VERSATILE SYMBOLIC MUSIC-FOR-MUSIC MODELING VIA **FUNCTION ALIGNMENT**

Anonymous Authors

Anonymous Affiliations

anonymous@ismir.net

41

44

45

46

47

48

49

50

52

53

54

55

56

57

59

60

61

63

64

65

66

67

ABSTRACT

1

19

Many music AI models learn a map between music con-2 tent and human-defined labels. However, many annota-3 tions, such as chords, can be naturally expressed within 4 the music modality itself, e.g., as sequences of symbolic 5 notes. This observation enables both understanding tasks 6 (e.g., chord recognition) and conditional generation tasks 7 (e.g., chord-conditioned melody generation) to be unified 8 9 under a music-for-music sequence modeling paradigm. In this work, we propose parameter-efficient solutions for a 10 variety of symbolic music-for-music tasks. The high-level 11 idea is that (1) we utilize a pretrained Language Model 12 (LM) for both the reference and the target sequence and $_{39}$ 13 (2) we link these two LMs via a lightweight adapter. Ex-14 40 periments show that our method achieves superior perfor-15 mance among different tasks such as chord recognition, 42 16 melody generation, and drum track generation. All demos, 43 17 code and model weights are publicly available¹. 18

1. INTRODUCTION

Many foundational tasks in music AI, such as music infor-20 mation retrieval (MIR) and conditional music generation, 21 have traditionally been formulated as mappings between 22 music and labels: either from music to task-specific anno-23 tations (e.g., chord recognition), or from descriptive con-24 ditions to music (e.g., chord-conditioned melody genera-25 tion). While these tasks have long been treated separately, 26 a key observation is that in many cases, the "labels" them-27 selves can also be represented in the same music modal-28 ity-for example, as note sequences. This suggests a uni-29 fying perspective: a wide range of MIR and generation 30 tasks can be reformulated as sequence-to-sequence prob-31 lems within the music domain. We refer to this formulation 32 as music-for-music modeling. 33

To achieve versatile music-for-music modeling in a 34 sample-efficient way, we apply knowledge transfer to 35 pretrained foundational Language Models (LMs) using a 36 light-parameterized adaptor. As illustrated in Fig. 1(a)-(b), 37 many existing methods such as probing [1-3] and prefix 38



Figure 1. Three types of sequence-to-sequence models by knowledge transfer from pretrained LMs. x and y are input sequences and $\hat{\mathbf{x}}$ and $\hat{\mathbf{y}}$ are predictions, possibly rightshifted due to the autoregressive targets. (a) Probing; (b) Prefix tuning; (c) Function alignment $(\mathbf{x} \rightarrow \mathbf{y})$.

tuning [4-6] transfer knowledge of foundation models to downstream tasks by adapting them to new input or output, but the knowledge resides in only one language-either the LM of source x or the LM of target y. In contrast, our method distills knowledge from both LMs via aligning them in a layer-wise manner, as shown in Fig. 1(c).

At the methodology level, our approach is inspired by function alignment [7], a recently proposed theory of mind that attributes the emergence of intelligence to the dynamic synergy among interacting agents. In our work, we contribute two concrete implementations of this idea-treating two language models (LMs) as agents and creating synergy through Parameter-Efficient Fine-Tuning (PEFT).

The first approach introduces a trainable cross-attention layer between two separately pretrained LMs. The second, more concise solution, uses a lightweight self-attentive adapter applied to concatenated input-output sequences within a single shared LM-a strategy applicable when both input and output share the same vocabulary. We show the effectiveness of both implementations using experiments on both generative and analysis tasks, including: (1) chord-conditioned melody generation, (2) melodyconditioned chord generation, (3) drum-conditioned song generation, (4) song-conditioned drum generation and (5) few-shot symbolic music analysis.

The main contribution of this paper is as follows:

- 1. We achieve versatile *music-for-music* modeling, unifying a broad range of music understanding and controllable generation tasks under a shared framework.
- 2. At the methodological level, we are the first to introduce function alignment-a recently proposed theory of mind emphasizing the synergy among agents-into the domain of music AI, offering a

¹ https://ismir2025submission25.github.io/function-alignment/

⁶⁸ (c) (i) © Anonymous Authors. Licensed under a Creative Commons Attribution 4.0 International License (CC BY 4.0). Attribution: 69 Anonymous Authors, "Versatile Symbolic Music-for-Music Modeling 70 via Function Alignment", submitted to ISMIR, 2025. 71

novel perspective on modeling music sequence-to-sequence tasks.

3. While function alignment remains at a theoretical 74 level, we present two concrete, parameter-efficient 75 implementations in the context of modern language 76 models: one via cross-attentive adapters across two 77 LMs, and another via a self-attentive adapter within 78 a shared LM. We demonstrate the effectiveness of 79 80 both approaches through theoretical analysis and empirical validation. 81

82 **2. RELATED WORKS**

83 2.1 Music Foundation Models

Since the invention of the Transformer architecture [8], 84 transformer-based language models have become the 85 mainstream of music foundation models on multiple 86 modalities, including audio [9-12], symbolic [13-21] and 87 text-based music representation [22]. In addition to au-88 127 toregressive models, masked language models [2, 23] and 89 diffusion models [24-29] and flow-based models [30] can ¹²⁸ 90 129 also be used as foundation models, but we focus on autore-91 gressive models in the literature review. 92

For symbolic music, music transformer [13] is an early ¹³⁰ 93 work to adopt the transformer architecture to music. Some 13194 follow-up works try to design a better representation of the 95 music content. For example, pop music transformer im- 132 96 poses a metrical structure in the data representation [15]. ¹³³ 97 MuPT trains transformers on their proposed synchronized 134 98 multi-track ABC notation [20]. Other works aim to in- 135 99 troduce controllability to the generative model. Musec- 136 100 oco generates the music score from text [14]. METEOR 137 101 performs melody-aware orchestral music generation with 138 102 texture control [16]. SymPAC trains symbolic generation 139 103 models from transcribed audio data with chord, section, 140 104 and instrument controls [17]. Zhang et al. improve gen-141 105 eration discriminators to better follow rhythm and melody 142 106 conditions [18]. The Theme Transformer [19] uses a short 143 107 theme condition for generation. MuseBarControl gener- 144 108 ates music with fine-grained control to the bar level [21]. 145 109

