

LEARNING INTERPRETABLE REPRESENTATION FOR CONTROLLABLE POLYPHONIC MUSIC GENERATION

Ziyu Wang

Dingsu Wang

Yixiao Zhang

Gus Xia

Music X Lab, Computer Science Department, NYU Shanghai

{ziyu.wang, dingsu.wang, yixiao.zhang, gxia}@nyu.edu

ABSTRACT


While deep generative models have become the leading methods for algorithmic composition, it remains a challenging problem to *control* the generation process because the latent variables of most deep-learning models lack good interpretability. Inspired by the content-style disentanglement idea, we design a novel architecture, under the VAE framework, that effectively learns two interpretable latent factors of polyphonic music: chord and texture. The current model focuses on learning 8-beat long piano composition segments. We show that such chord-texture disentanglement provides a controllable generation pathway leading to a wide spectrum of applications, including compositional style transfer, texture variation, and accompaniment arrangement. Both objective and subjective evaluations show that our method achieves a successful disentanglement and high quality controlled music generation.¹

1. INTRODUCTION

With the development of artificial neural networks, deep learning has become one of the most popular techniques for automated music generation. In particular, we see recurrent and attention-based models being able to generate creative and human-like music without heavily hand-crafted rules [1–3]. However, the main drawback of these deep generative models is that they behave like “black boxes”, and it is difficult to interpret the musical meaning of their internal latent variables [4]. Consequently, it remains a challenging task to control the generation process (i.e., to guide the music flow by manipulating the high-level compositional factors such as melody contour, accompaniment texture, style, etc.). This limitation restricts the application scenario of the powerful deep generative models.

In this paper, we improve the model interpretability for music generation via constrained representation learning. Inspired by the content-style disentanglement idea [5],

¹ Code and demos can be accessed via <https://github.com/ZZWang/polyphonic-chord-texture-disentanglement>

 © Z. Wang, D. Wang, Y. Zhang, G. Xia. Licensed under a Creative Commons Attribution 4.0 International License (CC BY 4.0). **Attribution:** Z. Wang, D. Wang, Y. Zhang, G. Xia, “Learning interpretable representation for controllable polyphonic music generation”, in *Proc. of the 21st Int. Society for Music Information Retrieval Conf.*, Montréal, Canada, 2020.

we enforce the model to learn two fundamental factors of polyphonic music: *chord* (content) and *texture* (style). The former refers to the representation of the underlying chord progression, and the latter includes chord arrangement, rhythmic pattern, and melody contour. The current design focuses on learning 8-beat long piano composition segments under a variational autoencoder (VAE) framework.

The core of the model design lies in the encoder. We incorporate the encoder with two **inductive biases** for a successful **chord-texture disentanglement**. The former applies a rule-based chord recognizer and embeds the information into the first half of the latent representation. The latter regards music as 2-D images and uses a chord-invariant convolutional network to extract the texture information, storing it into the second half of the latent representation. As for the decoder, we adopt the design from PianoTree VAE [6], an architecture that can reconstruct polyphonic music from the latent representation in a hierarchical manner.

We further show that the interpretable representations are **general-purpose**, empowering a wide spectrum of controllable music generation. In this study, we explore the following three scenarios:

Task 1: Compositional style transfer by swapping the chord and texture factors of different pieces of music, which can help us re-harmonize or re-arrange a music piece following the style of another piece.

Task 2: Texture variation by sampling the texture factor while keeping the chord factor, which is analogous to the creation of “Theme and Variations” form of composition.

Task 3: Accompaniment arrangement by predicting the texture factor given the melody using a downstream encoder-decoder generative model.

In sum, the contributions of our paper are as follows:

- We design a representation disentanglement method for polyphonic music, which learns two interpretable factors: chord and texture.
- We show that the interpretable factors are general-purpose features for controllable music generation, which reduces the necessity to design heavily-engineered control-specific model architectures. As far as we know, this is the first attempt to explicitly con-

trol the compositional texture feature for symbolic polyphonic music generation.

- We demonstrate that control methods are effective and the quality of generated music is high. Some style transferred pieces are rated even higher than the original ones composed by humans.

2. RELATED WORK

We review two techniques of automated music generation related to our paper: controlled generation (in Section 2.1) and representation disentanglement (in Section 2.2). For a more general review of deep music generation, we refer readers to [7, 8].

2.1 Controlled Music Generation

Most existing learning-based methods regard controlled music generation a *conditional estimation* problem. That is, to model $p(\text{music}|\text{control})$, in which both music and control are usually time-series features. Another approach that is closely related to conditional estimation is to first learn the joint distribution $p(\text{music}, \text{control})$ and later on *force* the value of control during the generation process.

The above two methods have been used in various tasks, including generating chords based on the melody [9], creating the melody based on the chords [10, 11], completing the counterparts or accompaniment based on the melody or chord [3, 12–16], and producing the audio waveform based on timbre features [17, 18].

