


# ChatMusician: Understanding and Generating Music Intrinsically with LLMs

 Multimodal Art Projection Research Community  
Skywork AI PTE. LTD.  
Hong Kong University of Science and Technology

## Abstract

While Large Language Models (LLMs) demonstrate impressive capabilities in text generation, we find that their ability has yet to be generalized to music, humanity’s creative language. We introduce **ChatMusician**<sup>1</sup>, an open-source LLM that integrates intrinsic musical abilities. It is based on continual pre-training and fine-tuning LLaMA2 on a text-compatible music representation, ABC notation, and the music is treated as a second language. ChatMusician can understand and generate music with a pure text tokenizer without any external multi-modal neural structures or tokenizers. Interestingly, endowing musical abilities does not harm language abilities, even achieving a slightly higher MMLU score. Our model is capable of composing well-structured, full-length music, conditioned on texts, chords, melodies, motifs, musical forms, etc, surpassing the GPT-4 baseline. On our meticulously curated college-level music understanding benchmark, **MusicTheoryBench**, ChatMusician surpasses LLaMA2 and GPT-3.5 in zero-shot setting by a noticeable margin. Our work reveals that LLMs can be an excellent compressor for music, but there remains significant territory to be conquered. We have open-sourced our 4B token music-language corpora **MusicPile**, the collected **MusicTheoryBench**, [code](#), [model](#) and [demo](#).

## 1 Introduction

The fusion of artificial intelligence and the arts, particularly music, has emerged as a pivotal area of research, for its profound implications on the essence of human creativity (Civit et al., 2022). Music holds a unique position due to its inherent structure and complexity, and Masataka (2009, 2007); Pino et al. (2023) suggest that language and music may have evolved from the same source.

Large Language Models (LLMs) have recently revolutionized various domains with their remark-

<sup>1</sup>See Contributions and Acknowledgments section for full author list.

able capacity for generating long sequences. Researchers have been exploring language modeling techniques for music generation (Vaswani et al., 2017; Huang et al., 2018; Payne, 2019; Lu et al., 2023; Dhariwal et al., 2020; Agostinelli et al., 2023; Copet et al., 2023; Margulis and Simchy-Gross, 2016; Dai et al., 2022; Jhamtani and Berg-Kirkpatrick, 2019). Although it seems that symbolic music can be treated in a similar way to natural language, research has shown that many distinct challenges are encountered in the realm of music. For example, even state-of-the-art models such as GPT-4 perform marginally better than random in music reasoning<sup>2</sup>. We argue that the main reason is that the intricacies of musical composition remain inadequately represented in current LLMs, including the long-term, contrapuntal context dependency and the complex connections between music notes and text descriptions.

Attempting to find solutions to these challenges, we propose ChatMusician, an open-source LLM that integrates intrinsic musical abilities, with a pipeline as shown in Figure 1. Our endeavors have focused on leveraging LLMs for symbolic music generation and understanding.

**Our contributions:** a) We introduce *ChatMusician*, a text-based LLM that unifies multiple symbolic music understanding and generation tasks, enriching their repertoire while maintaining or potentially enhancing their foundational general abilities. b) Empirical evaluations demonstrate our model’s superior musical composition capabilities, surpassing GPT-4 and established baselines in various music generation tasks, showcasing its prowess in generating coherent and structured musical pieces across diverse styles. c) We introduce the inaugural college-level symbolic music understanding

<sup>2</sup>The ability to estimate the varying harmonies, keys, rhythms, and other musical elements that are not explicitly annotated in a piece of music and are significant for music themes, progression, and styles is called **Music Reasoning**.