110 2.2 Parameter-Efficient Fine-Tuning

Parameter-Efficient Fine-Tuning (PEFT) methods add 148 111 lightly parameterized adapters to large pretrained models. 112 Compared to full-parameter fine-tuning, PEFT requires 113 significantly less computation and training data. Existing 114 151 methods include appending task-specific prefixes to input 115 sequences [4,31], injecting low-rank adaptation (LoRA) to 116 153 linear layers [32], and adding learnable hidden states to the 117 self-attention blocks [5, 33]. 118

PEFT has been applied to music foundation models to 119 support new tasks. Coco-Mulla [6] and MusiConGen [34] 154 120 both adapt MusicGen to follow content controls such as 155 121 chord and rhythm. Additionally, AirGen enables Mu- 156 122 sicGen to infill segments based on content controls [35]. 157 123 Instruct-MusicGen extends MusicGen for music editing 158 124 by text instructions [36]. Audio Prompt Adapter extends 159 125 AudioLDM2 for music editing following controls such as 160 126



Figure 2. The architecture of the foundation model. The left side shows the global decoder. The right side shows the encoding of a single time step $\mathbf{x}_3 = \{i_3^1, n_3^1, [eos]\}$ and the decoding of the next step $\mathbf{x}_4 = \{i_4^1, n_4^1, i_4^2, n_4^2, [eos]\}$.

genre, timbre, and melody [37]. Ou *et al.* finetunes a symbolic language model for band arrangement, piano reduction, drum arrangement, voice separation, and more [38].

3. METHODOLOGY

3.1 Base Model

For this study, we choose the base model (the pretrained symbolic LM) with two main considerations. First, we do not wish to introduce *any* control in the pretraining stage, since we want to demonstrate the controllability using PEFT. We refain from using any annotation or metadata (i.e., chord, bar or text annotations) to pretrain the base model. Second, we want to adopt a data representation that can help the model align multiple sequences in time easily. Instead of using a MIDI event-like representation [13, 15, 39] where two time-aligned sequences might have a significant length difference, we use a fixed time step (a 16th note unit) for the input sequence.

Since multiple notes can occur at the same time step, we use a hierarchical scheme to compress (decompress) the note lists on the same time step with a local encoder (decoder), as shown in Fig. 2.

3.1.1 Data Representation

146

147

Formally, we represent a score sequence $\mathbf{x} = {\mathbf{x}_1, ..., \mathbf{x}_T}$ with a fixed time step of a 16th note. Since each time step may contain multiple note onsets, each \mathbf{x}_t represents a list of N_t notes whose quantized onset time is the *t*-th 16th note (i.e., \mathbf{x}_t is a simu-note [40] at time step *t*). We define

$$\mathbf{x}_t = \{i_t^1, n_t^1, i_t^2, n_t^2, ..., i_t^{N_t}, n_t^{N_t}, [\text{eos}]\}$$
(1)

where i_t^k is the instrument ID for the k-th note. We use the MIDI program number 0...127 for pitched instruments and $i_t^k = 128$ for drums. $n_t^k = 24p_t^k + d_t^k$ is a flattened representation of the k-th note's pitch p_t^k and duration d_t^k . p_t^k denotes the MIDI pitch from 0 to 127. $d_t^k \in \{0, ..., 23\}$ is the note duration quantized into 24 possible bins, $d_t^k = j$ corrsponds to a duration of b_j sixteenth



Figure 3. The architecture of a cross-attentive function alignment adapter. The fire icon denotes trainable parameters, and the snowflake icon denotes frozen parameters.

notes where $\mathbf{b} = [1, 2, 3, 4, 6, 8, 12, 16, 24, \dots, 4096]$. [eos] 161 is a special token marking the end of the list. All notes in 162 \mathbf{x}_t are sorted primarily by i_t^k and secondarily by n_t^k . 163

3.1.2 Model Design 164

We use a Roformer [41], a popular transformer architecture 165 as the backbone model. The model architecture is shown 192 166 in Fig. 2. Since our input sequence contains nested lists, 167 193 we first encode each \mathbf{x}_t with a local Roformer encoder: 168

$$[\mathbf{h}_{t},_] = \text{LocalEncoder}([\text{cls}], \mathbf{x}_{t})$$
 (2) ¹⁹⁵

194

199

201

212

213

214

for all t = 1...T. Specifically, we prepend a [cls] token at ¹⁹⁷ 169 the beginning of \mathbf{x}_t and pass the sequence to the encoder. ¹⁹⁸ 170 \mathbf{h}_t is acquired from the output representation of the [cls] 171

token. We then use a global Roformer decoder to autore-172

gressively model the symbolic score: 173

$$\hat{\mathbf{h}}_t = \text{GlobalDecoder}(\mathbf{e}_{\text{sos}}, \mathbf{h}_{1...t-1})$$
 (3) ²⁰⁰

where \mathbf{e}_{sos} is a learnable start-of-sentence (sos) embedding. 174

Finally, a local Roformer decoder generates each note by 175

$$\hat{\mathbf{x}}_{t,j} = \text{LocalDecoder}([\text{sos}]_t, \mathbf{x}_{t,1\dots j-1})$$
(4)²⁰²₂₀₃

for all t = 1...T. Here, the embedding result of ²⁰⁴ 176 $\operatorname{Emb}([\operatorname{sos}]_t) := \hat{\mathbf{h}}_t$ passes the global state $\hat{\mathbf{h}}_t$ to the lo-²⁰⁵ 177 cal decoder. $\mathbf{x}_{t,j}$ denotes the *j*-th token of list \mathbf{x}_t (see ²⁰⁶ 178 Eqn. 1). The local decoder terminates when an end-of- 207 179 208 sentence (eos) token is generated. 180

