

Equipping Pretrained Unconditional Music Transformers with Instrument and Genre Controls

Weihan Xu* Julian McAuley† Shlomo Dubnov† Hao-Wen Dong†

*Duke University †University of California San Diego

weihan.xu@duke.edu jmcauley@ucsd.edu sdubnov@ucsd.edu hwdong@ucsd.edu

Abstract—The “pretraining-and-finetuning” paradigm has become a norm for training domain-specific models in natural language processing and computer vision. In this work, we aim to examine this paradigm for symbolic music generation through leveraging the largest ever symbolic music dataset sourced from the MuseScore forum. We first pretrain a large unconditional transformer model using 1.5 million songs. We then propose a simple technique to equip this pretrained unconditional music transformer model with instrument and genre controls by finetuning the model with additional control tokens. Our proposed representation offers improved high-level controllability and expressiveness against two existing representations. The experimental results show that the proposed model can successfully generate music with user-specified instruments and genre.

In a subjective listening test, the proposed model outperforms the pretrained baseline model in terms of coherence, harmony, arrangement and overall quality.

Index Terms—Music generation, music information retrieval, computer music, machine learning, deep learning

1. Introduction

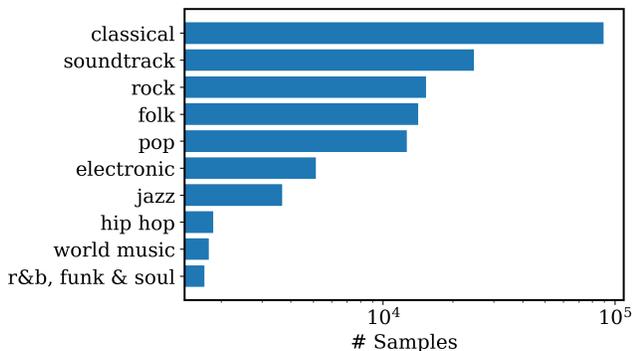
AI-driven music generation empowers people without formal music training to create music [10], [15]. On one hand, while many music generation models (e.g., [3], [16]) offer the capability for unconditional music generation, these systems offer limited controllability, e.g., users cannot specify the genre and instrumentation of the generated music. On the other hand, while some systems provide fine-grained controls for controllable music generation (e.g., [6], [8]), these controls are usually too complex for amateurs without a formal musical background. In this work, we aim to find a middle ground by equipping a pretrained unconditional music generation system with high-level controls such as genre and instrumentation to improve its controllability while keeping its accessibility to nonprofessional users. Various fine-grained control has been employed in music generation, including conditioning the model on preceding musical material, inpainting to generate music based on a partial subset of musical content, and requesting the model to fill in the gaps. For example, Music Sketchnet [11] presented a framework designed to create the missing

measures in incomplete monophonic musical compositions. This process was influenced by the surrounding context and can be enhanced by incorporating user-defined pitch and rhythm fragments if desired. PopMAG [6] generated multitrack pop music accompaniment by conditioning on chord and melody. Additionally, users can also specify artist, genre and emotion information as conditions. MuseCoco [17] extracted attributes such as genre and instruments from plain text to guide music generation. MuseCoco was trained on multiple MIDI datasets, including MMD and MetaMIDI datasets, whereas our model is trained on the larger MuseScore dataset. GetMusic [1] explored using a nonautoregressive model to generate music with any source-target track combinations. While GetMusic is trained on the MuseScore dataset, it lacks high-level controls over genre information. Multitrack Music Transformer (MMT) [12] can generate multitrack music for a set of instruments specified by the user but it lacks control over genres.

