

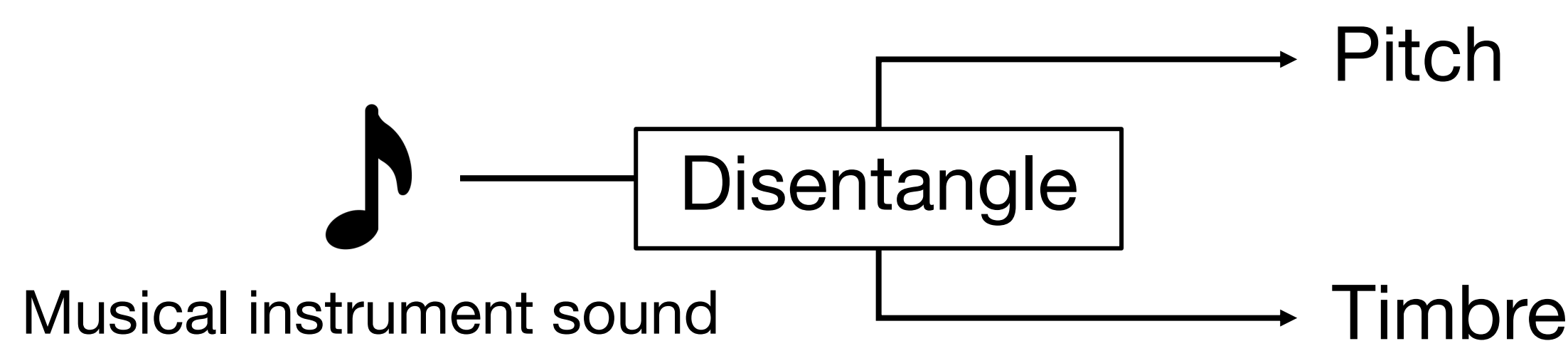
Pitch-Timbre Disentanglement of Musical Instrument Sounds Based on VAE-Based Metric Learning

Keitaro Tanaka¹ Ryo Nishikimi² Yoshiaki Bando³ Kazuyoshi Yoshii² Shigeo Morishima⁴

¹ Waseda University ² Kyoto University ³ National Institute of Advanced Industrial Science and Technology (AIST)
⁴ Waseda Research Institute for Science and Engineering, Japan

Backgrounds

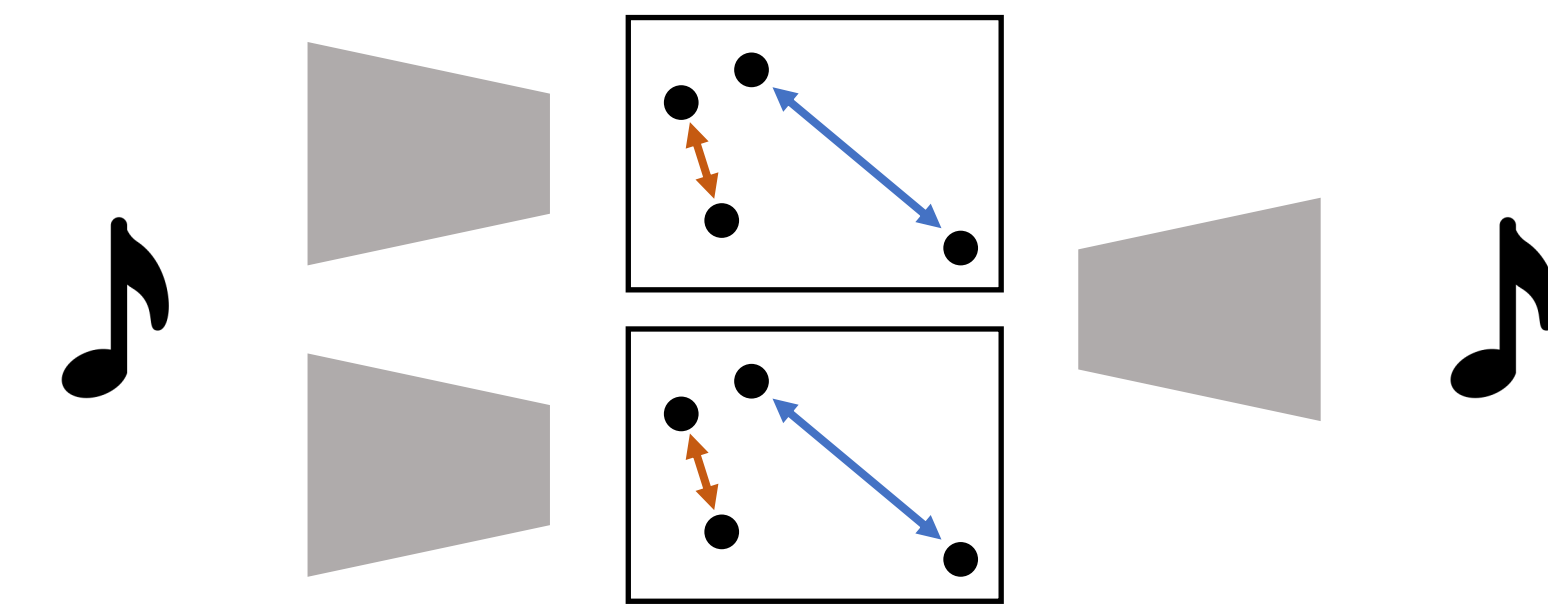
What is **Disentanglement**?