We will use $\mathbf{x}_t = \mathrm{LM}(\mathbf{x}_{0...t-1})$ (or simply $\mathrm{LM}(\mathbf{x})$)²⁰⁹ 181 as a shorthand for the autoregressive model of sequence 182 210 x through Eqs. 2-4. Here, x_0 denotes the global start-of-183 211

sentence embedding e_{sos} . 184

3.2 Parameter-Efficient Fine-Tuning 185

Our fine-tuning strategy leverages pretrained LMs for x 215 186 and y, connected via a parameter-efficient module. We 216 187 present two variants: cross-attentive adapters for separate 217 188 LMs, and self-attentive adapters for a shared LM. We apply 218 189 both adapters to the backbone of the foundation model (the 219 190 global decoder in Eqn. 3) only. 220 191



Figure 4. The architecture of a self-attentive function alignment adapter. Crossed vertical and horizontal arrows indicate the flow of information between the corresponding query and key/value pairs, while all other connections are masked by the autoregressive self-attention mechanism. The indices 0 through 4 represent the proposed positional embeddings for the concatenated sequence.

3.2.1 Cross-attentive Function Alignment

Our first approach is to use a cross-attention layer between the hidden layers of two LMs. A similar architecture has been adopted in language processing [42] and speech processing [43]. We refer to the design of [42] and show an adapted version in Fig. 3. For the *l*-th attention layer of $LM(\mathbf{y})$, the original self-attention is defined as:

$$\mathbf{h}_{\mathbf{p}}^{l} = \text{SelfAttn}(\mathbf{W}_{q}^{l}\mathbf{z}_{y}^{l}, \mathbf{W}_{k}^{l}\mathbf{z}_{y}^{l}, \mathbf{W}_{v}^{l}\mathbf{z}_{y}^{l})$$
(5)

where \mathbf{z}_{u}^{l} denote the *l*-th layer hidden states for LM(\mathbf{y}) and \mathbf{W}^{l} denotes pretrained weights. The adapted version can be written as:

$$\mathbf{h}_{\mathrm{a}}^{l} = \mathbf{h}_{\mathrm{p}}^{l} + g \cdot \operatorname{CrossAttn}(\mathbf{U}_{q}^{l}\mathbf{z}_{y}^{l}, \mathbf{U}_{k}^{l}\mathbf{z}_{x}^{l}, \mathbf{U}_{v}^{l}\mathbf{z}_{x}^{l}) \qquad (6)$$

where g is a zero-initialized trainable gate scaler. \mathbf{U}^{l} are trainable parameters. Intuitively, this allows the query from LM(y) to attend both to itself (self-attention) and to the condition from $LM(\mathbf{x})$ (cross-attention).

Besides the trainable cross-attention module, we also append LoRAs [32] to all \mathbf{W}_{q}^{l} and \mathbf{W}_{v}^{l} of both pretrained models LM(x) and LM(y), allowing the model to learn distinctive features of sequences x and y.

3.2.2 Self-attentive Function Alignment

When x and y share the same pretrained LM, alignment becomes a special case: it can be achieved by concatenating their sequences and feeding them into a single model. The LM will first model x and predict y given x as a prefix.

This implies that prior PEFT methods [35, 38], which structure the condition and generated sequence within a single language model, can be viewed as broader forms of function alignment. We show that a simpler configuration is also effective and explain why it realizes function alignment.

When we directly concatenate two sequences $[\mathbf{x}, \mathbf{y}]$ and 266 221

feed them to the decoder self-attention layer, we can de- 267 222 compose it into the self-attention of x and y, and an ex- 268 223

tra component influencing y from x, as shown in Fig. 4. 269 224 270

Specifically, we have 225

$$\begin{aligned} [\mathbf{h}_{xa}^{l}, \mathbf{h}_{ya}^{l}] &= \text{SelfAttn}(\mathbf{W}_{q}^{l}[\mathbf{z}_{x}^{l}, \mathbf{z}_{y}^{l}], \mathbf{W}_{k}^{l}[\mathbf{z}_{x}^{l}, \mathbf{z}_{y}^{l}], \mathbf{W}_{v}^{l}[\mathbf{z}_{x}^{l}, \mathbf{z}_{y}^{l}]) & \overset{271}{_{272}} \\ &= \text{SelfAttn}([\mathbf{Q}_{x}^{l}, \mathbf{Q}_{y}^{l}], [\mathbf{K}_{x}^{l}, \mathbf{K}_{y}^{l}], [\mathbf{V}_{x}^{l}, \mathbf{V}_{y}^{l}]) & \overset{273}{_{274}} \\ & (7)_{274} \end{aligned}$$

for every layer l. In a single-head setting, we have $_{275}$ 226 SelfAttn($\mathbf{Q}, \mathbf{K}, \mathbf{V}$) := softmax($\mathbf{Q}\mathbf{K}^{\top}/\sqrt{d} + \mathbf{M}$)V, where ₂₇₆ 227 d is the dimension of the key vectors and M is the autore- ₂₇₇ 228 gressive mask. We can rewrite Eqn. 7 by 229 278

$$\mathbf{h}_{xa}^{l} = \text{SelfAttn}(\mathbf{Q}_{x}^{l}, \mathbf{K}_{x}^{l}, \mathbf{V}_{x}^{l}) \tag{8}_{279}$$

$$\mathbf{h}_{y\mathbf{a}}^{l} = \frac{\mathbf{a}}{\mathbf{a} + \mathbf{b}} \text{SelfAttn}(\mathbf{Q}_{y}^{l}, \mathbf{K}_{y}^{l}, \mathbf{V}_{y}^{l}) \\ + \frac{\mathbf{b}}{\mathbf{a} + \mathbf{b}} \text{CrossAttn}(\mathbf{Q}_{y}^{l}, \mathbf{K}_{x}^{l}, \mathbf{V}_{x}^{l})$$
(9)

280

281

282

where $\mathbf{a} = \sum_{j} \exp\left[(\mathbf{Q}_{y}^{l} \mathbf{K}_{y}^{l})_{j} / \sqrt{d} + \mathbf{M} \right]$ and $\mathbf{b} = \frac{283}{284}$ 231 284 $\sum_{j} \exp\left[(\mathbf{Q}_{y}^{l} \mathbf{K}_{x}^{l})_{j} / \sqrt{d} \right]$. Note that Eqn. 9 closely mir-²⁸⁴₂₈₅ 232 rors Eqn. 6 in form. Although the gating vectors a and b 286 233 are not explicitly parameterized, we hypothesize that this 234 287 235 design remains effective.

