

MusIAC: An extensible generative framework for Music Infilling Applications with multi-level Control

Rui Guo¹[0000-0003-2907-0718], Ivor Simpson²[0000-0001-5605-6626], Chris Kiefer¹[0000-0002-3329-1938], Thor Magnusson¹[0000-0002-3731-0630], and Dorien Herremans³[0000-0001-8607-1640]

¹ Department of Music, University of Sussex r.guo@sussex.ac.uk

² Department of Informatics, University of Sussex

³ Information Systems Technology and Design, Singapore University of Technology and Design

Abstract. We present a novel music generation framework for music infilling, with a user friendly interface. Infilling refers to the task of generating musical sections given the surrounding multi-track music. The proposed transformer-based framework is extensible for new control tokens as the added music control tokens such as tonal tension per bar and track polyphony level in this work. We explore the effects of including several musically meaningful control tokens, and evaluate the results using objective metrics related to pitch and rhythm. Our results demonstrate that adding additional control tokens helps to generate music with stronger stylistic similarities to the original music. It also provides the user with more control to change properties like the music texture and tonal tension in each bar compared to previous research which only provided control for track density. We present the model in a Google Colab notebook to enable interactive generation.

Keywords: music generation · music transformer · music control · controllable generation · music representation · infilling

1 Introduction

Music composition by artificial intelligence (AI) methods, especially using deep learning, has been an active topic of research in recent years [3, 19]. In a recent survey of musical agents [30], twelve musical tasks such as accompaniment generation, melody/rhythm generation, continuation, arrangement, and style imitation are examined. In the deep learning era, all of these tasks have been explored to some extent.

When applying AI to music composition, however, an often ignored question is “why” one might wish computers to compose music. From the perspective of the deep learning practitioner, the answer may be to explore the limits of AI models for creative tasks, and investigate whether they can generate music as human composers. On the other hand, musicians and composers may want

to use AI as a source of inspiration, for instance, by rapidly offering several solutions. One such AI method is music infilling or inpainting [24, 12]. It is used to extend pre-existent music materials, such as filling in the missing bars or tracks given the surrounding music information. It can write a new melody line given the existing bass and accompaniment track, or rewrite a few bars in the middle given the beginning and the end. Many reasonable solutions may exist that match the surrounding music progression and harmony. Without efficient track and bar music property conditions, however, the user has to generate the music repeatedly until it satisfies user’s requirement.

Several research studies have used a transformer model [31] for symbolic music generation [17, 18, 27, 9, 15, 12] and the results are promising. However, controlling the generation process is still limited in these approaches.

One common control for the music infilling system is track density[9, 12], which is defined as the number of notes in a track divided by the total timesteps in that track. However, a sole density cannot easily change the accompaniment track from a monophonic style to a polyphonic style. A polyphony control can help to convert a monophonic track such as arpeggio to a more polyphonic texture such as a chord track or vice versa in a direct way, and that can be useful mostly for the accompaniment track. Another interesting control is the track occupation rate, which determines which ratio of a track is note occupied versus filled with rests. These track features may be useful as a composer may want to control the track texture.

Except for those track controls, a bar level tonal tension control[5, 13] can help to create music with specific tension movements, e.g. from low to high, high to low or any tension shape. One use case is to change the tension of the beginning and ending of a piece so as to set certain moods.

To implement these controls, the respective track/bar properties are calculated and added to the input. We deliberately choose to use higher level human interpretable parameters as controls, including six features: key, bar tensile strain, bar cloud diameter, track density, track polyphony, and track occupation, and they are calculated from the input sequence directly. It may be useful to generate music according to the track/bar control parameter template fit to a particular scenario, such as high track note density, low track polyphony rate and high track occupation. As the model learns the relationship between these control tokens and the music, the controls can be changed to generate variations of the original music. In the experiments, we observe that an additional benefit of including more music properties in the input is that the generated music is more similar to the original music measured by pitch and rhythm related metrics.

In this paper, we propose an extensible framework for music generation by infilling reconstruction. Six musically meaningful control tokens are calculated and added to the original input. The effect of adding this conditioning information is examined in an experiment that uses seven objective metrics selected from the literature. Our simple model design makes it extensible so that we can easily include additional tokens in the future. The music infilling task, which involves reconstructing a missing segment of music, is used to validate our results

by comparing the properties of original and generated examples. The results show that the model with added calculated music tokens to the input has more stylistic similarity to the original music. Google Colab notebook is shared for free exploration of this infilling system and gives a straightforward way to explore the effect of adding novel musically meaningful tokens to the input. The music generated by changing control tokens demonstrates the controllability of this method.