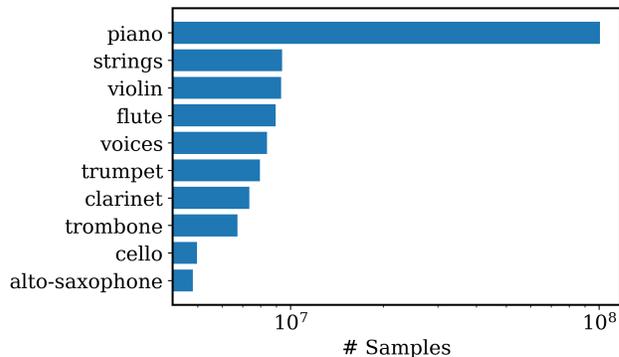
Prior work has proposed various representations for multitrack music. For example, REMI [2] and REMI+ [8], while encompassing a range of high-level musical information in a unified one-dimensional framework, do not allow for conditioning on genres and instruments. Moreover, Multitrack Music Transformer (MMT) [12] proposes a multitrack music representation with multi-dimensional inputs and outputs to tackle memory complexity issues. Although MMT offers instrumentation controls, it is not as expressive as REMI+. Therefore, we aim to design a new representation that brings together the expressiveness of REMI+ with the high-level controls available in MMT.

In this work, we develop a method to augment a pretrained unconditional music transformer with instrument and genre controls by adding new control tokens to the model during the finetuning stage. We train our proposed models on the MuseScore dataset, the largest-ever symbolic music dataset containing 1.5 million songs sourced from the MuseScore forum.¹ We evaluate the models through objective evaluation metrics and a subjective listening test. Our experiment results show that our proposed models can effectively generate music adhere to user-specified instruments and genre. The proposed model outperforms the pretrained baseline model in terms of coherence, harmony and arrangement as well as overall quality in the subjective

1. <https://musescore.com/>



(a) Most common genres



(b) Most common instruments

Figure 1. Most common genres and instruments in the MuseScore dataset

listening test. Audio samples can be found on an anonymous demo website.²

2. MuseScore Dataset

We use a collection of 1.5 million sheet music files scraped from the MuseScore forum, an online forum where users can upload sheet music for existing songs and their own compositions.¹ Table 1 summarizes the differences between several common multitrack music dataset. Moreover, we scraped the metadata from the MuseScore forum, which contains useful information such as genre, composer, rating, etc. There are 20 different genres and we use 64 different instruments in this dataset. Figure 1 shows the statistics of most common genres and instruments in the MuseScore dataset.

In this work, we consider three subsets of the MuseScore dataset that offer different levels of information for training models with different controls:

- **MuseScore-full** (1.5M): We extract information from the MuseScore dataset using MusPy package [7], producing MusPy JSON files. From these files, we further extract notes information and form MuseScore-full dataset. This dataset is used to train the unconditional pretrained model.
- **MuseScore-metadata** (821K): We notice that half of the samples in MuseScore-full do not come with a metadata file. The MuseScore-metadata subset contains only files that have a matching metadata file. This dataset is used to train the instrument conditioned model.
- **MuseScore-genre** (175K): We take a subset of MuseScore-metadata by filtering out those songs without genre information and form MuseScore-genre dataset. This dataset is used to train the genre conditioned and the genre-instrument conditioned model.

2. <https://goatlazy.github.io/MUSICAI/>

TABLE 1. LARGE DATASETS COMMONLY USED FOR MULTITRACK MUSIC GENERATION

Dataset	Publicly available	Needs scraping	Format	Songs
LMD [14]	✓	×	MIDI	175K
MetaMIDI [22]	✓	×	MIDI	436K
MMD [17]	×	-	MIDI	1.5M
MuseScore (ours)	✓	✓	XML	1.5M

3. Method

3.1. Data Representation

We represent a music piece as a one-dimensional array of integers using an event-based representation adapted from REMI+ [8] and MMT [12]. REMI+ represents notes with six consecutive tokens encoding note position, pitch, velocity, duration, instrument and time-signature information. However, it cannot provide control over genre and instruments. MMT represents a sequence of six-dimension events, with each event x_i encoded as a tuple of variables $(x^{type}, x^{beat}, x^{position}, x^{pitch}, x^{duration}, x^{instrument})$. Nevertheless, given that MMT processes the six output fields independently, it does not account for the potential interdependencies within these fields for a specific note. As a result, we adapt the REMI+ representation [8] to provide control over genre and instruments, while preserving the expressiveness offered by REMI+. Similar to MMT [12], we decompose note-on events to beat and position to reduce the size of the vocabulary and to help the model learn the rhythmic structure of music. In addition, we did not include the “tempo” and “chord” events as such information is not generally available in our dataset.