After concatenating x and y, we reset y's positional 288 236 embeddings to start from 0 to better preserve the pretrained ²⁸⁹ 237 behavior of LM(y). To avoid token indistinguishability $^{\rm 290}$ 238 due to overlapping positions, we also add zero-initialized, ²⁹¹ 239 learnable sentence embeddings e_x and e_y to their respec-²⁹² 240 tive positional encodings, as shwon in Fig. 4. 241 293

Similar to cross-attentive adapters, a trainable LoRA 242 294 module is also appended to the pretrained LM. 243 295

4. EXPERIMENTS

In the experiments, we first describe the hyperparame-245 ters and the pretraining scheme of our foundation model 246 299 (Sec. 4.1). We evaluate our adapters on both generative 247 300 and analysis tasks. We describe the tasks in Sec. 4.2 and 248 models in Sec. 4.3. We then show the setting for subjective 249 evaluation (Sec. 4.4) and objective evaluation (Sec. 4.5), 250 302

and analyze the results in Sec. 4.6. 251 303

4.1 Model Pretraining 252

305 We use a Roformer with a 12-layer global decoder (hidden 253

size 768, intermediate size 3072, 12 heads). The local en- $_{306}$ 254 coder and decoder are smaller 3-layer Roformers (hidden $_{307}$ 255 size 768, intermediate size 768, 8 heads). 256 308

We pretrain our foundation model on the Los Angeles 309 257 MIDI dataset [44], which contains approximately 405,000 310 258 MIDI files. As a score-based model, it relies on accurate 259 beat annotations (inferred from tempo change events) for ³¹¹ 260 correct quantization. However, many files in the pretrain- ³¹² 261 313 ing dataset contain incorrect tempo information. 262

To address this, we apply a rule-based filter. Normally, 314 263 note onsets are not uniformly distributed across odd and 315 264 even time steps. We compute the ratio of notes quantized 316 265

to odd vs. even time steps. If the ratio falls within $0.5 \pm$ 0.15 for every track, we assume it is poorly quantized and discard the song. This yields a cleaned subset of 357,279 files. During pretraining, we also apply a random pitch shift within [-5, 6] semitones for data augmentation.

We set the global sequence length to T = 384 and cap the maximal polyphony by $N_t < 16$, clipping excess notes per time step. A batch size of 48 is used for pretraining. We train the model for 2,000,000 iterations using AdamW [45] with $\beta = (0.9, 0.999)$ and weight decay 0.01. We use a OneCycleLR [46] scheduler with a maximum LR 10^{-4} and 10,000 warm-up steps. Pretraining takes around 12 days on $4 \times A100$ (40GB) GPUs.

4.2 Downstream Tasks

We evaluate the adaptor on different music generation and understanding task. Specifically, we have 3 sets of tasks:

- Melody to chord and chord to melody: we finetune the model on the Nottingham dataset [47] with a total of 1,020 songs. The model is asked to generate chords from a given melody or to generate a melody given a chord progression.
- Drum to others and others to drum: we fine-tune the model on a subset of 31,000 songs in the Los Angeles dataset with a drum track. The model is asked to generate the drum track given the full score of non-percussive instruments, or to generate other instruments given a drum track.
- Few-shot symbolic music analysis: we fine-tune the model on 93 songs in the RWC Pop dataset [48]. The model is asked to transcribe the chords and metrical structure given a symbolic pop music. We evaluate the results on symbolic chord recognition.

In each task, we perform a random 8:1:1 split for training, validation, and testing. For the drum-to-others and others-to-drum tasks, RWC Pop is used as an external test set.

4.3 Compared Models

We compare the performance of the following models, with slight hyperparameter adjustments to ensure comparable numbers of trainable parameters.

- FA-Cross: The base model is fine-tuned with a cross-attentive adapter (4 heads, hidden size 256), inserted every two layers of the global decoder. A LoRA with $r = 16, \alpha = 32$ is used on the query and value projectors of both LMs.
- FA-Self: The base model fine-tuned with a selfattentive adapter. A LoRA with $r = 64, \alpha = 128$ is used on the query and value projectors of both LMs.
- Coco-Mulla: The Cocomulla [6] adapter applied on the Roformer model. The adapter has a trainable positional encoding size of 384.

230

244

296

297

304



Figure 5. Subjective evaluation results. The error bars ³⁵³ show the 95% confidence intervals of the true mean. ³⁵⁴

- Prober: A 2-layer Multilayer Perceptron (MLP)
 prober as used in [2]. The MLP layer uses a
 weighted sum of all layers' hidden states and has a
 hidden dimension of 768.
 - Enc-Dec: A baseline trained from scratch with a 361 small Roformer encoder-decoder (3 layers, hidden 362 size 256, intermediate size 512, 4 heads for both en- 363 coder and decoder). 364
- **MelodyT5** [49]: an external baseline for the *melody* ³⁶⁵ *to chord* and *chord to melody* tasks. The model is ³⁶⁶ trained on 261K songs represented by ABC nota- ³⁶⁷ tions. We do not retrain this baseline. ³⁶⁸
- Assistant (Composers Assistant V2) [50]: an external state of the others to drum task. We do not retrain the baseline.
 370
 371
 372
 373
 374
 374
 375
 376
 377
 371
 371
 371
 372
 373
 374
 374
 375
 376
 376
 377
 371
 371
 371
 371
 372
 373
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 374
 <

336 4.4 Subjective Evaluation

321

322

323

324

For the three generative tasks (chord-to-melody, drum-to- 380 others, and others-to-drums), we conducted a subjective 381 evaluation via a user survey. We selected 8 songs from 382