2 Related Work

Over the years, many different generative models have been developed for symbolic music generation [19, 14]. Variational AutoEncoder(VAE) based models [29, 24, 11] usually generate short music pieces and explore different music features in the latent space. Generative Adversarial Network (GAN) based models [4] can generate longer music, but can be harder to train and may suffer mode collapse without careful parameter tuning [21]. Recursive Neural Networks [34], and more recently the powerful transformer based methods[17] can generate long music pieces but with less control explored so far compared to the VAE models.

Several improvements have been made since the transformer model was first used for music generation, related to both the input representation and the model structure. [18] uses “position”(timestep) and “duration” tokens to replace the “note on” and “note off” tokens [22]. This allows the model to learn to pair the “note on” and “note off” if they are far apart. [27] generates accompaniment given the melody track, and adds pitch, velocity, and duration embeddings in one timestep. [15] has a similar design and uses different output linear layers for different token types. The models by [9, 12] generate music infillings similar to the task tackled in this research. Both models take the track note density as the control parameter, without providing any other track/bar level control features, we will explore adding the latter features in this research.

Some interactive interfaces have previously been designed specifically for the music infilling task. [2] and [20]’s interfaces are based on music chorale generation [16].

3 Proposed model and representation

The existing transformer-based methods offer few semantic controls for the generation process, or focus on continuous generation rather than infilling. Given the right input controls, music generation models may be driven and steered by relevant musical concepts and ideas. Our work is based on the following assumptions:

1. Including additional derived musical features can improve the performance of music infilling.
2. Using human interpretable music features allows the user to control the generated music.

Because the music infilling region is the model’s prediction objective, it is natural to compare the generated music to the original. If the generated music has similar properties to the original infilled music region, then the model has performed well. Our model is versatile enough to allow multiple types of infilling. For instance, in pop music with multiple tracks, the infilling can work either by infilling a whole track or by infilling a bar across tracks, or both at the same time. Fig. 1 shows an example of how we can formulate the infilling task. The input music has three tracks, the yellow block region masks the first track, and the blue block region masks the second bar. The aim of the model here is to reconstruct the masked music region given the surrounding information. Providing input with multiple tracks makes it possible to have the track properties separately, and the control for different tracks can be tested separately.

Fig. 1: An example of original music with two infilled regions. The yellow block masks the melody track, and the blue block masks the second bar. The notes of those two masked regions are replaced by a “mask” token in the model input. The target output of the model is to reconstruct the missing track/bar in the infilled region.

3.1 Adding control features

We selected the following information to be added to the model input as controls from multiple levels. This is calculated from the MIDI data and provides high level musical concepts as conditions for the music generation.

1. Track level controls:
 - The track’s note density rate: $number_{note}/timesteps_{total}$. This is calculated by dividing the number of notes in a track by the maximum time steps in that track.
 - The track’s note polyphony rate: $timesteps_{polynote}/timesteps_{anynote}$. This is the number of timesteps with more than two notes divided by the total number of timesteps with any note.
 - The track’s note occupation rate: $timesteps_{anynote}/timesteps_{total}$. This is the total number of timesteps with any note divided by the total number of time steps, including those with rests.

2. Bar level controls:
 - The tensile strain [13] of the notes in that bar: $\sum_{i=1}^n (note_{pos}[i] - key_{pos})/n$, which is the average of the difference from the note position to the key position. The note and key position are calculated based on the spiral array theory [5]. This is a tonal tension measure.
 - The cloud diameter [13] of the notes in that bar: $\max_{i \in [1..n-1], j \in [i+1..n]} (note_{pos}[i] - note_{pos}[j])$. This is another tonal tension measure, which only calculates the largest distance between notes in that bar. The calculation of the note position is also based on the spiral array theory.

Except for the above controls, the following information is also added to the model’s input as auxiliary information. The key is calculated by [7, 10]. The tempo, time signature, and track instrument number are extracted directly from the MIDI files.

- The key of the song, which can be one of 24 keys (major and minor).
- The tempo of the song, categorised into seven different bins.
- The time signature of the song, including 4/4, 3/4, 2/4, and 6/8.
- The track’s instrument: The MIDI instrument number.