Following REMI+ [8], we use beat, position, pitch and duration events for representing musical notes. The beat event x^{beat} denotes the index of the beat that the note lies in. The value of position $x^{position}$ is determined by subtracting the product of the temporal resolution per beat and x^{beat} from the actual onset of the note. Moreover, we introduce the instrument and tag events for specifying the

Pretrained		MMT-G		MMT-I		MMT-GI	
<Start of song>	0	<Start of song>	0	<Start of song>	0	<Start of song>	0
<beat_0>	2	<Start of tags>	1311	<Start of program>	1311	<Start of tags>	1311
<poition_0>	1026	<tag_jazz>	1322	<program_piano>	1312	<tag_jazz>	1322
<instrument_piano>	1038	<Start of notes>	1310	<Start of notes>	1310	<Start of program>	1323
<pitch_69>	1160	<beat_0>	2	<beat_0>	2	<program_piano>	1312
<duration_24>	1295	<poition_0>	1026	<position_0>	1026	<Start of notes>	1310
<End of song>	1	<instrument_piano>	1038	<instrument_piano>	1038	<beat_0>	2
		<pitch_69>	1160	<pitch_69>	1219	<position_0>	1026
		<duration_24>	1295	<duration_6>	1283	<instrument_piano>	1038
		<End of song>	1	<End of song>	1	<pitch_69>	1219
						<duration_6>	1283
						<End of song>	1
Event	Index in Represent Dict.	Event	Index in Represent Dict.	Event	Index in Represent. Dict.	Event	Index in Represent. Dict.

Figure 2. An example of the baseline representation and our proposed representations—(a) pretrained, (b) genre conditioned generation used by MMT-G, (c) instrument conditioned generation used by MMT-I, and (d) genre-instrument conditioned generation used by MMT-GI. We highlight the newly added tokens in (b)–(d). The temporal resolution is 12 in this example.

instruments and tags. In addition to these data tokens, we have five special structural events:

- start-of-song event indicates the onset of a song.
- start-of-program event marks the beginning of Instruments, followed by a list of Instrument events.
- start-of-tags event marks the beginning of tags, followed by a list of tag events.
- start-of-notes event marks the beginning of notes, followed by a list of note events.
- end-of-song event indicates the end of the song, after which the model will stop predicting new tokens.

To facilitate controllability in the model, we propose to add ‘control tokens’ at the start of the data representation, as shown in Figure 2. The control tokens include tag and program. This model capitalizes on the autoregressive nature of the transformer model, enabling the integration of these tokens during the inference process, subtly enhancing the model’s controllability. Note that the start-of-notes token is used to mark the end of tag and instrument lists to prevent the model from continuing to generate instrument or tag events. The pretrained model has acquired a significant amount of features from the dataset. We show three examples in Figure 2 and highlight the new tokens that we add.

3.2. Models

We first train an unconditional model, as shown in Figure 3 then fine-tune the unconditionally pretrained model. We initialize the embeddings for tokens that already exist in the pretrained model with the corresponding pretrained weights, while assigning random weights to any newly introduced tokens. In order to equip the model with high level controllability over genre and instruments, the model inputs are one-dimensional input sequences with prefixing conditions. We propose three variants of models—MMT-I, MMT-G and MMT-GI, which offer different controls as specified below.