	Chord to melody	Melody to chord	Drum to others	Others to drum
FA-Cross	1.4204	1.4177	2.0459	1.8619
	± 0.0986	± 0.1042	± 0.5598	± 0.5633
EA Salf	1.4116	1.4104	2.0222	1.8402
ra-sen	± 0.1165	± 0.0994	± 0.6322	± 0.5676
Coco-	1.8016	1.5996	2.2027	1.9860
Mulla	± 0.1702	± 0.1437	± 0.6495	± 0.6818
Enc-	1.6113	1.5067	2.5830	1.8765
Dec	± 0.1780	± 0.1201	± 0.9094	± 0.5352
Ground	1.3917	1.3917	2.0730	2.0730
Truth	± 0.0982	± 0.0982	± 0.7118	± 0.7118

Table 1. Test set perplexity on different downstream tasks.

the test set (2 for chord-to-melody, 4 for drum-to-others, and 2 for others-to-drums). We asked participants to rate both the generated outputs and ground truth on a 5-point scale across the following metrics:

- Musicality: Does it sound good as music?
- Adherence: Does it respect and follow the input condition?
- **Creativity**: Given the input conditions, is it creative in its musical decisions?

We received a total of 54 answers, and the results are shown in Fig. 5.

4.5 Objective Evaluation

355

378

379

For the generative tasks by fine-tuned models, we report the generated results' perplexity on the Roformer base model on the test set. Since perplexity is inaccurate on long repetitive generations [51], we only calculate the perplexity using 8-bar generative results (128 steps) conditioned on 2-bar prompts (32 steps). The results are shown in Tab. 1.

For the melody to chord task, we report two additional metrics to compare with MelodT5. We first calculate the L_1 distance between the chromagram (chroma) of the predicted chords and the ground-truth chords. We also report the CTnCTR [52] metric between the melody and the generated chords. Since the test part of the Nottingham dataset has significant overlap with MelodyT5's training set, we perform a small pitch shift (up to 2 semitones) for all test songs to another commonly used key in the Nottingham dataset (e.g., C major to D major, A major to G major etc). The results are shown in Tab. 2.

For the music analysis task, we represent both the chord and the metrical labels by MIDI notes. The chord notes are represented by block notes using String Ensemble 1 (MIDI program 48). The bass note is placed in the range C3 to B3 (MIDI pitch 36-41), and other chord notes are stacked above them. We use a drum track to represent metrical labels. We use a bass drum note (MIDI pitch 35) to represent a downbeat and a snare drum note (MIDI pitch 38) for subsidiary strong beats. An 8-note infilling by closed hi-hat note (MIDI pitch 42) is also used.

For sequence-to-sequence modeling, the model predicts both tracks from the full MIDI input, and final chord labels are derived via template matching on the average of

	Chroma↓	CTnCTR ↑
Ground Truth	0.0000 ± 0.0000	0.9675 ± 0.0324
FA-Cross	1.5690 ± 0.7087	0.9113 ± 0.0750
FA-Self	$1.2685 {\pm} 0.5024$	0.9484±0.0415
Coco-Mulla [6]	3.4613 ± 0.5854	0.6647±0.1219
Seq2Seq	3.0044 ± 0.5613	$0.8387 {\pm} 0.0749$
MelodyT5 [54]	3.0428 ± 0.8694	0.8463 ± 0.1036

 Table 2.
 Objective evaluation results on unprompted melody to chord generation on the test split of the Nottingham dataset.

Model	Root ↑	Majmin ↑	Seventh ↑
Chorder [39]	0.7244	0.6760	0.3374
HMM [55]	0.8386	0.8169	0.6930
FA-Cross	0.8203	0.8455	0.6761
FA-Self	0.8275	0.8693	0.6986
Prober	0.8231	0.8370	0.6191
Seq2Seq	0.1786	0.1500	0.0378

Table 3. Evaluation results on symbolic chord recognition.
 416

 The table shows the median result among the test split of 417
 417

 the RWC Pop dataset.
 418

⁴¹⁹ ³⁸³ 16 generations. The exception is the prober, trained as a $^{420}_{420}$ ³⁸⁴ 25-class classifier (12 major, 12 minor, 1 no-chord). We $^{421}_{421}$

evaluate using chord metrics (root, majmin, seventh) from $\frac{421}{422}$

the mir_eval package [53]. Results are shown in Table 3. $_{423}$

387 4.6 Evaluation Results

In this subsection, we analyze the results for each downtransformation task.

390 4.6.1 Few-shot Symbolic Music Analysis

429 With only 74 training songs, our adapters outperform rule-391 430 based baselines on both majmin and seventh categories. 392 By comparing function alignment models (FA) with the ⁴³¹ 393 432 prober, we see that using a pretrained LM for the target 394 433 sequence y (chord+drums) improves performance on the 395 music understanding task. 396

Between the function alignment models, the self-⁴³⁵ attentive adapters achieve better performance compared ⁴³⁶ to cross-attentive implementation. Such trend is also ob-⁴³⁷ served in other tasks.

401 4.6.2 Chord to Melody

439

438

446

447

The results in subjective evaluation (Fig. 5a) shows that the $_{440}^{400}$ our proposed adapters (FA-Self, FA-Cross) achieve com- $_{441}^{404}$ parable performance compared to Melody T5. Cocomulla $_{442}^{442}$ is not effective on the task, achieving even lower perfor- $_{443}^{440}$ mance compared to a encoder-decoder training. is also $_{444}^{440}$ demonstrated in objective evaluation results (Tab. 1). $_{445}^{440}$

408 4.6.3 Melody to Chord

Both the perplexity results (Tab. 1) and chord consistency 448
results (Tab. 2) demonstrate the effectiveness of function 449
alignment, especially self-attentive adapters against other 450
baselines. We note that MelodyT5 shows low chroma 451
consistency. MelodyT5 often fails to generate music that 452



Figure 6. Case study of an others-to-drum example on RWC-Pop-003. The top displays the non-drum condition inputs with a piano roll (structure labels are shown for reference but not used by the model). The bottom shows the drum track by (a) FA-Cross; (b) FA-Self; (c) Ground truth.

meets the constraints of the condition melody (e.g., replaced by an improvised melody or inconsistent structures). This results in a misalignment between the generated chords and the ground truth.