However, many abstract music factors, such as texture and melody contour, could hardly be explicitly coded by labels. Even if such labels are provided, the control still does not allow continuous manipulation, such as sampling and interpolation. Consequently, it remains a challenging task to control music by more abstract factors without complex heuristics [19].

2.2 Music Representation Disentanglement

Learning disentangled representations is an ideal solution to the problem above, since: 1) representation learning embeds discrete music and control sequences into a continuous latent space, and 2) disentanglement techniques can further decompose the latent space into interpretable subparts that correspond to abstract music factors. Recent studies show that VAEs [20, 21] are in general an effective framework to learn the representations of discrete music sequences, and the key to a successful disentanglement is to incorporate proper inductive biases into the representation learning models [22].

Under a VAE framework, an inductive bias can be realized in various forms, including constraining the encoder [23–25], constraining the decoder [26], imposing multitask loss functions [27, 28], and enforcing transformation invariant results during the learning process [29, 30]. This study is based on our previous work Deep Music Analogy [27] in which we disentangle pitch and rhythm factors for monophonic segments. We extend this idea to polyphonic composition while the model design is more similar to [24].

3. MODEL

In this section, we introduce the model design and data representation in detail. The goal is to learn the representations of 8-beat long piano compositions (with $\frac{1}{4}$ beat as the shortest unit) and disentangle the representations into two interpretable factors: chord and texture.

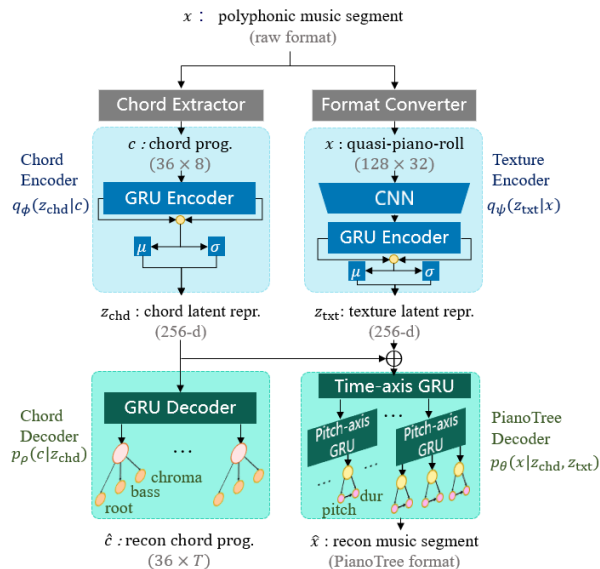


Figure 1: The model diagram.

Figure 1 shows the overall architecture of the model. It adopts a VAE framework and contains four parts: 1) a chord encoder, 2) a chord decoder, 3) a texture encoder, and 4) a PianoTree decoder. The chord encoder and chord decoder can be seen as a standalone VAE which extracts the latent chord representation z_{chd} . On the other hand, the texture encoder aims to extract the texture representation z_{txt} using a chord-invariant convolutional mapping. Finally, the PianoTree decoder takes in both z_{chd} and z_{txt} and outputs the original music in a tree-structured data format.

3.1 Chord Encoder

The chord encoder first applies rule-based methods [31, 32] to extract the chord progression under one-beat resolution. Each extracted chord progression is a 36 by 8 matrix, where each column denotes a chord of one beat. Each chord is a 36-D vector consisting of three parts: a 12-D one-hot vector for the pitch class of the *root*, a 12-D one-hot vector for the *bass*, and a 12-D multi-hot *chroma* vector.

The chord progression is then fed into a bi-directional GRU encoder [21], and the last hidden states on both ends of the GRU are concatenated and used to approximate the posterior distribution of z_{chd} . Following the assumption of a standard VAE, z_{chd} has a standard Gaussian prior and follows an isotropic Gaussian posterior.

Note that although the chord progression here is extracted using algorithms, it can also be provided by external labels, in which case the whole model becomes a conditional VAE [33].

3.2 Chord Decoder

The chord decoder reconstructs the chord progression from z_{chd} using another bi-directional GRU. The reconstruction loss of a chord progression is computed as a summation of 8 beat-wise chord loss using cross entropy functions [34]. For each beat, the chord loss is defined as the product of three parts: 1) the root loss, 2) the bass loss, and 3) the chroma loss. The root and bass are both considered 12-way categorical distributions and the chroma is regarded as 12 independent Bernoulli distributions.

3.3 Texture Encoder

The input of the texture encoder is an 8-beat segment of polyphonic piece represented by an image-like data format slightly modified from the piano-roll [14]. Each 8-beat segment is represented by a 128 by 32 matrix, where each row corresponds to a MIDI pitch and each column corresponds to $\frac{1}{4}$ beat. The data entry at (p, t) records the duration of the note if there is a note onset, and zero otherwise.