- To describe data as a combination of independent factors
- A **sound** can be disentangled into **pitch and timbre**
- Conventional approach cannot treat an arbitrary musical instrument sound **because of using the concrete category labels**

Approach

Introduce a **metric learning technique**

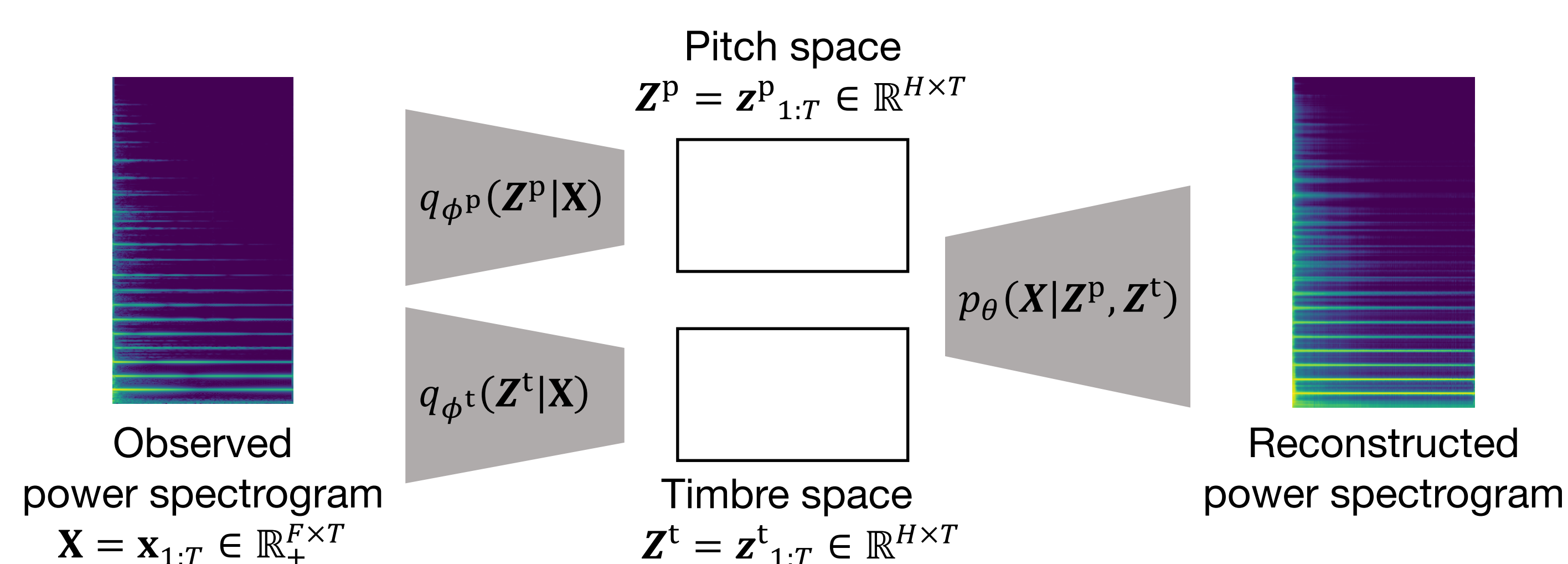


- A technique used for **representing the dissimilarities** of samples as the distances in a latent space
- Similar** samples are mapped **close** to each other
- Dissimilar** samples are mapped **far away** from each other