4.6.4 Drum to Others

414

415

424

425

427

428

FA models demonstrate strong performance in this category. Also, drum-to-others is the only task where Coco-Mulla outperforms Enc-Dec, highlighting the value of the pretrained LM for the target sequence y. However, Coco-Mulla does not utilize the knowledge stored in LM(x), leading to a worse performance compared to function alignment adapters.

4.6.5 Others to Drum

The others-to-drum task yields the interesting results: function alignment models outperform even the ground truth both subjectively (Fig. 5(c)) and objectively (Tab. 1). This is likely because RWC-Pop uses a limited drum set and regular patterns, while our training data (Los Angeles MIDI) includes diverse textures and instruments (e.g., Cuica). Function alignment models generate rich, varied drum patterns aligned with long-term structure, showing strong creativity and musicality (see Fig. 6 for an example). The baseline model Composer Assistant V2 [50] also produces less variation.

5. CONCLUSION AND FUTURE WORKS

In this paper, we address the problem of versatile musicfor-music modeling that unifies a broad range of music understanding and controllable generation tasks. Inspired by function alignment, we adopt a parameter-efficient approach by knowledge transfer from the pretrained LM of both the input and the output sequence. We introduce two implementations, the cross-attentive adapter and the self-attentive adapter. Both adapters show competitive results on analysis and generation tasks, with self-attentive adapters relatively outperforming.

There are mainly two future works. First, we need to refine the representations for different music-for-music tasks. We also plan to extend the framework to crossmodal adapters, such as text-to-music tasks. 453

6. REFERENCES

- 505 [1] C. Donahue, J. Thickstun, and P. Liang, "Melody tran-454 506 scription via generative pre-training," arXiv preprint 455 507 arXiv:2212.01884, 2022. 456 508
- [2] Y. Li, R. Yuan, G. Zhang, Y. Ma, X. Chen, H. Yin, 509 457
- C. Xiao, C. Lin, A. Ragni, E. Benetos et al., 510 458
- "Mert: Acoustic music understanding model with 459 large-scale self-supervised training," arXiv preprint 511 [14] 460 512
- arXiv:2306.00107, 2023. 461 513 [3] D. Li, Y. Ma, W. Wei, Q. Kong, Y. Wu, M. Che, 462
- F. Xia, E. Benetos, and W. Li, "Mertech: Instrument 514 [15] 463
- playing technique detection using self-supervised pre- 515 464 trained model with multi-task finetuning," in ICASSP 516
- 465 2024-2024 IEEE International Conference on Acous- 517 466
- tics, Speech and Signal Processing (ICASSP). IEEE, 518
- 467 2024, pp. 521-525. 468
- [4] X. L. Li and P. Liang, "Prefix-tuning: Optimizing 520 469 continuous prompts for generation," arXiv preprint 521 470 arXiv:2101.00190, 2021. 471
- [5] R. Zhang, J. Han, A. Zhou, X. Hu, S. Yan, P. Lu, H. Li, 523 472 P. Gao, and Y. Qiao, "Llama-adapter: Efficient fine- $\frac{1}{524}$ 473 tuning of language models with zero-init attention," 474
- arXiv preprint arXiv:2303.16199, 2023. 475
- [6] L. Lin, G. Xia, J. Jiang, and Y. Zhang, "Content-based ⁵²⁶ [18] 476 527 controls for music large language modeling," arXiv 477 528 preprint arXiv:2310.17162, 2023. 478 529
- [7] G. G. Xia, "Function alignment: A new theory for 530 479 mind and intelligence, part i: Foundations," 2025. 480
- 531 [Online]. Available: https://arxiv.org/abs/2503.21106 481 532
- [8] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, 533 482 L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, 534 483
- "Attention is all you need," Advances in neural infor-484
- mation processing systems, vol. 30, 2017. 485
- 536 [9] P. Dhariwal, H. Jun, C. Payne, J. W. Kim, A. Radford, 537 486 and I. Sutskever, "Jukebox: A generative model for 538 487 music," arXiv preprint arXiv:2005.00341, 2020. 488
- naeve, Y. Adi, and A. Défossez, "Simple and control-490 lable music generation," Advances in Neural Informa-⁵⁴¹ 491 tion Processing Systems, vol. 36, pp. 47704-47720, 542 492 2023. 493
- 494 [11] A. Agostinelli, T. I. Denk, Z. Borsos, J. Engel, 544 [22] M. Verzetti, A. Caillon, Q. Huang, A. Jansen, 545 495 A. Roberts, M. Tagliasacchi et al., "Musiclm: 546 496
- Generating music from text," arXiv preprint 547 497 arXiv:2301.11325, 2023. 498
- 499 [12] C. Zhang, Y. Ma, Q. Chen, W. Wang, S. Zhao, Z. Pan, 549 H. Wang, C. Ni, T. H. Nguyen, K. Zhou et al., "Inspire- 550 500 music: Integrating super resolution and large language 551 501 model for high-fidelity long-form music generation," 552 502 arXiv preprint arXiv:2503.00084, 2025. 553 503