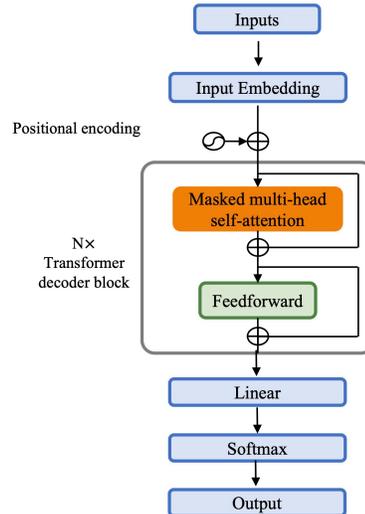


Figure 3. Illustration of the pretrained model structure

MMT-I for instrument conditioned generation.

In this scenario, this model is provided with three event types: start-of-song, start-of-program, and start-of-notes, which serve to indicate the type of information that follows them. Then the model generates note sequences based on the predefined program lists. end-of-notes marks the end of the note sequences. See Figure 2(b) for an example.

MMT-G for genre conditioned generation.

This model is provided with three event types: start-of-song, start-of-tags, and start-of-notes, which serve to denote the type of information that follows them. Then the model generates notes sequence based on the predefined tag lists. end-of-notes marks the end of the note sequences. See Figure 2 (c) for an example.

TABLE 2. OBJECTIVE EVALUATION RESULTS.
(Mean values and 95% confidence intervals are reported. A closer value to that of the ground truth is considered better.)

Model	Condition		Model size	Training samples	Pitch class entropy	Scale consistency (%)	Groove consistency (%)
	Instrument	Genre					
Pretrained	×	×	87.15M	1.35M	2.70 ± 0.08	95.92 ± 1.70	91.57 ± 2.10
MMT-I	✓	×	87.27M	739K	2.78 ± 0.08	94.36 ± 1.82	92.86 ± 1.40
Ground truth	-	-	-	-	2.48 ± 0.18	96.11 ± 1.36	92.58 ± 1.36
MMT-G	×	✓	87.18M	158K	2.88 ± 0.06	92.72 ± 2.08	92.58 ± 1.20
MMT-GI	✓	✓	87.28M	158K	2.86 ± 0.07	94.75 ± 1.90	92.48 ± 1.82
Ground truth	-	-	-	-	2.61 ± 0.13	96.10 ± 1.36	92.55 ± 1.36

MMT-GI for genre-instrument conditioned generation. This model is equipped with four event types: start-of-song, start-of-tags, start-of-program, and start-of-notes, which are used to identify the type of information that follows each event. We put the genre information ahead of the instrument list as we believe genre offers a higher-level control than instrumentation. end-of-notes marks the end of the note sequences. See Figure 2 (d) for an example.

4. Experiments & Results

4.1. Implementation Details

We use a model dimension of 768 and 12 attention heads for the transformer model. The maximum sequence length is 1024, and the maximum number of beats is 64. We use a temporal resolution of 12 time steps per quarter note for all the models. There are 87.15M, 87.27M, 87.18M and 87.28M parameters in the pretrained model, MMT-I, MMT-G and MMT-GI models, respectively. Different models are trained on different subsets based on the condition. We train the MMT-I model using the MuseScore-metadata dataset. We train the MMT-G and MMT-GI models with the MuseScore-genre dataset. We reserve 5% of each dataset for validation and 5% for test purpose for all the datasets. We validate the models every 10K steps. We pretrain the unconditional model on MuseScore-full for 1 million steps using a batch size of 32 on a NVIDIA RTX 2080 Ti GPU. The MMT-I, MMT-G and MMT-GI models are finetuned for 200K steps on a NVIDIA RTX A6000 GPU using a batch size of 16. For both the pretrained and finetuned models, we use the AdamW [21] optimizer with $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\lambda = 0.01$. For the pretrained model, the learning rate starts with 0.0005. For the finetuned model, the learning rate starts with 0.0001. For both pretrained and finetuned models, we apply a linear learning rate decay schedule so that the learning rate decays to 10% of its initial value in the first 100K steps.