The texture encoder aims to learn a chord-invariant representation of texture by leveraging both the translation invariance property of convolution and the blurry effect of max-pooling layers [35]. We use a convolutional layer with kernel size 12×4 and stride 1×4 , which is followed by a ReLU activation [36] and max-pooling with kernel size 4×1 and stride 4×1 . The convolutional layer has one input channel and 10 output channels. The convolutional layer design aims at extracting a blurry “concept sketch” of the polyphonic texture which contains minimum information of the underlying chord. Ideally, when such blurry sketch is combined with specific chord representation, the decoder can identify its concrete pitches in a musical way.

The output of the convolutional layer is then fed into a bi-directional GRU encoder to extract the texture representation z_{txt} , similar to how we encode z_{chd} introduced in Section 3.1.

3.4 PianoTree Decoder

The PianoTree decoder takes the concatenation of z_{chd} and z_{txt} as input and decodes the music segment using the same decoder structure invented in PianoTree VAE [6], a hierarchical model structure for polyphonic representation learning. The decoder works as follows. First, it generates 32 frame-wise hidden states (one for each $\frac{1}{4}$ beat) using a GRU layer. Then, each frame-wise hidden state is further decoded into the embeddings of individual notes using another GRU layer. Finally, the pitch and duration for each note are reconstructed from the note embedding using a fully-connected layer and a GRU layer, respectively. For more detailed derivation and model design, we refer the readers to [6].

3.5 Training Objective

Let x denote the input music piece and $c = f(x)$ denote the chord progression extracted by algorithm $f(\cdot)$. We assume standard Gaussian priors of $p(z_{\text{chd}})$ and $p(z_{\text{txt}})$, and denote the output posteriors of chord encoder and texture encoder by $q_{\phi}(z_{\text{chd}}|c)$, $q_{\psi}(z_{\text{txt}}|x)$, the output distributions

of chord decoder and PianoTree decoder by $p_{\rho}(c|z_{\text{chd}})$ and $p_{\theta}(x|z_{\text{chd}}, z_{\text{txt}})$. The objective of the model is:

$$\begin{aligned} \mathcal{L}(\phi, \psi, \rho, \theta; x) = & \\ & - \mathbb{E}_{\substack{z_{\text{chd}} \sim q_{\phi} \\ z_{\text{txt}} \sim q_{\psi}}} [\log p_{\rho}(c|z_{\text{chd}}) + \log p_{\theta}(x|z_{\text{chd}}, z_{\text{txt}})] \\ & + \text{KL}(q_{\phi}||p(z_{\text{chd}})) + \text{KL}(q_{\psi}||p(z_{\text{txt}})). \end{aligned} \quad (1)$$

4. CONTROLLED MUSIC GENERATION

In this section, we show some controlled generation examples of the three tasks mentioned in the introduction.

4.1 Compositional Style Transfer

By regarding chord progression *content* and texture *style*, we can achieve compositional style transfer by swapping the texture representations of different pieces. Figure 2 shows the transferred results ((c) & (d)) based on two 16-bar samples ((a) & (b)) in the test set by swapping z_{txt} every 2 bars (without overlap)².

We see that such long-term style transfer is successful: The generated segment (c) follows the chord progression of (b) while mimicking the texture of (a), while (d) follows the chord progression of (a) while mimicking the texture of (b). As shown in the marked scores, the style transfer is effective. E.g., the cut-offs, melody contours, and the shape of the left-hand accompaniment are all preserved.

4.2 Texture Variation by Sampling

We can make variations of texture by sampling from z_{txt} while keeping z_{chd} . Here, we investigate two sampling strategies: sampling from the posterior $q_{\psi}(z_{\text{txt}}|x)$, and sampling from the prior $p(z_{\text{txt}})$.

Sampling from the posterior distribution $q_{\psi}(z_{\text{txt}}|x)$ yields reasonable variations as shown in Figure 3a. The variations of the right-hand melody can be seen as an improvisation following the chord progression and the melody. On the contrary, there is only small variation in the left-hand part, showing that the model regards the left-hand accompaniment as the dominant feature of texture.

Sampling from the prior distribution $p(z_{\text{txt}})$ changes the texture completely. Figure 3b shows a series of examples of prior sampling under the same chord progression C-Am-F-G. The resulting generations follow exactly the chord progression but with new textures.

4.3 Accompaniment Arrangement

We use a downstream predictive model to achieve accompaniment arrangement. For this task, we provide extra vocal melody tracks paired with the piano samples, and the model learns to generate 16-bar piano accompaniment *conditioned* on melody in a supervised fashion.

We encode the music every 2 bars (without overlap) into latent representations. For the accompaniment, we use the proposed model to compute the latent chord and texture representation, denoted by $\mathbf{z}_{\text{chd}} = [z_{\text{chd}}^{(1)}, \dots, z_{\text{chd}}^{(4)}]$ and $\mathbf{z}_{\text{txt}} = [z_{\text{txt}}^{(1)}, \dots, z_{\text{txt}}^{(4)}]$. For the melody, we use the

² The presented excerpts are converted from MIDI by the authors. The chord labels are inferred from the original/generated samples.