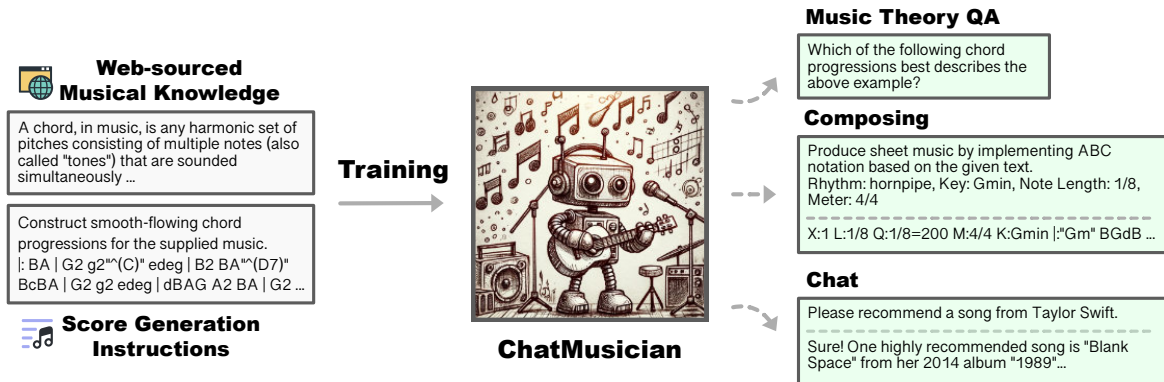


Figure 1: ChatMusician learns from web-sourced musical knowledge and handcrafted music score generation instructions, unifies music generation and music understanding, and can chat, compose, and answer college-level music theory questions.

benchmark, *MusicTheoryBench*, comprising facets of music understanding and reasoning. LLMs’ performance on this benchmark exposes their limitations, suggesting the uncharted territory of music as a domain demanding attention akin to code and mathematical reasoning. d) We open source the complete framework, including benchmark, codes, and a 4B-token music-language corpora *MusicPile*, fostering collaboration in this field.

## 2 Related Work

### 2.1 Issues in Music Generation and Understanding

The study of **music generation** is divided into acoustic (Dhariwal et al., 2020; Agostinelli et al., 2023; Copet et al., 2023) and symbolic (Sturm et al., 2015; Lu et al., 2023; Payne, 2022; Huang et al., 2018; Zhuo et al., 2023; Wu and Sun, 2022) modalities. However, the musical compositions generated by these models are still limited to a short context (for example, 30s in audio form) and are far from being fully musical and well-structured. Margulis and Simchy-Gross (2016) claims that "repetition" has a significantly positive effect on how listeners rate the "musicality" of an excerpt even if it is a random sequence. Early rule-based methods realize repetition with some pre-defined patterns which lack flexibility, whereas Dai et al. (2022) reveals that deep-learning-based works may lack repetition and music structure in the generated music.

The landscape of **music understanding** has traditionally centered on audio-focused tasks, exemplified by significant endeavors like the Music Information Retrieval Exchange (MIREX)<sup>3</sup> data chal-

<sup>3</sup>[https://www.music-ir.org/mirex/wiki/MIREX\\_HOME](https://www.music-ir.org/mirex/wiki/MIREX_HOME)

lenge and the MARBLE benchmark (?). For instance, they tackled various audio-based tasks such as genre classification, chord estimation, melody extraction, etc. In contrast, our contribution stands out with the introduction of the *MusicTheoryBench*, diverging from the conventional audio-centric focus by encompassing challenges in music verbal comprehension, advanced music theory understanding and symbolic music reasoning.

### 2.2 Music Representations

Figure 2 displays mainstream music representations with varying compression rates. Symbolic music includes formats such as MIDI, humdrum, and ABC notation (detailed in Appendix A). MIDI has been a research favorite (Lu et al., 2023; Huang and Yang, 2020a; Huang et al., 2019) with easily-accessible data due to its popularity in the music industry. However, to solve the challenge that MIDI’s lengthy sequences pose for transformer models, with their intensive training demands, sequences are typically segmented into shorter fragments, which do not represent a composition’s longer-term structure. Additionally, while MIDI encodes discrete notes and pitches, expressive performance timing in MIDI can lead to quantization errors and unstable rhythms when data is tokenized.

Therefore, we employ ABC notation, a score-oriented and plain text representation. Its high compression rate leads to shorter sequence lengths compared to MIDI and it intrinsically encodes musical repetition and structure (e.g. by the use of repeat symbols), enhancing processing efficiency using language models. It also includes detailed musical symbols denoting performance techniques and avoids quantization issues, ensuring rhythmic precision in music generation. ABC notation’s

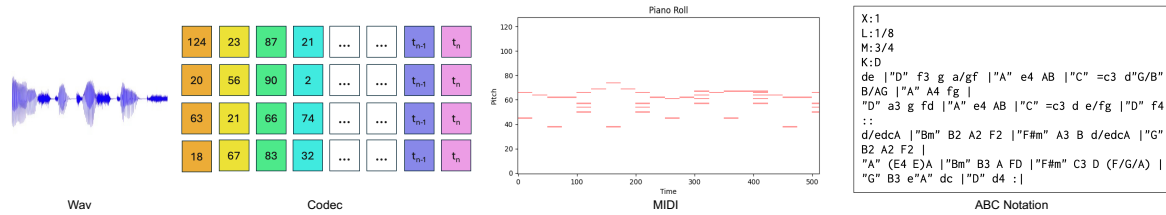


Figure 2: Commonly used music representations, including Wav, Codec, MIDI (visualized as piano roll), and ABC notation. From left to right, the compression rate gets higher.

compatibility with language models also facilitates its integration into LLM applications, allowing for advanced musical analysis and generation.