4.2. Objective Evaluation Metrics

To evaluate our proposed models, we follow [7] and compute the pitch class entropy, scale consistency and groove consistency:

- *Pitch class entropy* is determined by calculating the Shannon entropy of the normalized note pitch class histogram, with the exclusion of drum tracks.
- *Scale consistency* is measured as the highest pitch-in-scale rate observed across all major and minor scales, also excluding drum tracks.
- *Groove consistency* is assessed by computing the average Hamming distance between neighboring measures.

Our primary focus lies in discerning the distribution differences in these three criteria between genuine music and generated music.

4.3. Objective Evaluation Results

We randomly generate 50 songs for each model and show in Table 2 the computed objective metrics. When evaluating the quality of generated music, we consider it better if the values are closer to those of the ground truth. When comparing MMT-G to MMT-GI, we do not observe a noticeable difference in terms of pitch class entropy, scale consistency and groove consistency. However, in terms of pitch class entropy, a significant statistical disparity is observed between MMT-G and the ground truth, and similarly, between MMT-GI with ground truth, indicating that there is still room for improvements for both models from a machine learning perspective. Moreover, in practice, we notice that the proposed models perform notably better on generating music for more common instruments (e.g., “piano”) and genres (e.g., “classical”) seen in the dataset (see Figure 1). We encourage the readers to listen to the audio samples provided in the anonymous website.²

4.4. Subjective Listening Test

Within each music generation task, we seek to understand specific attributes of the produced tunes. To this end, we conduct a listening test where each participant is instructed to assess 5 songs for every given condition. Out of the 11 participants in our study, 10 people have experience in playing instruments, with one being a professional musician.

Music quality. First, we assess the quality of the generated music in terms of coherence, harmoniousness and arrangement. The participants are instructed to answer the following questions in a Likert scale of 1 to 5.

TABLE 3. SUBJECTIVE LISTENING TEST RESULTS.
(Mean values and 95% confidence intervals are reported.)

Model	Condition		Model size	Training samples	Coherence	Harmony	Arrangement	Condition adherence	Overall
	Instrument	Genre							
Pretrained	×	×	87.15M	1.35M	3.44 ± 0.51	3.28 ± 0.46	3.16 ± 0.48	-	3.12 ± 0.45
MMT-I	✓	×	87.27M	739K	3.62 ± 0.49	3.22 ± 0.56	3.40 ± 0.48	3.15 ± 0.53	3.73 ± 0.55
MMT-G	×	✓	87.18M	158K	3.83 ± 0.43	3.73 ± 0.46	3.66 ± 0.39	3.43 ± 0.47	3.37 ± 0.55
MMT-GI	✓	✓	87.28M	158K	3.90 ± 0.48	3.80 ± 0.57	3.56 ± 0.62	3.60 ± 0.57	3.42 ± 0.55

- *Coherence*: Is it temporally coherent? Is the rhythm steady? Are there many out-of-context notes?
- *Harmoniousness*: Is it harmonious?
- *Arrangement*: Are the instruments used reasonably? Are the instruments arranged properly?

Condition Adherence. We evaluate whether the generated music aligns with the predefined tags or Instruments. This is measured by a score for how relevant the generated samples to the desired genre/instrument list. We ask participants the following question: “How relevant the generate samples are to the desired genre and instruments?”

Overall Performance. We also ask the participants how they generally feel about the generated music with the following questions “Considering all criteria above, how much do you think it is a good song that can meet the condition while maintain coherence, harmony and arrangement? Is it pleasant to listen to? How close is it to real music? Does it match the condition we set?”