- 504 [13] C.-Z. A. Huang, A. Vaswani, J. Uszkoreit, N. M. Shazeer, I. Simon, C. Hawthorne, A. M. Dai, M. D. Hoffman, M. Dinculescu, and D. Eck, "Music transformer: Generating music with longterm structure," in International Conference on Learning Representations, 2018. [Online]. Available: https://api.semanticscholar.org/CorpusID:54477714
 - 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.
 - Y.-S. Huang and Y.-H. Yang, "Pop music transformer: Beat-based modeling and generation of expressive pop piano compositions," in Proceedings of the 28th ACM international conference on multimedia, 2020, pp. 1180-1188.
- 519 [16] D.-V.-T. Le and Y.-H. Yang, "Meteor: Melody-aware texture-controllable symbolic orchestral music generation," arXiv preprint arXiv:2409.11753, 2024.
- 522 [17] H. Chen, J. B. L. Smith, J. Spijkervet, J. Wang, P. Zou, B. Li, Q. Kong, and X. Du, "Sympac: Scalable symbolic music generation with prompts and constraints," pp. 1029-1036, 2024.
 - Z. Zhang, L. Li, J. Zhang, Z. Hu, H. Wang, C. Yan, J. Yang, and Y. Qi, "Generating high-quality symbolic music using fine-grained discriminators," in International Conference on Pattern Recognition. Springer, 2025, pp. 332-344.
 - 19] Y.-J. Shih, S.-L. Wu, F. Zalkow, M. Müller, and Y.-H. Yang, "Theme transformer: Symbolic music generation with theme-conditioned transformer," IEEE Transactions on Multimedia, vol. 25, pp. 3495-3508, 2022.
- 535 [20] X. Qu, Y. Bai, Y. Ma, Z. Zhou, K. M. Lo, J. Liu, R. Yuan, L. Min, X. Liu, T. Zhang et al., "Mupt: A generative symbolic music pretrained transformer," arXiv preprint arXiv:2404.06393, 2024.
- 489 [10] J. Copet, F. Kreuk, I. Gat, T. Remez, D. Kant, G. Syn- ⁵³⁹ [21] Y. Shu, H. Xu, Z. Zhou, A. v. d. Hengel, and L. Liu, "Musebarcontrol: Enhancing fine-grained control in symbolic music generation through pre-training and counterfactual loss," arXiv preprint arXiv:2407.04331, 2024.
 - R. Yuan, H. Lin, Y. Wang, Z. Tian, S. Wu, T. Shen, G. Zhang, Y. Wu, C. Liu, Z. Zhou et al., "Chatmusician: Understanding and generating music intrinsically with llm," arXiv preprint arXiv:2402.16153, 2024.
 - 548 [23] H. F. García, P. Seetharaman, R. Kumar, and B. Pardo, "Vampnet: Music generation via masked acoustic token modeling," in Proceedings of the 24th International Society for Music Information Retrieval Conference, ISMIR 2023, Milan, Italy, November 5-9, 2023, 2023, pp. 359-366.

- 554 [24] K. Chen, Y. Wu, H. Liu, M. Nezhurina, T. Berg- 606
- Kirkpatrick, and S. Dubnov, "Musicldm: Enhanc- 607 555
- ing novelty in text-to-music generation using beat- 608 556
- synchronous mixup strategies," in ICASSP 2024-2024 609 557
- IEEE International Conference on Acoustics, Speech 558 610 [35] and Signal Processing (ICASSP). IEEE, 2024, pp. 559 611 560

612

617

632

- 1206-1210.
- 561 [25] S.-L. Wu, C. Donahue, S. Watanabe, and N. J. Bryan, 613 "Music controlnet: Multiple time-varying controls for 562
- music generation," IEEE/ACM Transactions on Audio, 563 615
- Speech, and Language Processing, vol. 32, pp. 2692-564 616
- 2703, 2024. 565
- 566 [26] M. W. Lam, Q. Tian, T. Li, Z. Yin, S. Feng, M. Tu, 618
- Y. Ji, R. Xia, M. Ma, X. Song et al., "Efficient neu- 619 567
- ral music generation," Advances in Neural Information 568
- Processing Systems, vol. 36, pp. 17450–17463, 2023. 569 621
- 570 [27] F. Schneider, O. Kamal, Z. Jin, and B. Schölkopf, 622
- 571 "Moûsai: Efficient text-to-music diffusion models," in 623
- Proceedings of the 62nd Annual Meeting of the Associ-624 572
- ation for Computational Linguistics (Volume 1: Long 625 573 626
- Papers), 2024, pp. 8050-8068. 574
- 575 [28] S. Hou, S. Liu, R. Yuan, W. Xue, Y. Shan, M. Zhao, 627 [38] and C. Zhang, "Editing music with melody and text: 628 576
- Using controlnet for diffusion transformer," in ICASSP 629 577
- 2025-2025 IEEE International Conference on Acous- 630 578
- tics, Speech and Signal Processing (ICASSP). IEEE, 579
- 2025, pp. 1-5. 580
- Z. Wang, L. Min, and G. Xia, "Whole-song hierarchi-⁶³³ 291 581 634 cal generation of symbolic music using cascaded diffu-582
- 635 sion models," in The Twelfth International Conference 583
- on Learning Representations, 2024. 584
- O. Tal, A. Ziv, I. Gat, F. Kreuk, and Y. Adi, 637 585 [30] 638 "Joint audio and symbolic conditioning for temporally 586 controlled text-to-music generation," arXiv preprint 639 587
- arXiv:2406.10970, 2024. 588
- $_{\rm 589}$ [31] X. Liu, K. Ji, Y. Fu, W. Tam, Z. Du, Z. Yang, and $^{\rm 641}$ J. Tang, "P-tuning: Prompt tuning can be comparable 590 to fine-tuning across scales and tasks," in Proceedings 591
- of the 60th Annual Meeting of the Association for Com- 644 [42] 592
- putational Linguistics (Volume 2: Short Papers), 2022, 645 593 pp. 61-68. 594 646
- 647 595 [32] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, 648
- S. Wang, L. Wang, W. Chen et al., "Lora: Low-rank 596
- adaptation of large language models." ICLR, vol. 1, 649 [43] 597 no. 2, p. 3, 2022. 598 650
- 651 599 [33] P. Gao, J. Han, R. Zhang, Z. Lin, S. Geng, A. Zhou,
- W. Zhang, P. Lu, C. He, X. Yue et al., "Llama-adapter 652 [44] 600 v2: Parameter-efficient visual instruction model," 653 601 arXiv preprint arXiv:2304.15010, 2023. 602 654
- 603 [34] Y. Lan, W. Hsiao, H. Cheng, and Y. Yang, "Musicon- 655 [45] gen: Rhythm and chord control for transformer-based 656 604 text-to-music generation," in Proceedings of the 25th 657 605

International Society for Music Information Retrieval Conference, ISMIR 2024, San Francisco, California, USA and Online, November 10-14, 2024, 2024, pp. 311-318.