3.2 Data representation

We use an adapted version of the “REMI” [18] token representation in this work. The “REMI” format includes position tokens to mark the position of the note inside a bar. The number of the position is related to the minimum duration of the note. We select the 16th note as the minimum note length, and a bar in 4/4 metre is thus divided into 16 different start positions. The position tokens range from “e_0” to “e_15”, and the duration tokens range from “n_1” to “n_32”. The maximum note length “n_32” represents two whole notes in the time signature of 4/4. The pitch tokens range from “p_21” to “p_108”, representing A-1 to C7 respectively. There is a “bar” token to mark the start of a new bar. The velocity, tempo, and chord tokens proposed in [18] are discarded in the format used here. The dynamics of music is not the focus of this research, and by removing the velocity of each note, notes with the same duration can be grouped by using only one duration token after the pitch tokens. E.g. e_0, p_60, p_67, n_10 means note C3 and G3 have the same duration (10×16 th note), which equals the summation of a half note (8×16 th note) and an eighth note (2×16 th note). Because the tonal tension information is included, the chord information is also removed.

To represent the “track” concept, a “track” token is added to the vocabulary list, similar to [27]. Up to three tracks are used in this work: “track_0” is the melody track, “track_1” is the bass track, and “track_2” is an accompaniment track. The track token is the first token of all the tokens in that track. More tracks can be added in the future if they are arranged in the same order, e.g. track_3 for drum and track_4 for a second melody.

Fig. 2 shows a piece with three tracks. Before the calculated control information is added, the event list is: 4/4, t_3, i_0, i_32, i_48, bar, track_0, e_0, p_

79,n_4,e_4,p_76,n_4,e_8,p_74,n_6,track_1,e_0,p_45,n_8,e_8,p_41,n_8,track_2,e_0,p_64,p_67,n_8,e_0,p_60,n_16,e_8,p_65,n_8,bar,track_0,e_0,p_69,n_4,e_4,p_71,n_4,e_8,p_72,n_6,track_1,e_0,p_43,n_8,e_8,p_48,n_8,track_2,e_0,p_59,p_65,p_67,n_8,e_8,p_60,p_64,n_8.

The image shows a musical score for three instruments: Steinway Grand Piano, Fingerstyle Bass, and String Ensemble. The score is written in 4/4 time and consists of two measures. The Steinway Grand Piano part is in the treble clef, the Fingerstyle Bass part is in the bass clef, and the String Ensemble part is in the treble clef. The Steinway Grand Piano part features a melody of eighth notes. The Fingerstyle Bass part features a bass line of quarter notes. The String Ensemble part features a chordal accompaniment of quarter notes.

Fig. 2: Example of a musical segment in our dataset.

The control information that is included in our proposed framework is tensile strain (12 categories), cloud diameter (12 categories), track density/ polyphony/occupation rate (each for 10 categories) as per the previous subsection. Because the calculation of the bar tonal tension is based on a determined key, the key of the song is also determined and added to the music input. After those calculated control tokens are added, the data representation for Fig. 2 becomes: 4/4, t_3, k_0, d_0, d_0, d_0, o_8, o_9, o_9, y_0, y_0, y_9, i_0, i_32, i_48, bar, s_2, a_1, track_0, e_0, p_79, n_4, e_4, p_76, n_4, e_8, p_74, n_6, track_1, e_0, p_45, n_8, e_8, p_41, n_8, track_2, e_0, p_64, p_67, n_8, e_0, p_60, n_16, e_8, p_65, n_8, bar, s_5, a_6, track_0, e_0, p_69, n_4, e_4, p_71, n_4, e_8, p_72, n_6, track_1, e_0, p_43, n_8, e_8, p_48, n_8, track_2, e_0, p_59, p_65, p_67, n_8, e_8, p_60, p_64, n_8. The tokens at the start of the event list are time signature, tempo, and key tokens. The track control tokens appear after the key token, followed by the instrument tokens. A “bar” token follows the instrument token, immediately followed by tension control. The “track” token is followed by the “position”, “pitch” and “duration” tokens inside each track. The final vocabulary list is represented in Table 1.

3.3 Model architecture

As the core task here is music infilling rather than forward generation, the model should ideally use bidirectional information. The transformer encoder-decoder model [31] which was originally developed for the seq-seq translation task, is adapted in this work. The infilling task in music can be likened to the corrupted token reconstruction task in natural language processing[8]. In our proposed framework, a transformer encoder-decoder is used to reconstruct the masked input in the encoder[28]. The bi-directional encoder makes each token in the

Table 1: The event vocabulary, including all calculated control tokens.

token types	tokens	number
position	e_0...e_15	16
pitch	p_21...p_108	88
duration	n_1...n_32	32
structure tokens	bar, track_0, track_1, track_2	4
time signature	4/4, 3/4, 2/4, 6/8	4
tempo	t_0...t_6	7
instrument	i_0...i_127	128
key	k_0...k_23	24
tensile strain	s_0...s_11	12
cloud diameter	a_0...a_11	12
density	d_0...d_9	10
polyphony	y_0...y_9	10
occupation	y_0...o_9	10
model	mask, pad, eos	3
total		360

encoder attend to other positions in the input, while the token in a one stack decoder language model can only attend to the tokens before the current token[26]. Our model has the same structure as the vanilla transformer [31] with two stages of training. Firstly, music grammar is learned in the pretraining stage and then specific tasks are learned in the finetuning stage. This process is similar to the work of [8, 33, 6].