4.5. Subjective Listening Test Results

We report in Table 3 the mean opinion scores (MOS) and the 95% confidence intervals assuming a Gaussian distribution. We can see our models can generate qualified music. We do not observe a significant difference between MMT-G and MMT-GI in terms of coherence, harmony and arrangement. Moreover, we see that MMT-I significantly outperforms the pretrained baseline model in terms of overall quality. Interestingly, while MMT-GI has more conditions to follow, it achieves a high score in terms of condition adherence as compared to MMT-I and MMT-G, which have only one condition to follow.

5. Discussions & Future Work

In the subjective test, we observe an intriguing pattern: music generated with genre-instrument conditioning more closely adheres to the predefined conditions. This might stem from the inherent entanglement of genre and instrumentation information.

For future work, we want to explore the following three directions. First, the metadata of the MuseScore dataset contains other valuable attributes such as composer, key signature, time signature, popularity and rating information, which can serve as powerful conditioning signals. While we

consider only genre and instrument information in this work, we would like to extend our model to include additional controls to further improve its high-level controllability.

Second, we note that the music our model produces, especially when the genre is set to classical, does not exhibit a quality comparable to that of MMT [12]. This observation is from an informal, internal side-by-side listening comparison between samples from our model and those from MMT [12]. A potential reason is the presence of noisy and low-quality entries in MuseScore. To enhance the quality of our dataset, we intend to remove low-quality songs such as practice sessions and incomplete songs as well as those without enough metadata. Furthermore, we plan to integrate other high-quality datasets such as POP909 [19] and SOD [20].

Third, we find that most of the generated music leans heavily towards certain types of music, notably classical and piano compositions. As illustrated in Figure 1, MuseScore is an imbalanced dataset that biases towards classical and piano compositions. The long-tailed nature of the dataset poses challenges to generating music of less common categories. Therefore, it is crucial to adopt a smarter sampling strategy that compensates underrepresented genres and instruments.

We plan to downsample common genres and instruments, and adapt advanced sampling and weighting method from other field [18]. Additionally, it might be helpful to include music pieces in underrepresented genres from other genre-specific datasets.

6. Conclusion

In this work, we first pretrained a large and unconditional transformer model. Then we finetuned the pretrained model with new control tokens to equip the model with genre and instrument controls. We have shown through our experiments that the proposed model can successfully produce music of the specified genre and instruments. The model also outperforms the pretrained baseline model in the subjective listening test. Our study has shown that the “pretraining-and-finetuning” paradigm can be a promising direction worth further exploration for music generation.

7. Acknowledgements

Hao-Wen thanks Taiwan Ministry of Education for supporting his PhD study. This project has received funding from

the European Research Council (ERC REACH) under the European Union’s Horizon 2020 research and innovation programme (Grant agreement #883313).