### 2.3 LLMs for Complex Problem-solving Tasks in Non-language Domain

To better understand and generate music, a model needs to handle complex sequential modeling concerning motifs, harmonies, rhythms, texture, etc., compromising between a well-organized structure and divergent creativity. Based on language sequence modeling, recent LLMs’ advancements have showcased their generalization ability in complex decision-making and problem-solving tasks across non-language domains such as mathematics, programming, and games, but have not considered music yet. MAMmoTH (Yue et al., 2023) leverage a hybrid approach of chain-of-thought and program-of-thought rationales to process and solve structured logical tasks, bridging language understanding with mathematical reasoning. CodeL-LaMA (Roziere et al., 2023), a suite of LLMs for programming tasks, exemplifies LLMs’ capabilities in applying textual instructions to generate coherent and functional code sequences. OthelloGPT (Li et al., 2022) applies a variant of GPT, using nonlinear probe representations, layerwise interventions, and latent saliency maps, to predict legal moves in the Othello game. ChessGPT (Feng et al., 2023) integrates historical chess game data and analytical insights in natural language, showcasing the fusion of policy learning with language modeling.

## 3 Method

### 3.1 Language Corpora Curation

To the best of our knowledge, there is currently no publicly available music-related natural language corpus. Fortunately, there are many large-scale corpora available from which we can curate our own.

To enable our model to interact and conversationally receive instructions, we use data from various domains. In this section, we introduce our dataset MusicPile, a first-of-its-kind pretraining dataset for injecting musical abilities into LLMs.

**General corpora.** Representative public datasets, including Pile (Gao et al., 2020), Falcon RefinedWeb (Penedo et al., 2023) and Wikipedia (Wikipedia contributors, 2023) are used. To curate a musically relevant corpus, we list a set of music-related words as a criterion to filter Pile, based on music terminology<sup>4</sup>. We only include documents with at least 10 music terms and where the terms represent at least 0.5% of the text.

**Instruction and chat data.** The instruction datasets Conover et al. (2023); Peng et al. (2023); Wang et al. (2023b) are diverse and representative enough to adapt the LLM to potential downstream usage. To enable multiple rounds of conversations, chat corpora (Wang et al., 2023a) are included.

**Music knowledge and music summary.** We crawl the metadata corresponding to 2 million music tracks from YouTube, including metadata such as song title, description, album, artist, lyrics, playlist, etc. 500k of them are extracted. We generate summaries of these metadata using GPT-4. We generate music knowledge QA pairs following Self-instruct (Wang et al., 2022). According to our topic outline in Appendix B, 255k instructions are generated, with corresponding answers generated with GPT-4.

**Math and code data.** The computational music community lacks symbolic music datasets, and we hypothesize that including math (Cobbe et al., 2021; Kenney, 2023; Yue et al., 2023; Li et al., 2023) and code (Li et al., 2023; Wang et al., 2023a)

<sup>4</sup>[https://en.m.wikipedia.org/wiki/Glossary\\_of\\_music\\_terminology](https://en.m.wikipedia.org/wiki/Glossary_of_music_terminology)