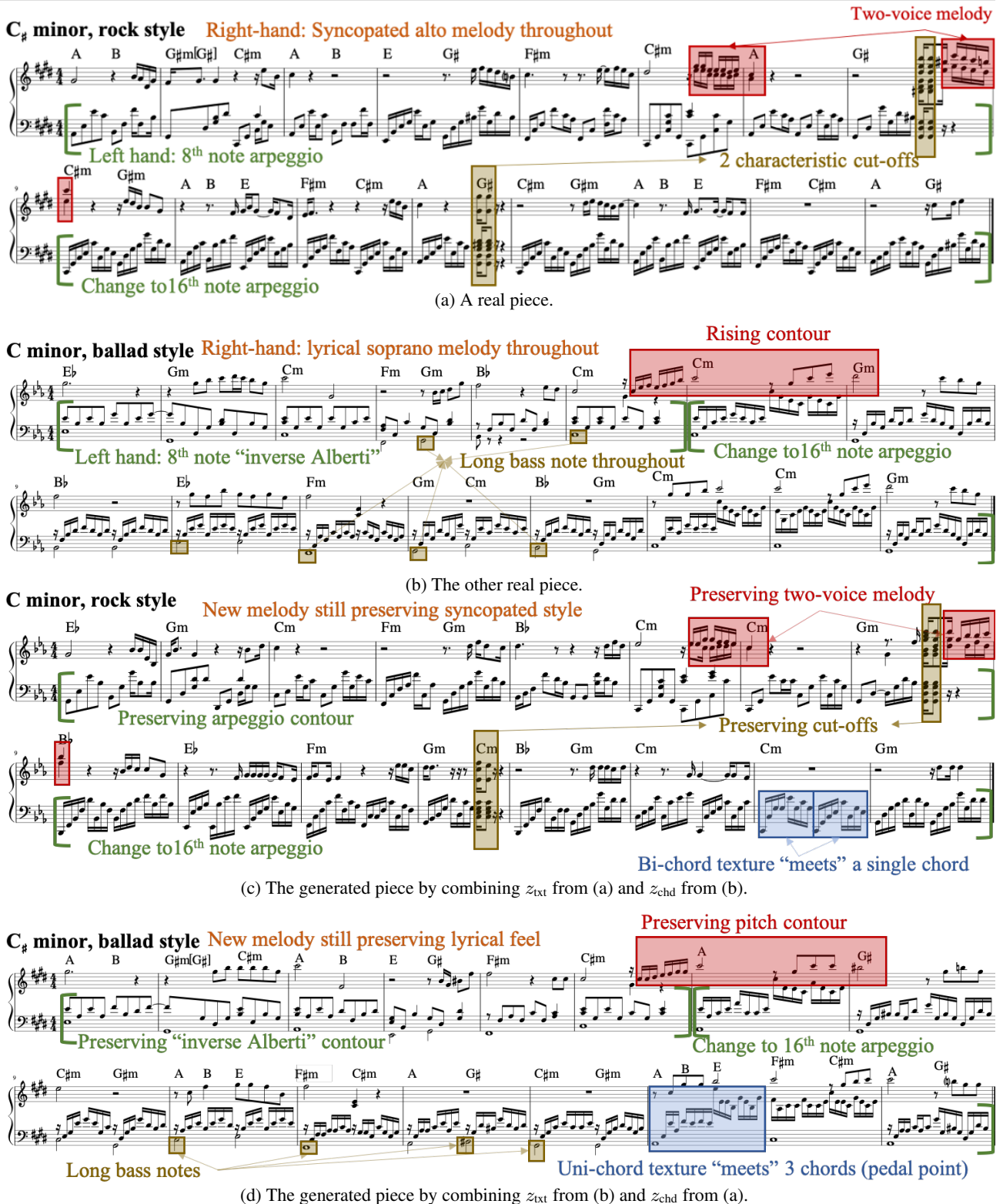
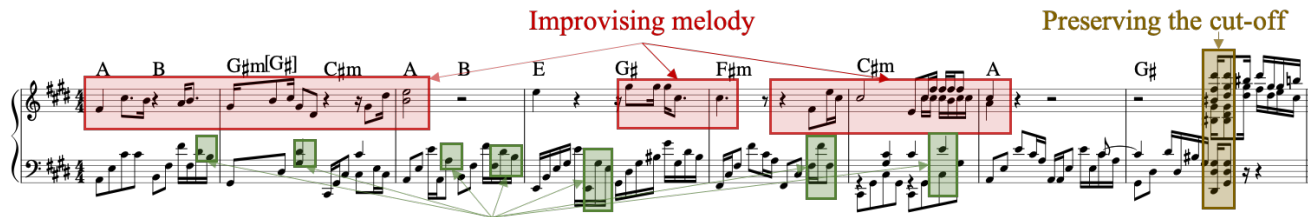


Figure 2: An example of compositional style transfer of 16-bar-long samples when $k = 2$.

