

Lisa Kawai<sup>1</sup>, Philippe Esling<sup>2</sup>, Tatsuya Harada<sup>1,3</sup>  
<sup>1</sup> The University of Tokyo, <sup>2</sup> IRCAM, <sup>3</sup> RIKEN

## Research Objective

**Continuous control** of musical attributes is important.

→ To transform music with a model that

1. can control **how much** it changes the musical attributes



2. is applicable to **any set of musical annotations**

→ No need to change the model structure or implementation for new attributes

## Musical Attributes

Annotations from **jSymbolic** [McKay & Fujinaga, 2006]

- Calculates **musical statistics** (e.g. pitch, melody, chord, and rhythm)
- Selected 12 interpretable features, which are **continuous** values

## Experimental Settings

Dataset: **Nottingham** dataset [Foxley, 2011] (monophonic)

Input representation: piano-roll with {rest, continue} tokens

Baselines: Naïve VAE (MusicVAE [Simon+, 2018]) and

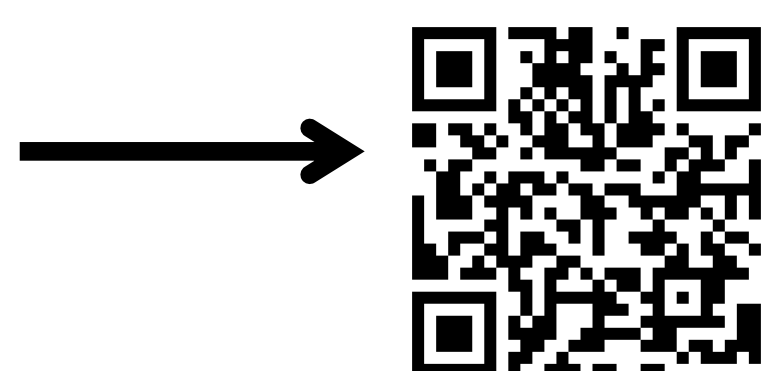
GLSR-VAE [Hadjeres+, 2017]

Network: Composed of Gated Recurrent Unit

## Project Page and Contact

Other results with audio samples are available at

Email: [kawai@mi.t.u-tokyo.ac.jp](mailto:kawai@mi.t.u-tokyo.ac.jp)



## Method

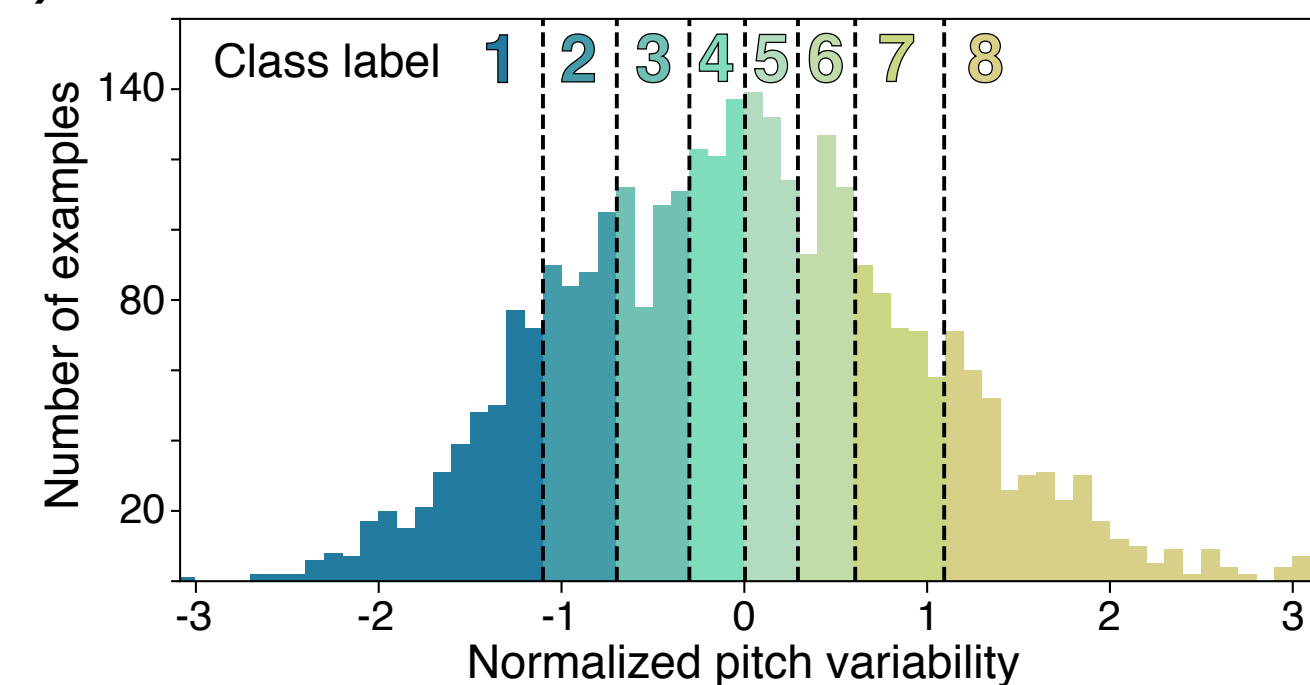
- VAE-based model with attribute conditioning to decoder
- Adversarial discriminator to force  $z$  to abstract from  $a$
- Usually  $a$  is **binary** [Lample+, 2017] ( $a, 1 - a \in \{0, 1\}$ )