During pretraining, we accustom the model to small masked sections: one “mask” token can replace up to three tokens. If the input x position from u to v is masked, and the $l = u - v$ is the masked token span length, the loss function is calculated as in Eq. (1):

$$L(\theta) = \log P(x^{u:v} | x^{\setminus u:v}; \theta), \quad 0 < u - v \leq 3. \quad (1)$$

Up to 15% of the tokens in the input are randomly masked with a “mask” in pretraining. We only use one “mask” token to replace each span, which differs from other work [26] which uses a different mask token for each span masked. The lengths of the spans of the masked token are 3, 1, 2 and the frequency of the masked tokens with those span lengths is in the ratio of 2:1:1 in the training respectively.

After pretraining, the finetuning stage is used to train the model for the real application task with larger masked areas). The finetuning task includes three masking types corresponding to the application. For each song: 1. randomly select a bar, and mask all tracks in that bar. 2. randomly select a track, and mask all the bars in selected tracks. 3. randomly select bars, and randomly select tracks in that bar.

One “mask” token represents a track in a bar, and the decoder target is to reconstruct that masked bar track. Each “mask” in the encoder input is matched

with a “mask” input in the decoder, and the decoder target output will end with an “eos” token. A “pad” token is also added to pad sequences of different lengths to match the batch size. Fig. 3 shows masked encoder input and the decoder input and target output during pretraining/finetuning. During finetuning, if the first bar of Fig. 2 is infilled, the encoder input becomes: 4/4, t_3, k_0, d_0, d_0, d_0, o_8, o_9, y_0, y_0, y_9, i_0, i_32, i_48, bar, s_2, a_1, mask, mask, mask, bar, s_5, a_6, track_0, e_0, p_69, n_4, e_4, p_71, n_4, e_8, p_72, n_6, track_1, e_0, p_43, n_8, e_8, p_48, n_8, track_2, e_0, p_59, p_65, p_67, n_8, e_8, p_60, p_64, n_8. The decoder input is: mask, mask, mask, and the decoder target output is track_0, e_0, p_79, n_4, e_4, p_76, n_4, e_8, p_74, n_6, eos, track_1, e_0, p_45, n_8, e_8, p_41, n_8, eos, track_2, e_0, p_64, p_67, n_8, e_0, p_60, n_16, e_8, p_65, n_8, eos. We omitted the second bar’s tokens to save page space.

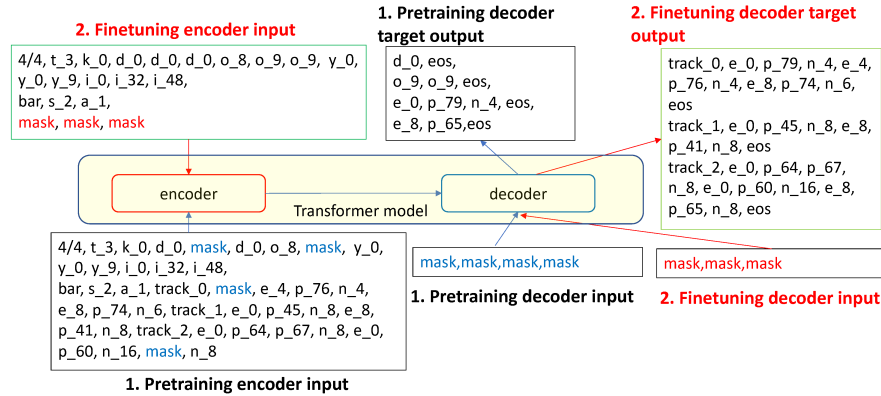


Fig. 3: The model encoder input, decoder input and decoder target output during pretraining and finetuning. The maximum masked span length is three for pretraining, and a “mask” token will replace a track in a bar during the finetuning stage.

4 Experimental setup

We conducted an experiment to validate the musical quality of the output as well as the influence of the control features. Two models with the same vocabulary size were trained in the experiment: one with controls and one without. The model without control will not add the six calculated controls to the input.