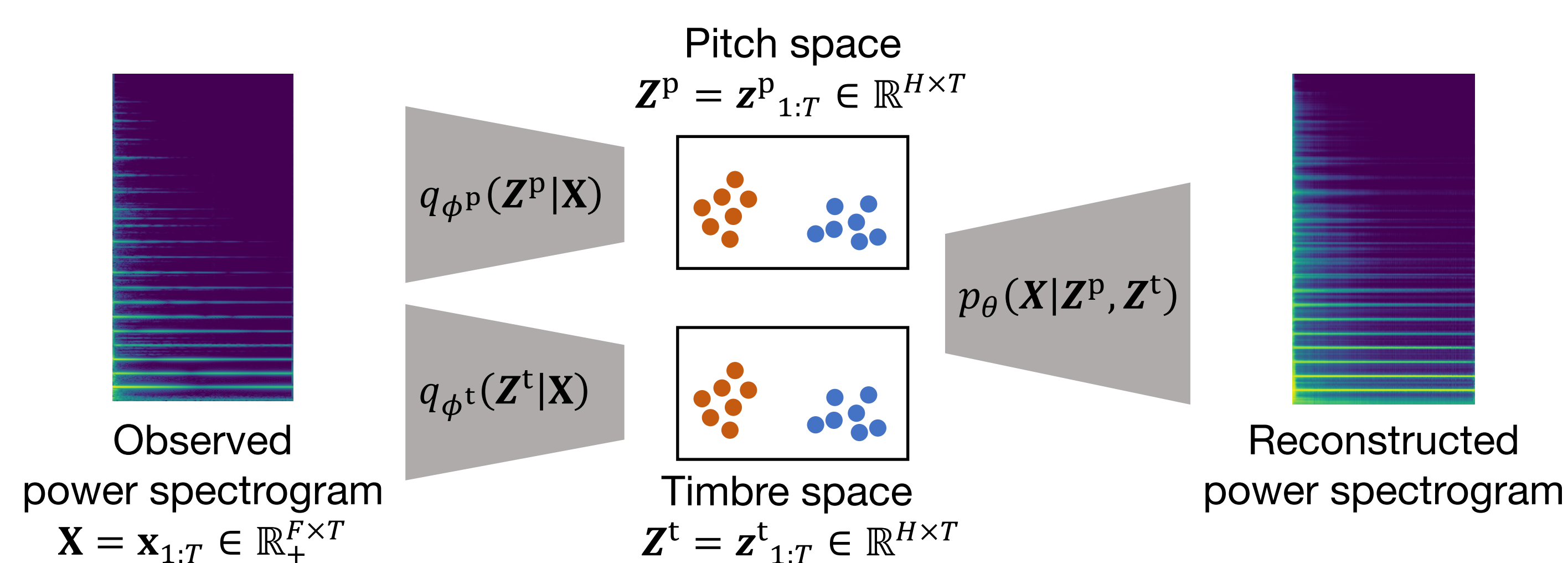
Method

Generative Model of the Observed Spectrogram

Formulate a probabilistic model of the observed spectrogram \mathbf{X} with latent representations \mathbf{Z}^p and \mathbf{Z}^t



Transform two observed spectrograms \mathbf{X}_1 and \mathbf{X}_2 into the latent variables \mathbf{Z}^p and \mathbf{Z}^t independently



- Conduct pairwise metric learning for these $2T$ samples

Pairwise Metric Learning for Disentanglement

Contrastive loss functions \mathcal{L}_c^p (for pitch) and \mathcal{L}_c^t (for timbre)

$$\mathcal{L}_c^p = \mathcal{D}_{11}^p + \mathcal{D}_{22}^p \pm \mathcal{D}_{12}^p \quad (\mathcal{L}_c^t \text{ is calculated in a similar way})$$