↔ We have **continuous annotations** of attributes which are needed to achieve continuous control.

How did we solve this problem?

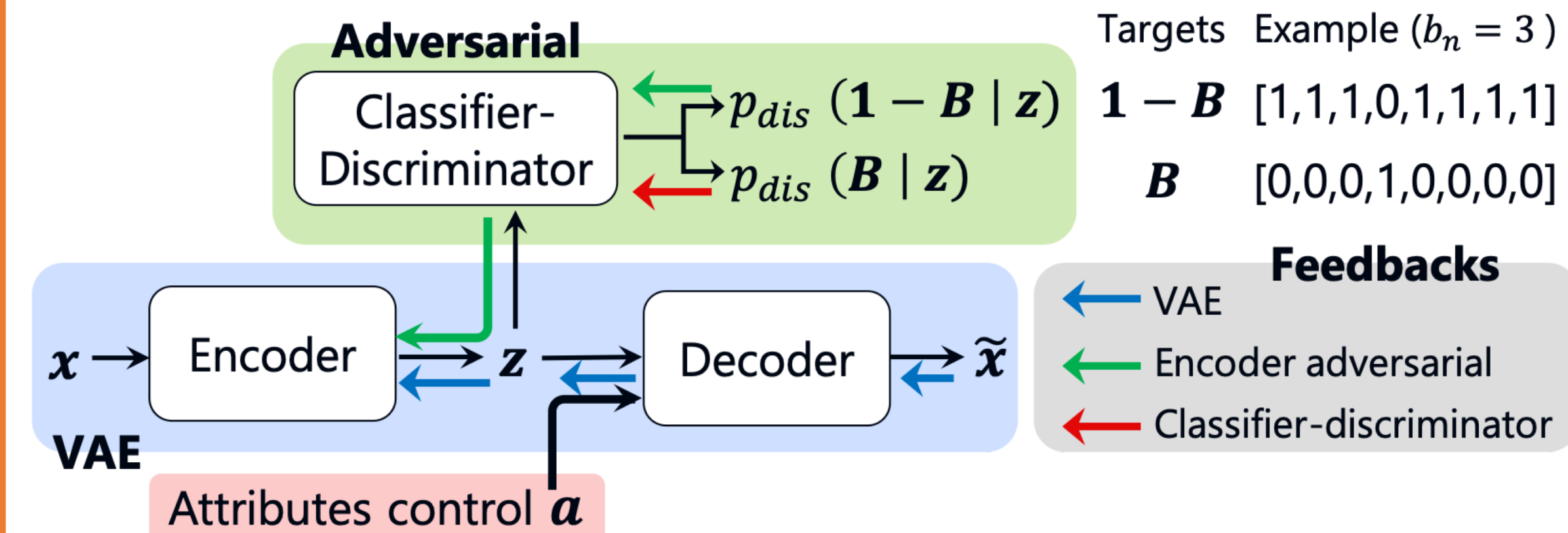
1. **Quantize** data into  $K (= 8)$  classes