EC²-VAE [27] to compute the latent pitch and rhythm representations, denoted by $\mathbf{z}_p = [z_p^{(1)}, \dots, z_p^{(4)}]$ and $\mathbf{z}_r = [z_r^{(1)}, \dots, z_r^{(4)}]$. Then, we adopt a vanilla Transformer [37] to model $p(\mathbf{z}_{\text{txt}}, \mathbf{z}_{\text{chd}} | \mathbf{z}_p, \mathbf{z}_r)$, in which the encoder takes in the condition and the decoder's input is a shifted right version $[z_{\text{chd}}, \mathbf{z}_{\text{txt}}]$. Both encoder and decoder inputs are incorporated with a *positional encoding* indicating the time po-

sitions and a learned *factor embedding* indicating the representation type (i.e., pitch, rhythm, chord or texture).

Figure 4 shows an example of accompaniment arrangement, where the first staff shows the melody and the second staff shows the piano accompaniment. In this case, the whole melody, together with the complete chord progression and the first 2 bars of accompaniment are given. The chord conditioning is done by forcing the decoded chord



(a) An example of posterior sampling of z_{txt} of the first 8 bars of the segment (a) in Figure 2



(b) An example of prior sampling of z_{txt} under given chord progression C-Am-F-G. Each two-bar segment is independently sampled, having different texture.

Figure 3: Examples of texture variations via posterior sampling and prior sampling.

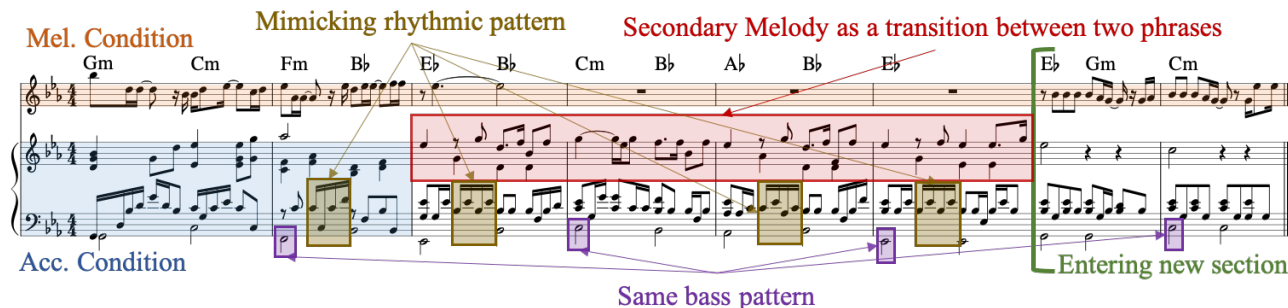


Figure 4: An example of accompaniment arrangement conditioned on melody, chord progression, and first 2 bars of accompaniment.

representation to match the given input during inference time. (A similar trick is used in [15].) From Figure 4, we see that the model predicts a similar texture to the given accompaniment. Moreover, it fills in a secondary melody line as a transition when the lead melody is rest.

Note that the arrangement can be generated in a flexible way by conditioning on different sets of latent factors. Much longer examples and more conditioning settings are available on our github page.

5. EXPERIMENTS

5.1 Dataset and Training

We train our model on the POP909 dataset [38], which contains about 1K MIDI files of pop songs (including paired vocal melody and piano accompaniment). We further extract the chord annotations using [31, 32]. We only keep the pieces with $\frac{2}{4}$ and $\frac{4}{4}$ meters and cut them into 8-beat music segments (so that each data sample in our experiment contains 32 time steps under 16th note resolution). In all, we have 66K samples. We randomly split the dataset (at song-level) into training set (90%) and test set (10%). All training samples are further augmented by transposing to all 12 keys.

In our experiment, the VAE model uses 256, 512, and 512 hidden dimensions for the GRUs in chord encoder, chord decoder and texture encoder respectively. The latent dimension of z_{chd} and z_{txt} are both 256. The model size of the PianoTree decoder is the same as the implementation

in the original paper [6]. The transformer model has the following size: hidden dimension = 256, number of layers = 4 and number of heads = 8.

For both models, we use Adam optimizer [39] with a scheduled learning rate from 1e-3 to 1e-5. Moreover, for the VAE model, we use KL-annealing [40], i.e. setting a weight parameter for the KL-divergence loss starting from 0 to 0.1. We set batch size to be 128 and the training converges within 6 epochs. For the downstream transformer model, we use 12K warm-up steps for learning rate update [41]. We use the same batch size and the model converges within 40 epochs.

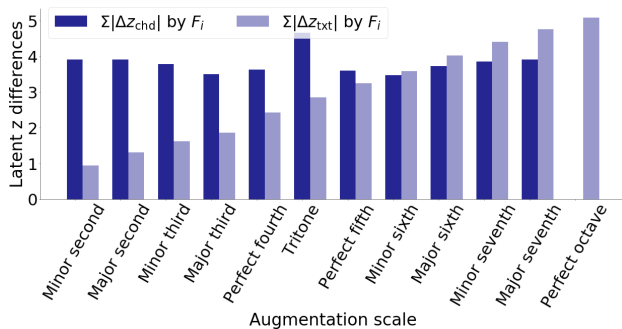
5.2 Objective Measurement

When z_{chd} and z_{txt} are well disentangled, small variations over the note pitches of the original music should lead to a larger change on z_{chd} , while variations of rhythm will influence more on z_{txt} . Following this assumption, we adopt a *disentanglement evaluation via data augmentation* method used in [42] and further developed in [27].