The sum of distances between all latent variable pairs of the same spectrogram (\mathbf{X}_1 and \mathbf{X}_2)

The sum of distances between all latent variable pairs from different spectrograms
+ : if \mathbf{X}_1 and \mathbf{X}_2 have the same pitch
- : otherwise

- \mathcal{L}_c^p and \mathcal{L}_c^t pull similar samples close to each other and keep dissimilar samples far from each other

Training with the Weakly Supervised Learning

The training is conducted with the total loss function $\mathcal{L}^{\text{total}}$

$$\mathcal{L}^{\text{total}} = -\mathcal{L}^{\text{vae}} + \alpha \mathcal{L}_c^p + \beta \mathcal{L}_c^t$$

$$\mathbb{E}_{q_{\phi^p} q_{\phi^t}(\mathbf{Z}^p, \mathbf{Z}^t | \mathbf{X})} [\log p_{\theta}(\mathbf{X} | \mathbf{Z}^p, \mathbf{Z}^t)] - \mathcal{D}_{\text{KL}}(q_{\phi^p}(\mathbf{Z}^p | \mathbf{X}) \| p(\mathbf{Z}^p)) - \mathcal{D}_{\text{KL}}(q_{\phi^t}(\mathbf{Z}^t | \mathbf{X}) \| p(\mathbf{Z}^t))$$

Hyperparameters to control the weights

- Only information on **whether pitches and timbres of a pair of observed spectrograms are identical or not is required**
- Their actual labels are not necessary
- An arbitrary musical instrument sound can be treated**

Evaluation

Evaluate **denseness** and **divergence** for unseen musical instruments

- Denseness** shows how **close** the latent variables with **the same pitch or timbre** label are
- Calculated as:

$$\frac{1}{M} \sum_{m=1}^M \frac{1}{9N_m} \sum_{n=1}^{N_m} \sum_{t=1}^9 \| \mathbf{z}_{mnt}^p - \boldsymbol{\eta}_m^p \| \quad (\text{Timbre in a similar way})$$

$$\boldsymbol{\eta}_m^p = \frac{1}{9N_m} \sum_{n=1}^{N_m} \sum_{t=1}^9 \mathbf{z}_{mnt}^p$$

- Divergence** shows how **far** the latent variables with **different pitch or timbre** labels are
- Calculated as:

