

S1. Supplementary Material

S1.1. Details of the FHVAE model

We follow the terminology used in [22] to describe the details of the FHVAE model. Let $\mathcal{D} = \{\mathbf{X}^i\}_{i=1}^M$ denote a speech dataset with M sequences. The i -th sequence \mathbf{X}^i contains N^i speech segments $\{\mathbf{x}^{(i,n)}\}_{n=1}^{N^i}$, where $\mathbf{x}^{(i,n)}$ is a segment of a fixed number of frames. The FHVAE model formulates the generation process of a sequence \mathbf{X} as¹¹ [22],

1. A vector $\boldsymbol{\mu}_2$ is drawn from a prior distribution $p_\theta(\boldsymbol{\mu}_2) = \mathcal{N}(\mathbf{0}, \sigma_{\boldsymbol{\mu}_2}^2 \mathbf{I})$;
2. Latent segment variables \mathbf{z}_1^n and latent sequence variables \mathbf{z}_2^n are drawn from $p_\theta(\mathbf{z}_1^n) = \mathcal{N}(\mathbf{0}, \sigma_{\mathbf{z}_1}^2 \mathbf{I})$ and $p_\theta(\mathbf{z}_2^n | \boldsymbol{\mu}_2) = \mathcal{N}(\boldsymbol{\mu}_2, \sigma_{\mathbf{z}_2}^2 \mathbf{I})$;
3. Speech segment \mathbf{x}^n is drawn from

$$p_\theta(\mathbf{x}^n | \mathbf{z}_1^n, \mathbf{z}_2^n) = \mathcal{N}(f_{\mu_x}(\mathbf{z}_1^n, \mathbf{z}_2^n), \text{diag}(f_{\sigma_x^2}(\mathbf{z}_1^n, \mathbf{z}_2^n))). \quad (\text{S.1})$$

Here \mathcal{N} denotes the standard normal distribution, $f_{\mu_x}(\cdot, \cdot)$ and $f_{\sigma_x^2}(\cdot, \cdot)$ are parameterized by two DNNs. Based on Equation (S.1), the joint probability for generating \mathbf{X} is formulated as (same as Equation (2)),

$$p_\theta(\boldsymbol{\mu}_2) \prod_{n=1}^N p_\theta(\mathbf{z}_1^n) p_\theta(\mathbf{z}_2^n | \boldsymbol{\mu}_2) p_\theta(\mathbf{x}^n | \mathbf{z}_1^n, \mathbf{z}_2^n). \quad (\text{S.2})$$

The FHVAE introduces an inference model to approximate the true posterior as follows (same as Equation (3)),

$$p_\phi(\boldsymbol{\mu}_2) \prod_{n=1}^N p_\phi(\mathbf{z}_2^n | \mathbf{x}^n) p_\phi(\mathbf{z}_1^n | \mathbf{x}^n, \mathbf{z}_2^n). \quad (\text{S.3})$$

Here $p_\phi(\boldsymbol{\mu}_2)$, $p_\phi(\mathbf{z}_2^n | \mathbf{x}^n)$ and $p_\phi(\mathbf{z}_1^n | \mathbf{x}^n, \mathbf{z}_2^n)$ are all diagonal Gaussian distributions. The mean and variance values of $p_\phi(\mathbf{z}_2^n | \mathbf{x}^n)$ and $p_\phi(\mathbf{z}_1^n | \mathbf{x}^n, \mathbf{z}_2^n)$ are parameterized by DNNs.

¹¹For simplicity, the superscript i in \mathbf{X}^i and subsequent equations is omitted. This does not cause confusion.

The FHVAE is trained to optimise the *discriminative segmental variational lower bound* $\mathcal{L}(\theta, \phi; \mathbf{x}^{(i,n)})$ [22], which is defined as,

$$\begin{aligned}
& \mathbb{E}_{q_\phi(\mathbf{z}_1^{(i,n)}, \mathbf{z}_2^{(i,n)} | \mathbf{x}^{(i,n)})} [\log p_\theta(\mathbf{x}^{(i,n)} | \mathbf{z}_1^{(i,n)}, \mathbf{z}_2^{(i,n)})] \\
& - \mathbb{E}_{q_\phi(\mathbf{z}_2^{(i,n)} | \mathbf{x}^{(i,n)})} [\text{KL}(q_\phi(\mathbf{z}_1^{(i,n)} | \mathbf{x}^{(i,n)}, \mathbf{z}_2^{(i,n)}) || p_\theta(\mathbf{z}_1^{(i,n)}))] \\
& - \text{KL}(q_\phi(\mathbf{z}_2^{(i,n)} | \mathbf{x}^{(i,n)}) || p_\theta(\mathbf{z}_2^{(i,n)} | \tilde{\boldsymbol{\mu}}_2^i)) \\
& + \frac{1}{N^i} \log p_\theta(\tilde{\boldsymbol{\mu}}_2^i) + \alpha \log p(i | \mathbf{z}_2^{(i,n)}),
\end{aligned} \tag{S.4}$$

where $\tilde{\boldsymbol{\mu}}_2^i$ denotes the posterior mean of $\boldsymbol{\mu}_2$ for the i -th sequence and α denotes the discriminative weight. The discriminative objective $\log p(i | \mathbf{z}_2^{(i,n)})$ is formulated as,

$$\log p(i | \mathbf{z}_2^{(i,n)}) := \log p_\theta(\mathbf{z}_2^{(i,n)} | \tilde{\boldsymbol{\mu}}_2^i) - \log \sum_{j=1}^M p_\theta(\mathbf{z}_2^{(j,n)} | \tilde{\boldsymbol{\mu}}_2^j). \tag{S.5}$$

After FHVAE training, \mathbf{z}_1 representation is extracted as the desired speaker-invariant representation of speech.