4.1 Dataset

Any dataset used by our proposed model should have proper track/bar numbers. The maximum number of tracks in a song is limited to three, which includes mandatory melody and bass tracks, and an optional accompaniment track. The maximum bar length in a song is set to 16, which is enough for the infilling reconstruction and not too long for the model to compute.

To build our symbolic dataset, we filter the Lakh LMD-matched dataset [25] for songs that have both a melody and bass track, as well as an optional accompaniment track. After that, the key of the song and the tension metrics are calculated using [7, 10]. A total of 32,352 songs remain after this step. To tackle the imbalance of the time signature in the remaining files, a subset with metre 2/4, 3/4 and 6/8 is pitch-shifted to the other 11 keys in the same mode. The same strategy is applied to a subset of songs with minor keys. A moving window size of 8 bars is used to create multiple dataset entries from a single song. All the calculated control features are added in this step.

4.2 Model configuration and training

One model is trained with all input and control tokens, the second model is trained without adding the control tokens. Both the encoder and decoder of the transformer have 4 layers, 8 heads, and the model dimension is 512. Both the models are trained for 10 epochs with 2 epochs of pretraining and the remaining 8 epochs for fine-tuning. The learning rate is 0.0001, and the training takes around 2 days per model on a Tesla V100 GPU. The training/validation/test data split ratio is 8:1:1.

4.3 Inference strategy

Our token representation allows us to guide the model to only generate output that adheres to the representation’s grammar. The grammar of the notes in a track in the regular expression format is $(\text{step pitch} + \text{duration})^*$. In the inference stage, the tokens not in this grammar are removed by setting those notes’ logit to -100, and then weighted sampling is applied to sample from the tokens. This makes sure that the output will not sample the incorrect tokens and the result always has a correct grammar.

5 Evaluation

To evaluate the generated music infillings generated by the model with and without controls, we select seven objective metrics based on pitch and rhythm similarity. We compare the difference of those features between the generated and the original music in the masked position. Then we check if our model can really control features of the generated music by changing track/bar controls through our developed Google Colab interface. Our experiment evaluates if the generated music follows the desired control features and is musically meaningful.

5.1 Objective evaluation using selected metrics

To compare the quality of the generated infilling by those two models, we selected five pitch-related metrics and two duration-related metrics inspired by [32]. The infilling generation task makes it meaningful to compare the metrics’ difference between the generated infilling and the original music. A smaller difference means the generated infilling has more stylistic similarity to the original music. Note that there is not only one optimal way to infill music, and we assume the original one is the target here. In future work, this assumption may be tested by allowing for a human listening experiment to evaluate the generated infillings. Both the track and bar infilling are evaluated.

We selected 1,000 songs randomly from the testset and masked a random track/bar, to test each of the two models. The models then generate the infilling for the masked track/bar. Seven objective metrics are selected inspired by [32] including five pitch related metrics: 1) pitch number: the number of used pitches. 2) note number: the number of used notes. 3) pitch range: $pitch_{max} - pitch_{min}$. 4) chromagram histogram: the histogram of 12 pitch groups after categorising all the pitches 5) pitch interval histogram: the histogram of the pitch difference between two consecutive notes. Two duration features: 6) duration histogram. 7) onset interval histogram: the histogram of time between two consecutive notes. These seven features are calculated for the generated/original infilled track/bar. For the first three features we calculate the absolute difference between the feature for the generated and original music, normalised by the feature of the original music: $abs(feature_{gen} - feature_{ori})/feature_{ori}$. For the last four histogram features we calculate the sum of the square difference between the features of the generated and the original music, normalised by the sum of the square of the feature of the original music: $sum(square(feature_{gen} - feature_{ori}))/sum(square(feature_{ori}))$.

The mean and the standard deviation are calculated on those difference features and reported in Table 2. The left value in each cell is the result for the model without added control tokens, and the right value is the result for the model with added control tokens. All of the values, except the track pitch number standard deviation, show that the model with added control generates music more similar to the original music, especially in terms of melody, accompaniment track, and bar infilling. The added control work much like a template, and the generated music follows these conditions well.

5.2 The interactive interface and controllability

A Google Colab notebook has been prepared for the exploration of this application ⁴. The user can upload MIDI files or select MIDI files from the test dataset. Here “Imagine” from John Lennon is selected from the test dataset as an example.