Datasets	Sourced from	Tokens	# Samples	Category	Format
Pile (Gao et al., 2020)	public dataset	0.83B	18K	general	article
Falcon-RefinedWeb (Penedo et al., 2023)	public dataset	0.80B	101K	general	article
Wikipedia (Wikipedia contributors, 2023)	public dataset	0.39B	588K	general	article
OpenChat (Wang et al., 2023a)	public dataset	62.44M	43K	general	chat
LinkSoul (LinkSoul-AI, 2023)	public dataset	0.6B	1.5M	general	chat
GPT4-Alpaca (Peng et al., 2023)	public dataset	9.77M	49K	general	chat
Dolly (Conover et al., 2023)	public dataset	3.12M	14K	general	chat
IrishMAN (Wu and Sun, 2023)	public dataset + Human-written Instructions	0.23B	868K	music score	chat
KernScores (CCARH at Stanford University, 2023)	public dataset + Human-written Instructions	2.76M	10K	music score	chat
JSB Chorales (Wu et al., 2023)	public dataset + Human-written Instructions	0.44M	349	music score	chat
synthetic music chat <sup>★</sup>	public dataset + Human-written Instructions	0.54B	50K	music score	chat
music knowledge <sup>★</sup>	Generated w/ GPT-4	0.22B	255K	music verbal	chat
music summary <sup>★</sup>	Generated w/ GPT-4	0.21B	500K	music verbal	chat
GSM8k (Cobbe et al., 2021)	public dataset	1.68M	7K	math	chat
math (Kenney, 2023)	public dataset	7.03M	37K	math	chat
MathInstruct (Yue et al., 2023)	public dataset	55.50M	188K	math	chat
Camel-Math (Li et al., 2023)	public dataset	27.76M	50K	math	chat
arxiv-math-instruct-50k (Kenney, 2023)	public dataset	9.06M	50K	math	chat
Camel-Code (Li et al., 2023)	public dataset	0.13B	366K	code	chat
OpenCoder (Wang et al., 2023a)	public dataset	36.99M	28K	code	chat
Total		4.16B	5.17M		

Table 1: Overview of MusicPile. <sup>★</sup> means synthesis from music score data and general data. <sup>★</sup> means with NEW rationales curated by us by prompting GPT-4.

may enhance the reasoning power of symbolic music. Empirically, we find this helps to improve the performance of music LLMs.

Except for the general corpora, all the other datasets were constructed as conversation forms for one or more rounds. The overall content consists of percentage of musical verbal (10.42%), code (2.43%), music score (18.43%), math (4.05%), and general (64.68%). Table 1 shows an overview of all data.

### 3.2 Music Score Corpora Curation

Although symbolic music datasets are scarce in the computational music community, we have made an effort to include music from various regions of the world. A subset of music scores, featuring regional information, has been mapped onto the world map. As depicted in Figure 3, our music scores showcase significant regional diversity. We designed a total of eight representative musical tasks on the collected corpora, including six for generating music scores and two for music understanding. The generative tasks involve generating music scores conditioned on the chord, melody, motif<sup>5</sup>, musical form<sup>6</sup>, and style. The understanding tasks involve extracting motifs and forms from the user input scores. For each task, we have created multiple instructions, which are listed in Table 2, each with

<sup>5</sup>In music, motif is a short musical idea, a salient recurring figure, musical fragment or succession of notes that has some special importance in or is characteristic of a composition

<sup>6</sup>In music, form refers to the structure of a musical composition or performance.

one example. The process of curating music instructions and algorithms is described in detail in Appendix D.

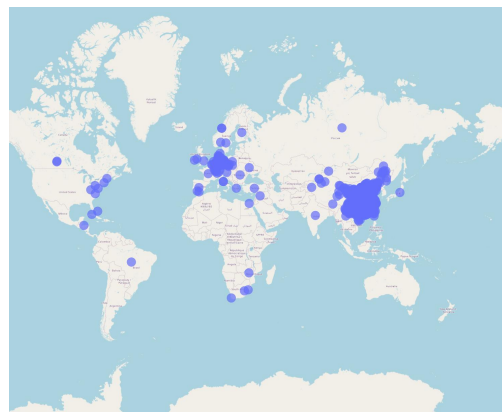


Figure 3: We included diverse music scores from around the world in MusicPile. The distribution of a portion of music scores containing regional information has been marked with blue points on the world map.

### 3.3 MusicTheoryBench

Despite significant progress in music information retrieval, the definition of advanced music understanding capabilities remains unclear in current research. In this study, to measure the advanced understanding abilities of existing LLMs in music, we first define two critical elements of music understanding: music knowledge and music reasoning. We then introduce MusicTheoryBench, a benchmark designed to assess the advanced music understanding capabilities of current LLMs.