e.g.  $a$   $b$   
 -1.0  $\rightarrow$  2  
 -0.5  $\rightarrow$  3  
 1.0  $\rightarrow$  7

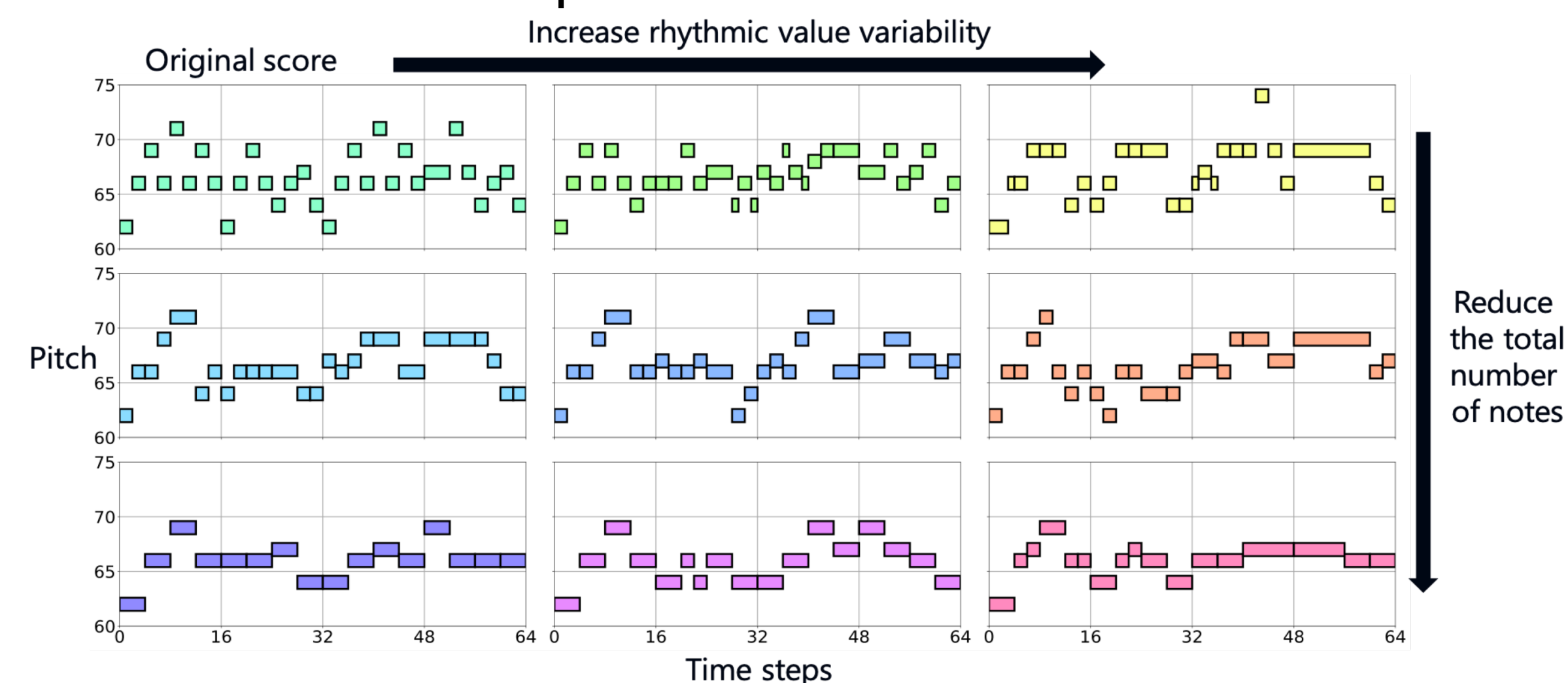


2. Extend the discriminator

to be **multivariate**

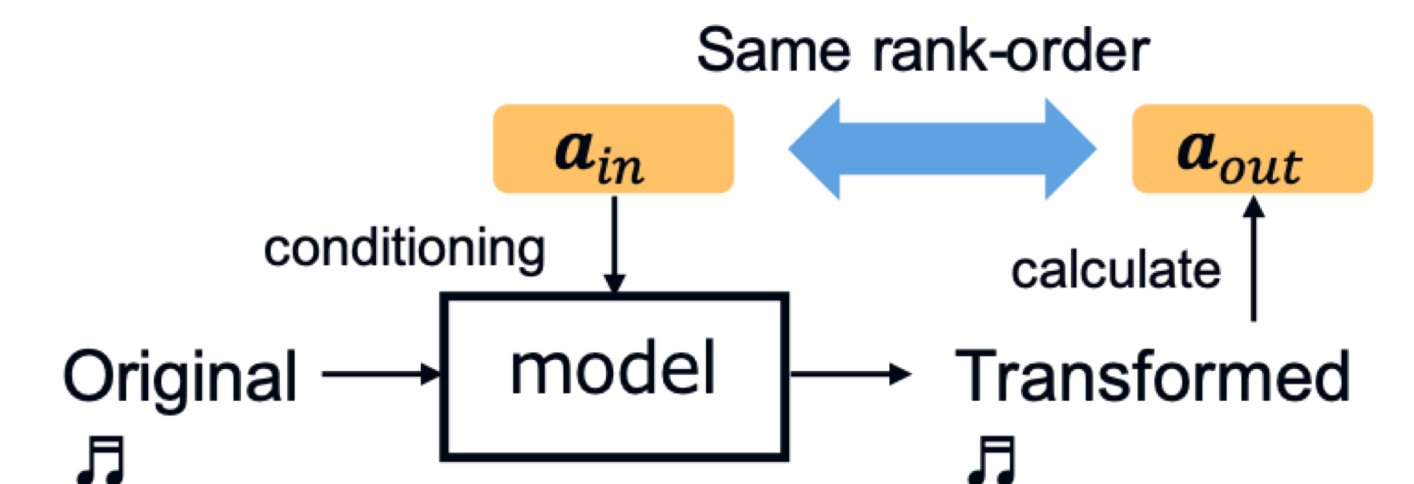


## Generated Samples



## Results

Spearman's rank-order **correlation coefficient**



attribute	Naive	GLSR	Ours
total number of notes	0.973	0.975	<b>0.981</b>
pitch variability	0.807	-	<b>0.938</b>
rhythmic value variability	0.830	-	<b>0.938</b>
pitch kurtosis	0.528	-	<b>0.723</b>
pitch skewness	0.366	-	<b>0.492</b>
prevalence of most common rhythmic value	0.795	-	<b>0.851</b>
average note duration	0.968	-	<b>0.983</b>
note density variability	0.677	-	<b>0.855</b>
amount of arpeggiation	0.126	-	<b>0.386</b>
chromatic motion	0.284	-	<b>0.622</b>
direction of melodic motion	0.428	-	<b>0.702</b>
average interval spanned by melodic arcs	0.262	-	<b>0.523</b>
chord duration	0.920	-	<b>0.944</b>

Our model outperforms the others on all the attributes.

## Reconstruction accuracy

Our proposed model does not deteriorate the reconstruction quality.

	Naive	GLSR	Ours
NLL	2.269	<b>1.002</b>	1.679
accuracy	0.759	<b>0.808</b>	0.790

## Chord consistency: Cosine similarity of chroma feature

12-dim vector that shows duration ratio of pitch classes

Our method outperforms the baselines.

GLSR-VAE is better for  $\delta \approx 0$ , but the cosine similarity significantly drops for  $\delta > 0$ .