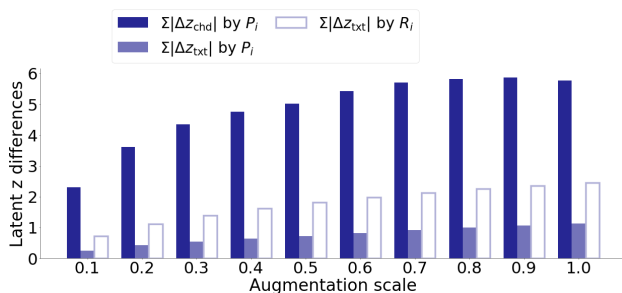
We define F_i as the operation of transposing all the notes by i semitones, and use the L_1 -norm to measure the change of latent z after augmentation. Figure 5a shows a comparison between $\sum|\Delta z_{\text{chd}}|$ and $\sum|\Delta z_{\text{txt}}|$ when we apply F_i to all the music pieces in the test set (where $i \in [1, 12]$).

It is conspicuous that when augmenting pitch in a small range, the change of z_{chd} is much larger than the change of z_{txt} . At the same time, the change of z_{txt} gets higher as the

augmentation scale increases. Similar to the result in [27], the change of z_{chd} reflects human pitch perception as z_{chd} is very sensitive to a tritone transposition, and least sensitive for a perfect octave.



(a) A comparison between Δz_{chd} , Δz_{txt} after pitch transposition on all notes.



(b) A comparison among Δz_{chd} , Δz_{txt} after beat-wise pitch transposition and texture augmentation with different probabilities.

Figure 5: Results of objective measurement.

We further define P_i as the function to randomly transpose all the notes in one beat either up or down one semitone under a certain probability i , and R_i as the function to randomly reduce the note duration by half. Figure 5b shows a comparison between $\Sigma|\Delta z_{\text{chd}}|$ and $\Sigma|\Delta z_{\text{txt}}|$ when we apply P_i and R_i to all the music pieces in our test set (where $i \in [0.1, 1.0]$).

For each value of i in the figure 5b, the first and second bars demonstrate $\Sigma|\Delta z_{\text{chd}}|$ and $\Sigma|\Delta z_{\text{txt}}|$ caused by P_i function, while the third bar indicates $\Sigma|\Delta z_{\text{txt}}|$ caused by R_i function. (We did not show $\Sigma|\Delta z_{\text{chd}}|$ caused by R_i since they are all zero.) It again proves that the chord representation is more sensitive than texture representation under pitch variations, and conversely, texture representation is more sensitive than chord representation under rhythm variations.

5.3 Subjective Evaluation

Besides objective measurement, we conduct a survey to evaluate the musical quality of compositional style transfer (see Section 4.1). Each subject listens to ten 2-bar pieces with different chord progressions, each paired with 5 style-transfer versions generated by swapping the texture representation with a random sample from the test set. In other words, each subject evaluates 10 groups of samples, each of which contains 6 versions of textures (1 from the original piece and 5 from other pieces) under the same chord progression. Both the order of groups and the sample order

within each group are randomized. After listening to each sample, the subjects rate them based on a 5-point scale from 1 (very low) to 5 (very high) according to three criteria: *creativity*, *naturalness* and *musicality*.

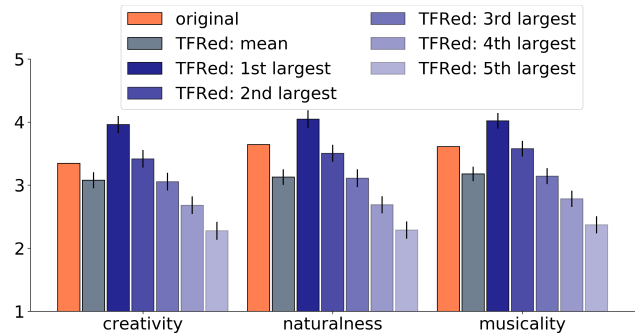


Figure 6: Subjective evaluation results. Here “TFRed: x^{th} largest” denotes the x^{th} (largest) order statistic of the transferred segments.

A total of 36 subjects (26 females and 10 males) participated in the survey. Figure 6 shows the comparison result among the original pieces (indicated by the orange bars) and the transferred pieces in terms of their mean and *order* statistics. The heights of bars represent averaged ratings across the subjects and the error bars represent the confidence intervals computed via paired t-test [43]. The result shows if we randomly transfer a piece’s texture 5 times, the best result is significantly better than the original version (with p -value < 0.005), and there are only marginal differences between the second-largest statistics and the original (with p -value > 0.05) in terms of creativity and musicality. We also see that on average the transferred results are still rated lower than the original ones. How to automatically decide the quality of a transferred result is considered a future work.