After selecting/uploading the song, the user can choose to infill a track/bar without changing the control tokens, or change the track/bar controls first, and

⁴ <https://github.com/ruiguo-bio/MusIAC>

Table 2: The mean and standard deviation of the difference for the seven objective metrics between the generated and original music. The left value in each cell is the result from the model without added control tokens, and the right value is the result from the model with added control. The column header shows was infilled: melody track, bass track, accompaniment track, or a random bar (all tracks in this bar).

features	Melody		Bass		Accompaniment		Bar	
	mean	std	mean	std	mean	std	mean	std
pitch number	0.45/0.39	0.57/0.46	0.55/0.52	0.76/0.78	0.57/0.41	0.69/0.52	0.41/0.33	0.79/0.58
note number	0.82/0.29	2.71/0.85	0.63/0.29	1.25/0.71	0.97/0.41	2.41/0.73	0.42/0.32	0.90/0.71
pitch range	0.59/0.49	0.93/0.74	0.62/0.58	0.92/0.87	0.80/0.55	1.05/0.76	0.55/0.38	2.55/1.54
chroma hist	0.59/0.53	0.49/0.42	0.45/0.38	0.50/0.44	0.34/0.27	0.42/0.31	0.73/0.58	0.79/0.70
pitch itv hist	0.44/0.39	0.46/0.40	0.50/0.45	0.55/0.53	1/0.84	0.75/0.63	0.66/0.60	0.77/0.75
duration hist	0.55/0.42	0.61/0.46	0.77/0.61	0.77/0.68	1.03/0.74	0.93/0.78	0.52/0.45	0.64/0.61
onset itv hist	0.45/0.35	0.55/0.41	0.74/0.61	0.71/0.70	0.86/0.69	0.85/0.77	0.44/0.38	0.63/0.58

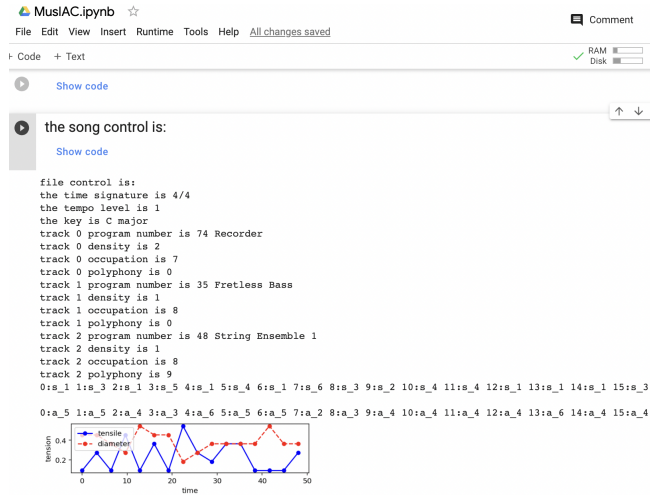


Fig. 4: The original music’s control information including track/bar control and song key, tempo and time signature.

then generate the corresponding track/bar. The original control tokens of one section of “Imagine” are calculated and shown as in the left figure of Fig. 6. The melody and accompaniment tracks have low note density, which means there are not many note onsets in that track. The accompaniment track is composed of mainly chord notes. The track/bar control can be changed by selecting the specific control type and value as shown in Fig. 5(only a section is shown due to page limitations). To add more notes to those tracks, and make the accompaniment track less polyphonic, we first change the melody track density to level 8 from level 1. After the new melody track is generated, the accompaniment track

is generated with density set to level 5 from level 1 and the polyphony level set to level 2 from level 9. The generated result is shown in the right figure in Fig. 6. The resulting music matches the desired control with a tolerance threshold of 1 (which means level 3 is accepted if the target level is 4 or 2). The resulting accompaniment track’s texture is similar to Alberti bass, and both of the tracks have more notes added following the increase of the track density level control.



Fig. 5: The track/bar controls can be changed separately

Based on the previous track’s infilling result, the first bar’s tensile strain is changed from level 1 to level 6 to increase the tension of the beginning. The infilled result is shown in Fig. 7. The first bar contains the subdominant of the F major chord, which is the second bar’s chord. This new first bar, together with the following two bars gives the progression of $IV/IV \rightarrow IV \rightarrow I$, which is musically meaningful (from subdominant to tonic chord), and it also increases the tension of the first bar. The full 16 bars music/sheet of the original/generated music are in the supplement material.

The track/bar infilling operation can be repeated several times until the result is satisfactory. The generated MIDI/rendered wav file can be downloaded for further editing or imported in a digital audio workstation (DAW).

6 Conclusion

In this work, we propose an pretraining-finetuning transformer framework for the music infilling task with multiple levels of control, together with an intuitive interface. We selected track density, polyphony, and occupation level as track level controls to increase the user’s controllability. This offers a greater steerability compared to existing systems with just density control[9, 12]. We also added tensile strain and cloud diameter features per bar as controls for the tonality (tonal tension) of each bar. Control tokens work as a template on which the generated music can be conditioned. The generated result from the input with those added controls as conditions has a higher stylistic similarity to the original music, versus a model without controls.