Task Name	Type	Example Instruction
Chord Conditioned Music Generation	G	Develop a musical piece using the given chord progression. [CHORDS]
Musical Form Conditioned Music Generation	G	Craft a musical work that incorporates the given musical pattern as a central element. [MUSICAL FORMS]
Alphabetic Musical Form and Motif Conditioned Music Generation	G	Develop a musical piece employing the provided motif and an alphabet-based structure. [MUSICAL FORMS A] [MOTIF]
Terminology Musical Form and Motif conditioned Music Generation	G	Create tunes by incorporating the provided motif in the specified composition structure. [MUSICAL FORMS T] [MOTIF]
Melody Harmonization	G	Formulate chord combinations to increase the harmonic complexity of the specified musical excerpt. [MELODY]
Bach’s Style Music Generation	G	Provide a musical piece that draws inspiration from Bach’s compositions.
Motif Extraction	U	Analyze the musical work and pinpoint the consistent melodic element in every section. [MUSIC]
Musical Form Extraction	U	Investigate the attributes of this musical creation and identify its arrangement using suitable music-related terms. [MUSIC]

Table 2: Handcrafted musical tasks in MusicPile, including 6 generation tasks (Type:G) and 2 understanding tasks (Type:U), and provide an example prompt for each task. In the examples, we use tokens in square brackets to represent information other than natural language instruction ([MUSICAL FORM A] represents musical form in alphabets and [MUSICAL FORM T] represents musical form in terminology. [MOTIF], [MUSIC] and [MELODY] are represented in ABC notation. [CHORD] is represented in chord symbols.)

### Definition of Music Knowledge and Reasoning.

*Reasoning* refers to the process of making inferences based on existing knowledge and observations, usually associated with math. (Yu et al., 2023; Luo et al., 2023). Music is often likened to mathematics, where composers meticulously calculate the principles of form, harmony, scales, rhythm, tonality, and structural organization. This process ensures that the organization of notes across rhythmic and pitch domains meets established norms, yielding pleasing auditory experiences. The composition process frequently employs complex rules, including symmetry, transposition, repetition, variation, and contrast.

We define *Music Reasoning* as the capacity to infer the varying harmonies, keys, rhythms, and other musical elements that, although not explicitly annotated in a musical piece, are crucial for understanding its themes, progression, style, and internal logic. *Music knowledge*, on the other hand, is defined as the accumulated understanding of musical commonsense, e.g. notions in music theory, history, instrument characteristics, and cultural context, which informs the analytical and creative processes involved in music composition, performance, and appreciation. Examples can be found in Figure 4.

**Curation Process.** We hired a professional college music teacher to craft MusicTheoryBenchmark according to college-level textbooks and exam papers, to ensure consistency with human testing standards. The content underwent multiple rounds of discussions and reviews by a team of musicians. The team carefully selected questions and manually compiled them into JSON and ABC notation. The questions are then labeled into music knowledge and music reasoning subsets. Since the teacher is from China, half of the questions are delivered in Chinese, and later translated into English with GPT-4 Azure API and proofread by the team.

The resulting benchmark consists of 372 questions, formatted as multiple-choice questions, each with 4 options, among which only one is correct. There are 269 questions on music knowledge and 98 questions on music reasoning, along with 5 questions held out for enabling few-shot evaluation.

**Knowledge Subset.** In the music knowledge subset, the questions span Eastern and Western musical aspects. It includes 30 topics such as notes, rhythm, beats, chords, counterpoint, orchestration and instrumentation, music-related culture, history, etc (see Appendix B). Each major area undergoes targeted examination under the guidance of experts and is divided into various subcategories. For example, in the triads section, the test set specifically

**Question:** "Which of the following statements about triads is correct?",  
**Options:**  
A: "A triad can only be composed of three notes."  
B: "All triads are consonant chords."  
C: "When a triad is inverted, its properties and consonance remain consistent with its original position."  
D: "A triad is defined as a chord formed by stacking three notes in a third relationship."  
**Answer:** D

---

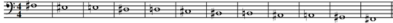
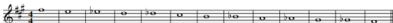

**Question:** "Which of the following is a descending natural minor scale with E as the leading tone?",  
**Options:**  
A: "L:1/4mM:4/4nK:Cn ^F,4 | ^E,4 | =E,4 | ^D,4 | =D,4 | ^C,4 | ^B,,4 | =B,,4 | ^A,,4 | =A,,4 | ^G,,4 | ^F,,4 ] %12",  
  
B: "L:1/4mM:4/4nK:A'n f4 | e4 | \_e4 | d4 | \_d4 | c4 | B4 | \_B4 | A4 | \_A4 | G4 | \_G4 | F4 ] %13",  
  
C: "L:1/4mM:4/4nK:F'n f4 | e4 | \_e4 | d4 | \_d4 | c4 | B4 | A4 | \_A4 | G4 | \_G4 | F4 ] %12",  
  
D: "None of the first three options are correct".  
**Answer:** D

Figure 4: Simple examples of (a) music knowledge and (b) music reasoning from MusicTheoryBench. Question a. mainly includes concepts that can be answered through memorizing them. Question b. requires the knowledge of *descending*, *natural minor scale* and *leading tone*, and inference based on the musical score.

examines the definition, types, and related technical details of triads. This test also features different levels of difficulty, corresponding to the high school and college levels of music major students.