References

- [1] A. Lv, X. Tan, P. Lu, W. Ye, S. Zhang, J. Bian, and R. Yan, “GET-Music: Generating Any Music Tracks with a Unified Representation and Diffusion Framework,” *arXiv preprint arXiv:2305.10841*, 2023.
- [2] Y.-S. Huang and Y.-H. Yang, “Pop Music Transformer: Beat-based Modeling and Generation of Expressive Pop Piano Compositions,” *Proc. MM*, 2020.
- [3] 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,” *Proc. AAAI*, 2018.
- [4] R. Huang, M. Li, D. Yang, J. Shi, X. Chang, Z. Ye, Y. Wu, Z. Hong, J. Huang, J. Liu, Y. Ren, Z. Zhao, and S. Watanabe, “AudioGPT: Understanding and Generating Speech, Music, Sound, and Talking Head,” *arXiv preprint arXiv:2304.12995*, 2023.
- [5] A. Jagannathan, B. Chandrasekaran, S. Dutta, U. R. Patil, and M. Eirinaki, “Original Music Generation using Recurrent Neural Networks with Self-Attention,” *Proc. AITest*, 2022.
- [6] Y. Ren, J. He, X. Tan, T. Qin, Z. Zhao, and T.-Y. Liu, “PopMAG: Pop Music Accompaniment Generation,” *Proc. MM*, 2020.
- [7] H.-W. Dong, K. Chen, J. McAuley, and T. Berg-Kirkpatrick, “MusPy: A Toolkit for Symbolic Music Generation,” *Proc. ISMIR*, 2020.
- [8] D. von Rütte, L. Biggio, Y. Kilcher, and T. Hofmann, “FIGARO: Controllable Music Generation using Learned and Expert Features,” *Proc. ICLR*, 2023.
- [9] C.-Z. A. Huang, A. Vaswani, J. Uszkoreit, N. Shazeer, I. Simon, C. Hawthorne, A. M. Dai, M. D. Hoffman, M. Dinculescu, and D. Eck, “Music transformer: Generating music with long-term structure generating music with long-term structure,” *Proc. ICLR*, 2019.
- [10] J. P. Briot, G. Hadjeres, and F. Pachet, “Deep Learning Techniques for Music Generation - A Survey,” *arXiv preprint arXiv:1709.01620*, 2017.
- [11] K. Chen, C. Wang, T. Berg-Kirkpatrick, and S. Dubnov, “Music SketchNet: Controllable Music Generation via Factorized Representations of Pitch and Rhythm,” *Proc. ISMIR*, 2020.
- [12] H.-W. Dong, K. Chen, S. Dubnov, J. McAuley, and T. Berg-Kirkpatrick, “Multitrack Music Transformer,” *Proc. ICASSP*, 2023.
- [13] C. Donahue, H. H. Mao, and J. McAuley, “The NES Music Database: A multi-instrumental dataset with expressive performance attributes,” *Proc. ISMIR*, 2018.
- [14] Colin Raffel, “Learning-Based Methods for Comparing Sequences, with Applications to Audio-to-MIDI Alignment and Matching,” PhD thesis, Columbia University, 2016.
- [15] C.-Z. A. Huang, H. V. Kooops, E. Newton-Rex, M. Dinculescu, and C. J. Cai, “AI Song Contest: Human-AI Co-Creation in Songwriting,” *Proc. ISMIR*, 2020.
- [16] B. Yu, P. Lu, R. Wang, W. Hu, and X. Tan, “Museformer: Transformer with Fine- and Coarse-Grained Attention for Music Generation,” *Proc. NeurIPS*, 2022.
- [17] P. Lu, X. Xu, C. Kang, B. Yu, C. Xing, X. Tan, and J. Bian, “MuseCoco: Generating Symbolic Music from Text,” *arXiv preprint arXiv:2306.00110*, 2023.
- [18] M. Zhang, X. Zhao, J. Yao, C. Yuan, and W. Huang, “When Noisy Labels Meet Long Tail Dilemmas: A Representation Calibration Method,” *Proc. ICCV*, 2023.
- [19] Z. Wang, K. Chen, J. Jiang, Y. Zhang, M. Xu, S. Dai, G. Bin, and G. Xia, “POP909: A Pop-song Dataset for Music Arrangement Generation,” *Proc. ISMIR*, 2020.
- [20] L. Crestel, P. Esling, L. Heng, and S. McAdams, “A database linking piano and orchestral MIDI scores with application to automatic projective orchestration,” *Proc. ISMIR*, 2017.
- [21] Ilya Loshchilov and Frank Hutter, “Decoupled Weight Decay Regularization,” *Proc. ICLR*, 2019.
- [22] J. Ens and P. Pasquier, “Building the MetaMIDI dataset: Linking Symbolic and Audio Musical Data,” *Proc. ISMIR*, 2021.
- [23] Yin-Cheng Yeh, Wen-Yi Hsiao, Satoru Fukayama, Tetsuro Kitahara, Benjamin Genchel, Hao-Min Liu, Hao-Wen Dong, Yian Chen, Terence Leong, and Yi-Hsuan Yang, “Automatic melody harmonization with triad chords: A comparative study,” *JNMR*, vol. 50, no. 1, pp. 37–51, 2021.