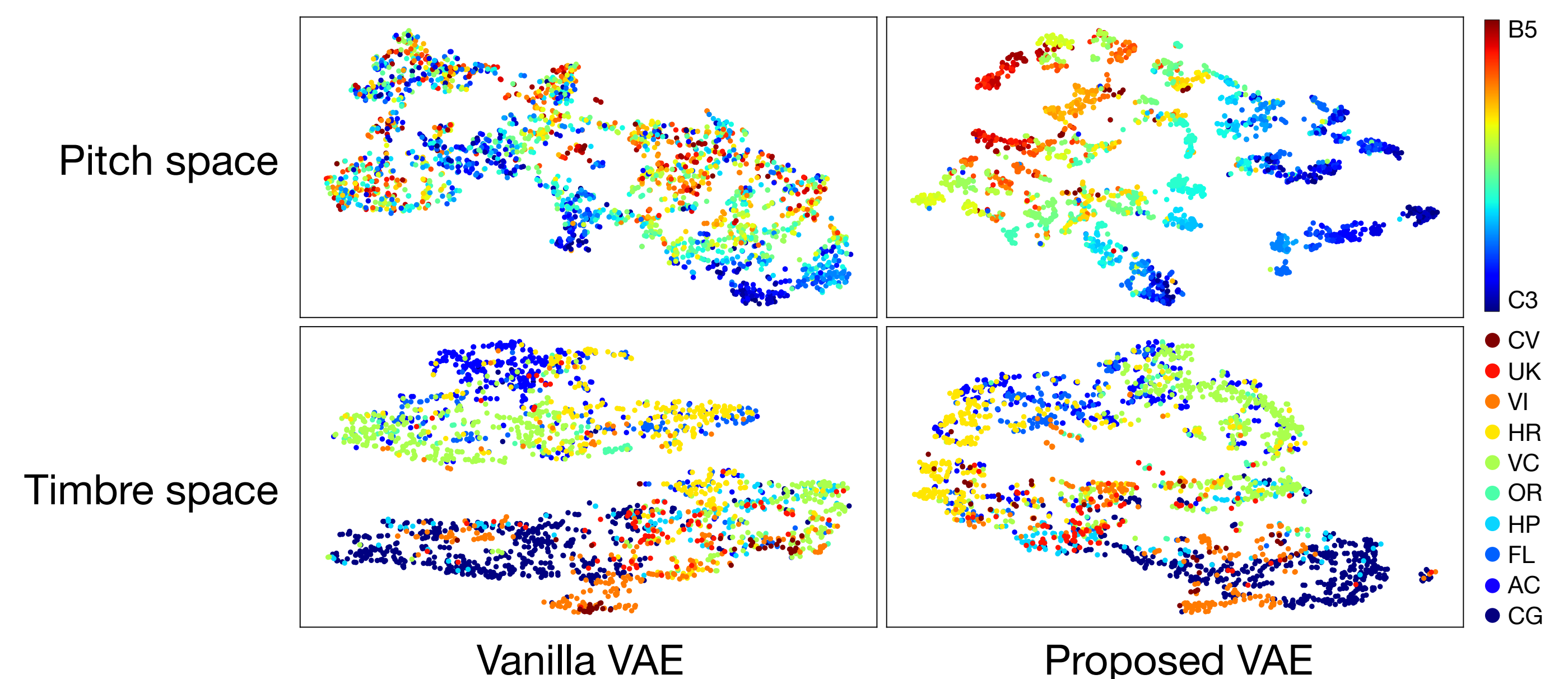
$$\frac{2}{M(M-1)} \sum_{m_1=1}^{M-1} \sum_{m_2=1}^M \| \boldsymbol{\eta}_{m_1}^p - \boldsymbol{\eta}_{m_2}^p \| \quad (\text{Timbre in a similar way})$$

Result

Methods	Pitch representations		Timbre representations	
	Denseness ↓	Divergence ↑	Denseness ↓	Divergence ↑
Vanilla VAE	3.334	2.279	3.640	1.541
Proposed VAE	2.891	3.551	3.420	2.654

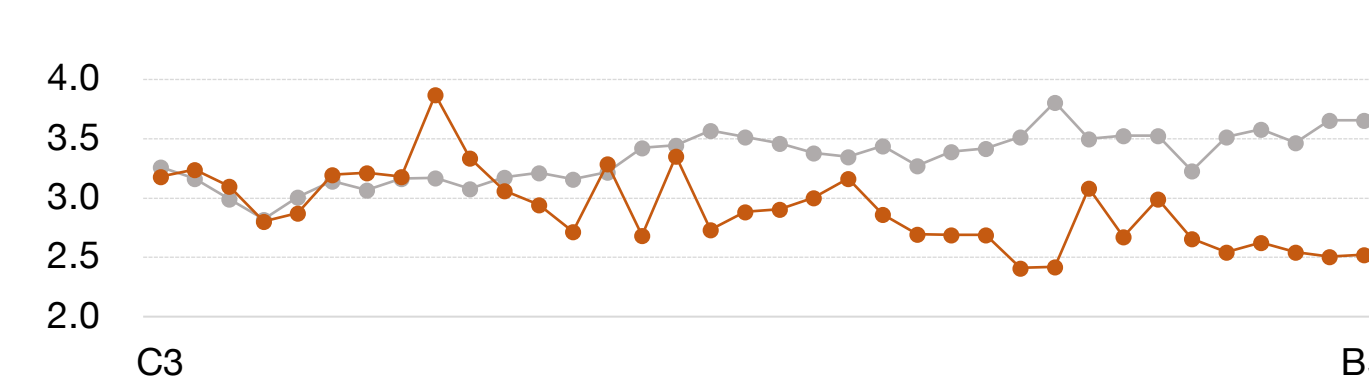
- The denseness got smaller**, and **the divergence got larger** in both latent spaces by introducing the metric learning

Visualizations of the pitch and timbre spaces



Denseness and divergence in details

- Denseness**



- Divergence**

