t=1}^T \left\| 2 (y_t - \hat{y}_t) \right\| \right)$$

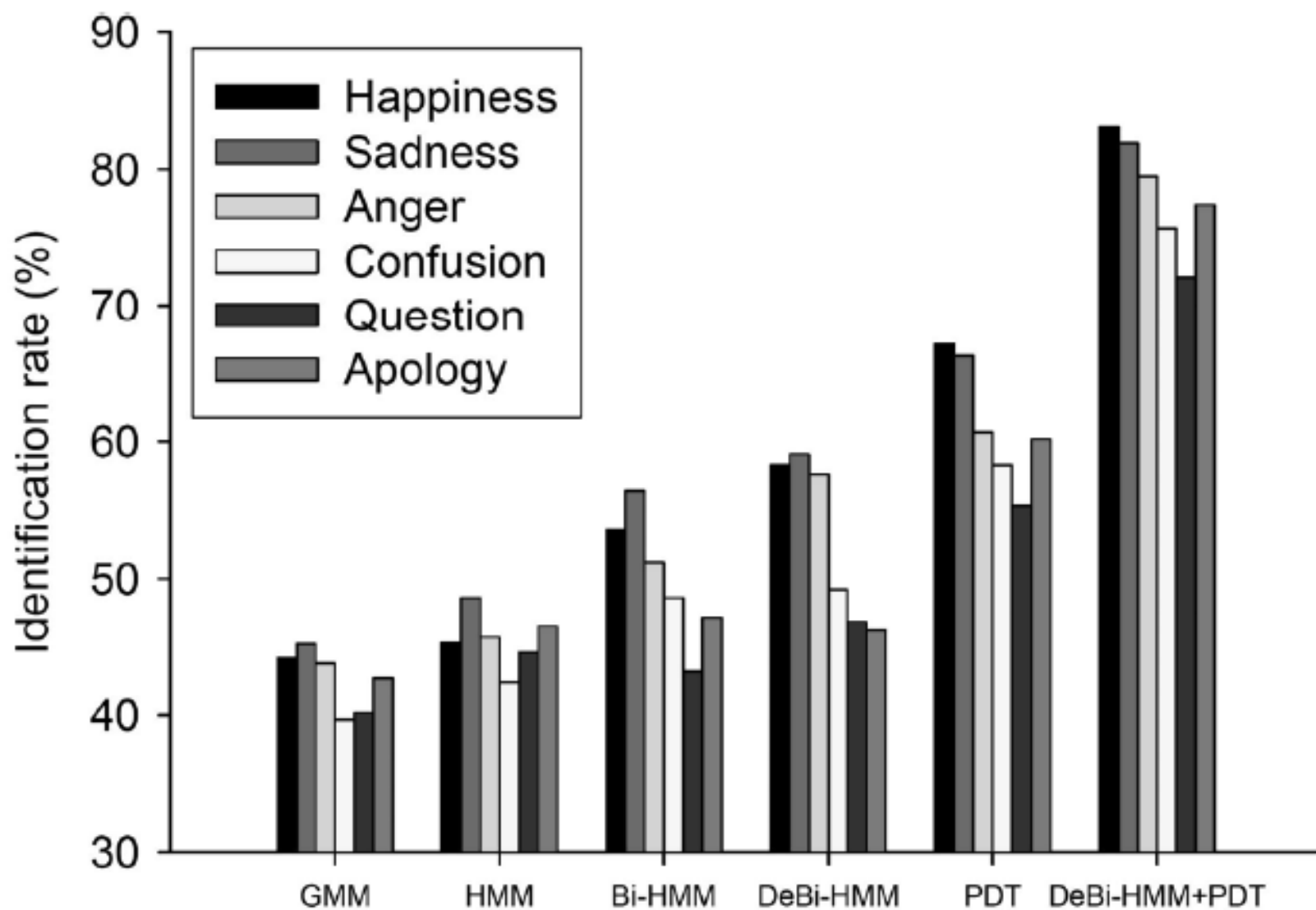
- where y_t and \hat{y}_t denote the target and the spectrum-converted feature vectors at time.





Subjective Test

- For each test sentence randomly selected from the test set, 36 converted utterances processed by each conversion method to each expressive style were randomly output to the human subjects.
- Twenty adult subjects, around 22–31 years of age, were asked to classify each utterance as one of the six expressive styles.





Subjective Test

- The naturalness of the converted utterance was also evaluated, according to a 5-scale scoring method excellent very poor.

