

Sound & Music Computing Lab

# Domain Adversarial Training on Conditional VAE for Controllable Music Generation



MUSIC X LAB

Jingwei Zhao Gus Xia Ye Wang

## MOTIVATION

- The conditional VAE architecture can suffer from "condition collapse", because  $z_x$  is often too informative so the the decoder tends to ignore condition c.
- Domain adversarial training (DAT) can induce representation  $z_x$  to be disentangled from *c*, but it only applies to simple scenarios with categorical condition.
- We propose an adversarial condition de-

# MODEL ARCHITECTURE



noising objective and generalize DAT to controllable music generation with complex sequential condition (e.g., melody).



# **CHORD REPRESENTATION**



# **OBJECTIVES & TRAINING**

• VAE Objective

# **DEMO: CHORD GENERATION WITH VARIED MELODY CONDITIONS**







$$\mathcal{L}(\theta_{\text{enc}}, \theta_{\text{dec}}) = -\mathbb{E}_Q \left[ \log P_{\theta_{\text{dec}}} \left( x \mid z_x, c \right) \right] + \alpha \mathbb{KL}(Q_{\theta_{\text{enc}}}(z_x \mid x, c) \parallel \mathcal{N}(\mathbf{0}, \mathbf{1})),$$
(1)

Adversarial Objective

 $\mathcal{L}(\theta_{\rm dis}) = -\mathbb{E}_Q \left[ \log R_{\theta_{\rm dis}} \left( c \mid z_x, c^* \right) \right], \qquad (2)$  $\mathcal{L}(\theta_{\rm enc} \mid \theta_{\rm dis}) = -\mathbb{E}_Q \left[ \log R_{\theta_{\rm dis}} \left( \mathbf{1} - c \mid z_x, c^* \right) \right] \\ + \alpha \mathbb{KL}(Q_{\theta_{\rm enc}}(z_x \mid x, c) \parallel \mathcal{N}(\mathbf{0}, \mathbf{1})), \qquad (3)$ 

• Training Procedure

#### 1 while training do

- 3 Optimize VAE with  $\mathcal{L}(\theta_{enc}, \theta_{dec})$ ,
- 4 **for** *j iterations* **do**
- 5 **for** *k iteration* **do** 
  - Optimize discriminator with  $\mathcal{L}(\theta_{dis})$ ,
- **for** *l iterations* **do**
- Optimize encoder with  $\mathcal{L}(\theta_{enc} \mid \theta_{dis})$ .
- Training Loss Curve



## EXPERIMENTS

#### • Ablation Models

*Non-DAT*: Same VAE framework but without a discriminator. It does not explicitly try to disentangle  $z_x$  from *c* using domain adversarial training (DAT);

*Mask-CR*: Applying a general masking corruption instead of pitch transposition for condition corruption;

*Non-CR*: Using the conventional DAT objective without condition corruption. It predicts c directly from  $z_x$  using a GRU-based discriminator.

#### • Objective Evaluation



#### Disentanglement



#### • Subjective Evaluation



**Figure 7**: Subjective evaluation on the harmonization performance of our model and baseline models.

**Figure 8**: Object evaluation on representation similarity (invariance) against pitch transposition. A higher value denotes better disentanglement.



**Figure 9**: Objective evaluation on harmony histogram upon melody swapping. A higher ratio in root, 3rd, and 5th notes indicates a higher degree of controllability.