- L. Lin, G. Xia, Y. Zhang, and J. Jiang, "Arrange, inpaint, and refine: Steerable long-term music audio generation and editing via content-based controls," arXiv preprint arXiv:2402.09508, 2024.
- 614 [36] Y. Zhang, Y. Ikemiya, W. Choi, N. Murata, M. A. Martínez-Ramírez, L. Lin, G. Xia, W.-H. Liao, Y. Mitsufuji, and S. Dixon, "Instruct-musicgen: Unlocking text-to-music editing for music language models via instruction tuning," arXiv preprint arXiv:2405.18386, 2024.
- F. Tsai, S. Wu, H. Kim, B. Chen, H. Cheng, and 620 [37] Y. Yang, "Audio prompt adapter: Unleashing music editing abilities for text-to-music with lightweight finetuning," in Proceedings of the 25th International Society for Music Information Retrieval Conference, IS-MIR 2024, San Francisco, California, USA and Online, November 10-14, 2024, 2024, pp. 634-641.
 - L. Ou, J. Zhao, Z. Wang, G. Xia, and Y. Wang, "Unlocking potential in pre-trained music language models for versatile multi-track music arrangement," arXiv preprint arXiv:2408.15176, 2024.
- 631 [39] W.-Y. Hsiao, J.-Y. Liu, Y.-C. Yeh, and Y.-H. Yang, "Compound word transformer: Learning to compose full-song music over dynamic directed hypergraphs," in Proceedings of the AAAI Conference on Artificial Intelligence, vol. 35, no. 1, 2021, pp. 178-186.
- Z. Wang, Y. Zhang, Y. Zhang, J. Jiang, R. Yang, 636 [40] J. Zhao, and G. Xia, "Pianotree vae: Structured representation learning for polyphonic music," arXiv preprint arXiv:2008.07118, 2020.
- 640 [41] J. Su, M. Ahmed, Y. Lu, S. Pan, W. Bo, and Y. Liu, "Roformer: Enhanced transformer with rotary position embedding," Neurocomputing, vol. 568, p. 127063, 2024.
 - R. Bansal, B. Samanta, S. Dalmia, N. Gupta, S. Vashishth, S. Ganapathy, A. Bapna, P. Jain, and P. Talukdar, "Llm augmented llms: Expanding capabilities through composition," arXiv preprint arXiv:2401.02412, 2024.
 - V. Zayats, P. Chen, M. Ferrari, and D. Padfield, "Zipper: A multi-tower decoder architecture for fusing modalities," arXiv preprint arXiv:2405.18669, 2024.
 - A. Lev, "Los angeles midi dataset: Sota kilo-scale midi dataset for mir and music ai purposes," in GitHub, 2024.
 - I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," arXiv preprint arXiv:1711.05101, 2017.

658 [46] L. N. Smith and N. Topin, "Super-convergence: Very

659 fast training of neural networks using large learning

- rates," in Artificial intelligence and machine learning
- *for multi-domain operations applications*, vol. 11006.
 SPIE, 2019, pp. 369–386.

⁶⁶³ [47] "Nottingham database," http://ifdo.ca/~seymour/
⁶⁶⁴ nottingham/nottingham.html, accessed: 2025-03-26.

- M. Goto, H. Hashiguchi, T. Nishimura, and R. Oka,
 "Rwc music database: Popular, classical and jazz music databases." in *ISMIR 2002, 3rd International Con- ference on Music Information Retrieval, Paris, France,*
- 669 October 13-17, 2002, Proceedings, vol. 2, 2002, pp.
 670 287–288.
- 671 [49] S. Wu, Y. Wang, X. Li, F. Yu, and M. Sun, "Melodyt5:
 672 A unified score-to-score transformer for symbolic
 673 music processing," *arXiv preprint arXiv:2407.02277*,
 674 2024.
- 675 [50] M. Malandro, "Composer's Assistant 2: Interactive
 676 Multi-Track MIDI Infilling with Fine-Grained User
 677 Control," in *Proc. 25th Int. Society for Music Informa-*678 *tion Retrieval Conf.*, San Francisco, CA, USA, 2024,
- 679 pp. 438–445.
- K. Wang, J. Deng, A. Sun, and X. Meng, "Perplexity from plm is unreliable for evaluating text quality," *arXiv preprint arXiv:2210.05892*, 2022.

⁶⁸³ [52] Y.-C. Yeh, W.-Y. Hsiao, S. Fukayama, T. Kita⁶⁸⁴ hara, B. Genchel, H.-M. Liu, H.-W. Dong, Y. Chen,
⁶⁸⁵ T. Leong, and Y.-H. Yang, "Automatic melody harmo⁶⁸⁶ nization with triad chords: A comparative study," *Jour-*⁶⁸⁷ nal of New Music Research, vol. 50, no. 1, pp. 37–51,

- 688 2021.
- 689 [53] C. Raffel, B. McFee, E. J. Humphrey, J. Salamon,
 690 O. Nieto, D. Liang, D. P. Ellis, and C. C. Raffel,
 691 "Mir_eval: A transparent implementation of common
 692 mir metrics." in *Proceedings of the 15th International*693 Society for Music Information Retrieval Conference,
 694 ISMIR 2014, Taipei, Taiwan, October 27-31, 2014,
 695 vol. 10, 2014, p. 2014.
- 696 [54] S. Wu, Y. Wang, X. Li, F. Yu, and M. Sun, "Melodyt5:
 A unified score-to-score transformer for symbolic
 music processing," *arXiv preprint arXiv:2407.02277*,
 2024.
- Z. Wang, D. Wang, Y. Zhang, and G. Xia, "Learning interpretable representation for controllable polyphonic
 music generation," *arXiv preprint arXiv:2008.07122*, 2020.