Fig. 6: The first three bars of a section of “Imagine”. The left figure is the original, and the right figure shows the infilled melody and accompaniment track with changed track density level from 1 to 5 and polyphony rate from 9 to 2. The original melody track has a low note density level of 1. The accompaniment track has low note density level 1 and high polyphony level 9. The infilled melody/accompaniment track match the selected controls, and the accompaniment is similar to Alberti bass, with more notes and less polyphony.

To optimally demonstrate our proposed framework with a user-friendly interactive interface, we have made it available through Google Colab. In this interface, the user can modify the music while it is being generated.

In the future work, we will systematically check the controllability of each of the six control tokens and further evaluate the musicality with quantitative metrics. A listening test would also be useful to evaluate the musical quality, as there may be more good sounding possibilities than just the original music. We would also like to explore how to further increase the controllability of this model. Currently, our model learns to follow controls (i.e., features) that are already present or easy to calculate from our dataset. It is hard for the model to generate music with “unseen” musical features, i.e. hard to capture, implicit characteristics. In recent research, a transformer model was combined with a prior model to model the latent space [1]. If different music properties can be disentangled in the latent space[23], this will allow for direct manipulation of the generated music’s properties even though these features were not explicit in the dataset.

7 Acknowledgement

This work is funded by Chinese scholarship Council and Singapore Ministry of Education Grant no. MOE2018-T2-2-161.

The musical score is presented in two systems. The first system contains the first bar of music, with a '1' above the staff and 'IV/IV' above the staff. The second system contains the second and third bars, with 'IV' above the first bar and '1' above the second bar. The Medieval Recorder part is in the treble clef, the Fretless Bass part is in the bass clef, and the String Ensemble part is in the grand staff (treble and bass clefs).

Fig. 7: The first bar tonal tension is changed from 1 to 6. Here the “tensile strain” is changed, and the result shows that the first bar is the subdominant of the IV chord of C major. The second bar is subdominant and goes to I in the third bar. This result increases the tension but also progresses smoothly with the surrounding of the music.

References

1. Akama, T.: A Contextual Latent Space Model: Subsequence Modulation in Melodic Sequence. In: Proceedings of the 22nd International Society for Music Information Retrieval Conference. pp. 27–34. Online (2021)
2. Bazin, T., Hadjeres, G.: NONOTO: A Model-agnostic Web Interface for Interactive Music Composition by Inpainting. arXiv:1907.10380 (2019)
3. Briot, J.P., Hadjeres, G., Pachet, F.: Deep learning techniques for music generation. Springer (2020)