6. CONCLUSION AND FUTURE WORK

In conclusion, we contributed an effective algorithm to disentangle polyphonic music representation into two interpretable factors, chord and texture, under a VAE framework. Such interpretable representations serve as an intuitive human-computer co-creation interface, by which we can precisely manipulate individual factors to control the flow of the generated music. In this paper, we demonstrated three ways to interact with the model, including compositional style transfer via swapping the latent codes, texture variation by sampling from the latent distribution, accompaniment arrangement using downstream conditional prediction, and there are potentially many more. We hope this work can shed light on the field of controllable algorithmic composition in general, especially on the paradox between model complexity and model interpretability.

We acknowledge that the learned music factors are still very basic. In the future, we plan to extract more abstract and longer-range features using hierarchical models. We also plan to explore more ways to control the music generation for practical usage.

7. REFERENCES

- [1] K. Chen, W. Zhang, S. Dubnov, G. Xia, and W. Li, "The effect of explicit structure encoding of deep neural networks for symbolic music generation," in *2019 International Workshop on Multilayer Music Representation and Processing (MMRP)*. IEEE, 2019, pp. 77–84.
- [2] C. A. H. et al., "Music transformer: Generating music with long-term structure," in *7th International Conference on Learning Representations (ICLR), New Orleans, LA, USA, 2019*.
- [3] Y.-S. Huang and Y.-H. Yang, "Pop music transformer: Generating music with rhythm and harmony," *arXiv preprint arXiv:2002.00212*, 2020.
- [4] J.-P. Briot and F. Pachet, "Deep learning for music generation: challenges and directions," *Neural Computing and Applications*, vol. 32, no. 4, pp. 981–993, 2020.
- [5] S. Dai, Z. Zhang, and G. G. Xia, "Music style transfer: A position paper," *arXiv preprint arXiv:1803.06841*, 2018.
- [6] Z. Wang, Y. Zhang, Y. Zhang, J. Jiang, R. Yang, J. Zhao, and G. Xia, "Pianotree vae: Structured representation learning for polyphonic music," in *Proceedings of 21st International Conference on Music Information Retrieval (ISMIR), virtual conference, 2020*.
- [7] J.-P. Briot, G. Hadjeres, and F.-D. Pachet, "Deep learning techniques for music generation—a survey," *arXiv preprint arXiv:1709.01620*, 2017.
- [8] J.-P. Briot, "From artificial neural networks to deep learning for music generation—history, concepts and trends," *arXiv preprint arXiv:2004.03586*, 2020.
- [9] I. Simon, D. Morris, and S. Basu, "Mysong: automatic accompaniment generation for vocal melodies," in *Proceedings of the SIGCHI conference on human factors in computing systems*, 2008, pp. 725–734.
- [10] L.-C. Yang, S.-Y. Chou, and Y.-H. Yang, "Midinet: A convolutional generative adversarial network for symbolic-domain music generation," *arXiv preprint arXiv:1703.10847*, 2017.
- [11] K. Chen, W. Zhang, S. Dubnov, G. Xia, and W. Li, "The effect of explicit structure encoding of deep neural networks for symbolic music generation," in *2019 International Workshop on Multilayer Music Representation and Processing (MMRP)*. IEEE, 2019, pp. 77–84.
- [12] G. Hadjeres, F. Pachet, and F. Nielsen, "Deepbach: a steerable model for bach chorales generation," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017, pp. 1362–1371.
- [13] H. Zhu, Q. Liu, N. J. Yuan, C. Qin, J. Li, K. Zhang, G. Zhou, F. Wei, Y. Xu, and E. Chen, "Xiaoice band: A melody and arrangement generation framework for pop music," in *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2018, pp. 2837–2846.
- [14] H.-W. Dong, W.-Y. Hsiao, L.-C. Yang, and Y.-H. Yang, "Musegan: Multi-track sequential generative adversarial networks for symbolic music generation and accompaniment," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [15] C. Donahue, H. H. Mao, Y. E. Li, G. W. Cottrell, and J. McAuley, "Lakhnes: Improving multi-instrumental music generation with cross-domain pre-training," *arXiv preprint arXiv:1907.04868*, 2019.
- [16] I. Simon, A. Roberts, C. Raffel, J. Engel, C. Hawthorne, and D. Eck, "Learning a latent space of multitrack measures," *arXiv preprint arXiv:1806.00195*, 2018.
- [17] S. Huang, Q. Li, C. Anil, X. Bao, S. Oore, and R. B. Grosse, "Timbretron: A wavenet (cyclegan (cqt (audio))) pipeline for musical timbre transfer," *arXiv preprint arXiv:1811.09620*, 2018.
- [18] O. Kwon, I. Jang, C. Ahn, and H.-G. Kang, "Emotional speech synthesis based on style embedded tacotron2 framework," in *2019 34th International Technical Conference on Circuits/Systems, Computers and Communications (ITC-CSCC)*. IEEE, 2019, pp. 1–4.
- [19] S. Lattner, M. Grachten, and G. Widmer, "Imposing higher-level structure in polyphonic music generation using convolutional restricted boltzmann machines and constraints," *arXiv preprint arXiv:1612.04742*, 2016.
- [20] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [21] A. Roberts, J. Engel, C. Raffel, C. Hawthorne, and D. Eck, "A hierarchical latent vector model for learning long-term structure in music," *arXiv preprint arXiv:1803.05428*, 2018.
- [22] F. Locatello, S. Bauer, M. Lucic, G. Rätsch, S. Gelly, B. Schölkopf, and O. Bachem, "Challenging common assumptions in the unsupervised learning of disentangled representations," *arXiv preprint arXiv:1811.12359*, 2018.
- [23] Y.-J. Luo, K. Agres, and D. Herremans, "Learning disentangled representations of timbre and pitch for musical instrument sounds using gaussian mixture variational autoencoders," *arXiv preprint arXiv:1906.08152*, 2019.
- [24] Y. Wu, T. Carsault, E. Nakamura, and K. Yoshii, "Semi-supervised neural chord estimation based on a variational autoencoder with discrete labels and continuous textures of chords," *arXiv preprint arXiv:2005.07091*, 2020.