4. Brunner, G., Wang, Y., Wattenhofer, R., Zhao, S.: Symbolic Music Genre Transfer with CycleGAN. In: 2018 IEEE 30th Int. Conf. on Tools with Artificial Intelligence (ICTAI). pp. 786–793. Volos, Greece (2018)
5. Chew, E.: The Spiral Array: An Algorithm for Determining Key Boundaries. In: Proc. of the Second Int. Conf., ICMAI 2002. pp. 18–31 (2002)
6. Chou, Y.H., Chen, I., Chang, C.J., Ching, J., Yang, Y.H., et al.: Midibert-piano: Large-scale pre-training for symbolic music understanding. arXiv:2107.05223 (2021)
7. Cuthbert, M.S., Ariza, C.: music21: A toolkit for computer-aided musicology and symbolic music data. In: Proceedings of the 11th International Society for Music Information Retrieval Conference. pp. 637–642. Utrecht, Netherlands (2010)
8. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv:1810.04805 (2021)
9. Ens, J., Pasquier, P.: Mmm : Exploring conditional multi-track music generation with the transformer. arXiv:2008.06048 (2020)
10. Guo, R., Herremans, D., Magnusson, T.: Midi miner - a python library for tonal tension and track classification. arXiv:1910.02049 (2019)
11. Guo, R., Simpson, I., Magnusson, T., Kiefer, C., Herremans, D.: A variational autoencoder for music generation controlled by tonal tension. arXiv preprint arXiv:2010.06230 (2020)
12. Hadjeres, G., Crestel, L.: The piano inpainting application. arXiv:2107.05944 (2021)
13. Herremans, D., Chew, E.: Tension ribbons: Quantifying and visualising tonal tension. In: 2nd Int. Conf. on Technologies for Music Notation and Representation. pp. 8–18. Cambridge, UK (2016)
14. Herremans, D., Chew, E.: Morpheus: generating structured music with constrained patterns and tension. IEEE Transactions on Affective Computing **10**(4), 510–523 (2017)
15. Hsiao, W.Y., Liu, J.Y., Yeh, Y.C., Yang, Y.H.: Compound word transformer: Learning to compose full-song music over dynamic directed hypergraphs. arXiv:2101.02402 (2021)
16. Huang, C.A., Cooijmans, T., Roberts, A., Courville, A.C., Eck, D.: Counterpoint by convolution. In: Proc. of the 18th Int. Society for Music Information Retrieval Conf. pp. 211–218. Suzhou, China (2017)
17. Huang, C.A., Vaswani, A., Uszkoreit, J., Simon, I., Hawthorne, C., Shazeer, N., Dai, A.M., Hoffman, M.D., Dinculescu, M., Eck, D.: Music transformer: Generating music with long-term structure. In: 7th Int. Conf. on Learning Representations. New Orleans, USA (2019)
18. Huang, Y.S., Yang, Y.H.: Pop music transformer: Beat-based modeling and generation of expressive pop piano compositions. In: Proc. of the 28th ACM Int. Conf. on Multimedia. pp. 1180–1188. Seattle, USA (2020)
19. Ji, S., Luo, J., Yang, X.: A comprehensive survey on deep music generation: Multi-level representations, algorithms, evaluations, and future directions. arXiv:2011.06801 (2020)
20. Louie, R., Coenen, A., Huang, C.Z., Terry, M., Cai, C.J.: Novice-ai music creation via ai-steering tools for deep generative models. In: Proc. of the 2020 CHI Conf. on Human Factors in Computing Systems. pp. 1–13. Honolulu, USA (2020)
21. Muhamed, A., Li, L., Shi, X., Yaddanapudi, S., Chi, W., Jackson, D., Suresh, R., Lipton, Z., Smola, A.J.: Transformer-gan: Symbolic music generation using a learned loss. 4th Workshop on Machine Learning for Creativity and Design at NeurIPS 2020 (2020)

22. Oore, S., Simon, I., Dieleman, S., Eck, D., Simonyan, K.: This time with feeling: Learning expressive musical performance. *Neural Computing and Applications* **32**(4), 955–967 (2020)
23. Pati, A., Lerch, A.: Is Disentanglement enough? On Latent Representations for Controllable Music Generation. In: *Proceedings of the 22nd International Society for Music Information Retrieval Conference*. pp. 517–524. Online (2021)
24. Pati, A., Lerch, A., Hadjeres, G.: Learning to traverse latent spaces for musical score inpainting. In: *Proc. of the 20th Int. Society for Music Information Retrieval Conf.* pp. 343–351. Delft, The Netherlands, (2019)
25. Raffel, C.: Learning-based methods for comparing sequences, with applications to audio-to-midi alignment and matching. Ph.D. thesis, Columbia University (2016)
26. Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J.: Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research* **21**(140), 1–67 (2020)
27. Ren, Y., He, J., Tan, X., Qin, T., Zhao, Z., Liu, T.Y.: Popmag: Pop music accompaniment generation. In: *Proc. of the 28th ACM Int. Conf. on Multimedia*. pp. 1198–1206. Seattle, USA (2020)
28. Song, K., Tan, X., Qin, T., Lu, J., Liu, T.: MASS: masked sequence to sequence pre-training for language generation. In: *Proc. of the 36th Int. Conf. on Machine Learning*. vol. 97, pp. 5926–5936. Long Beach, USA (2019)
29. Tan, H.H., Herremans, D.: Music FaderNets: Controllable Music Generation Based On High-Level Features via Low-Level Feature Modelling. In: *Proc. of the 21st Int. Society for Music Information Retrieval Conf.* pp. 109–116. Montréal, Canada (2020)
30. Tatar, K., Pasquier, P.: Musical agents: A typology and state of the art towards musical metacreation. *Journal of New Music Research* **48**, 105 – 56 (2019)
31. Vaswani, A., Shazeer, N.M., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. *arXiv:1706.03762* (2017)
32. Yang, L.C., Lerch, A.: On the evaluation of generative models in music. *Neural Computing and Applications* **32**(9), 4773–4784 (2020)
33. Zeng, M., Tan, X., Wang, R., Ju, Z., Qin, T., Liu, T.Y.: Musicbert: Symbolic music understanding with large-scale pre-training. *arXiv:2106.05630* (2021)
34. Zixun, G., Makris, D., Herremans, D.: Hierarchical recurrent neural networks for conditional melody generation with long-term structure. In: *2021 International Joint Conference on Neural Networks (IJCNN)*. pp. 1–8. IEEE (2021)