- [25] T. Akama, “Controlling symbolic music generation based on concept learning from domain knowledge,” in *ISMIR*, 2019, pp. 816–823.
- [26] K. Choi and K. Cho, “Deep unsupervised drum transcription,” *arXiv preprint arXiv:1906.03697*, 2019.
- [27] R. Yang, D. Wang, Z. Wang, T. Chen, J. Jiang, and G. Xia, “Deep music analogy via latent representation disentanglement,” *arXiv preprint arXiv:1906.03626*, 2019.
- [28] G. Brunner, A. Konrad, Y. Wang, and R. Wattenhofer, “Midi-vae: Modeling dynamics and instrumentation of music with applications to style transfer,” *arXiv preprint arXiv:1809.07600*, 2018.
- [29] S. Lattner, M. Dörfler, and A. Arzt, “Learning complex basis functions for invariant representations of audio,” *arXiv preprint arXiv:1907.05982*, 2019.
- [30] M. F. Mathieu, J. J. Zhao, J. Zhao, A. Ramesh, P. Sprechmann, and Y. LeCun, “Disentangling factors of variation in deep representation using adversarial training,” in *Advances in neural information processing systems*, 2016, pp. 5040–5048.
- [31] B. Pardo and W. P. Birmingham, “Algorithms for chordal analysis,” *Computer Music Journal*, vol. 26, no. 2, pp. 27–49, 2002.
- [32] C. Raffel, B. McFee, E. J. Humphrey, J. Salamon, O. Nieto, D. Liang, D. P. Ellis, and C. C. Raffel, “mir_eval: A transparent implementation of common mir metrics,” in *In Proceedings of the 15th International Society for Music Information Retrieval Conference, ISMIR*. Citeseer, 2014.
- [33] X. Yan, J. Yang, K. Sohn, and H. Lee, “Attribute2image: Conditional image generation from visual attributes,” in *European Conference on Computer Vision*. Springer, 2016, pp. 776–791.
- [34] P.-T. De Boer, D. P. Kroese, S. Mannor, and R. Y. Rubinstein, “A tutorial on the cross-entropy method,” *Annals of operations research*, vol. 134, no. 1, pp. 19–67, 2005.
- [35] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [36] V. Nair and G. E. Hinton, “Rectified linear units improve restricted boltzmann machines,” in *Proceedings of the 27th international conference on machine learning (ICML-10)*, 2010, pp. 807–814.
- [37] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in neural information processing systems*, 2017, pp. 5998–6008.
- [38] Z. Wang, K. Chen, J. Jiang, Y. Zhang, M. Xu, S. Dai, X. Gu, and G. Xia, “Pop909: A pop-song dataset for music arrangement generation,” in *Proceedings of 21st International Conference on Music Information Retrieval (ISMIR), virtual conference*, 2020.
- [39] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [40] S. R. Bowman, L. Vilnis, O. Vinyals, A. M. Dai, R. Jozefowicz, and S. Bengio, “Generating sentences from a continuous space,” *arXiv preprint arXiv:1511.06349*, 2015.
- [41] R. Xiong, Y. Yang, D. He, K. Zheng, S. Zheng, C. Xing, H. Zhang, Y. Lan, L. Wang, and T.-Y. Liu, “On layer normalization in the transformer architecture,” *arXiv preprint arXiv:2002.04745*, 2020.
- [42] H. Kim and A. Mnih, “Disentangling by factorising,” *arXiv preprint arXiv:1802.05983*, 2018.
- [43] H. Hsu and P. A. Lachenbruch, “Paired t test,” *Encyclopedia of Biostatistics*, vol. 6, 2005.