**Reasoning Subset.** Most of the questions in the reasoning subset require both music knowledge and reasoning capabilities. Correctly answering these questions requires detailed analysis of the given information and multi-step logical reasoning, calculating chords, melodies, scales, rhythms, etc.

## 4 Experiments

### 4.1 Training Settings

We initialized a fp16-precision ChatMusician-Base from the LLaMA2-7B-Base weights (Touvron et al., 2023a,b), and applied a continual pre-training plus fine-tuning pipeline. The data settings will be introduced later. LoRA adapters (Hu et al., 2021) were integrated into the attention and MLP layers, with additional training on embeddings and all linear layers. The maximum sequence length was 2048. We utilized 16 80GB-A800 GPUs for one epoch pre-training and two epoch fine-tuning. DeepSpeed (Rasley et al., 2020) was employed for memory efficiency, and the AdamW optimizer was used with a 1e-4 learning rate and a 5% warm-up cosine scheduler. Gradient clipping was set at

1.0. The LoRA parameters dimension, alpha, and dropout were set to 64, 16, and 0.1, with a batch size of 8.

The model that has been pre-trained and supervised fine-tuned is called ChatMusician(CM), and the model that has only been pre-trained is called ChatMusician-Base.

### 4.2 Data Settings

During the pretraining, we combined all training data in Section 3 and performed one epoch training. To explore the effect of different data on the pre-trained model, in the supervised finetuning, we investigated different ratios of data, and empirically determined a 2:1 ratio between music knowledge vs. music summary data and music scores. We found that this ratio performed excellently in music generation as well as music understanding while guaranteeing a good MMLU performance. According to the 2:1 ratio, we first sampled 78K samples from the training set and trained for 10 epochs. Then, we maintained the ratio and utilized all available music scores data, which includes 1.1M samples, and trained for 2 epochs. The data mixture settings are summarized in Appendix G.

### 4.3 Evaluation and Baseline Systems

**Baseline Systems.** There are currently few LLMs with capabilities in symbolic music. However, observations from (Bubeck et al., 2023) suggest that the ChatGPT series possesses musical abilities. Therefore, we selected several popular LLM systems, including GPT-3.5, GPT-4, and LLaMA-2, as our baselines.

**Evaluation of General Language Abilities.** In order to evaluate general language abilities, we adopt the Massive Multitask Language Understanding (MMLU) dataset (Hendrycks et al., 2020), a pioneering benchmark designed to evaluate the knowledge acquired during pretraining of language models. To achieve a fair comparison, we evaluate our models under a 5-shot setting, which keeps the same as our selected baselines.

**Evaluation of Music Understanding Abilities.** As introduced in Section 3.3, MusicTheoryBench is a music benchmark proposed in this paper, aiming to inspect the understanding and reasoning capabilities over music knowledge for LLMs. For the MusicTheoryBench, we report the average accuracy after shuffling the option five times as the final results under a zero-shot setting.

**Evaluation of Music Generation Abilities.** Our evaluation of musicality primarily depends on human judgment. Additionally, we have developed two specific metrics: a phrase-level repetition metric and a parsing success rate metric, aimed at assessing the structure and format accuracy of the generated music. Furthermore, we introduce an average percentile score metric to gauge the models’ controllability.

## 5 Results and Discussion

### 5.1 Music Understanding

We use the proposed MusicTheoryBench to evaluate our model and the baseline systems’ music understanding abilities. We report the zero-shot performance of GPT3.5, GPT4, LLaMA2-7B-Base, ChatMusician-Base, and ChatMusician on MusicTheoryBench, as shown in Figure 5. The blue bar represents the performance on the music knowledge metric, and the red bar represents the music reasoning metric. A random baseline corresponds to a score of 25%, denoted as a dashed line.

**Music Knowledge.** According to Figure 5, all systems significantly surpassed the random baseline in the music knowledge metric. GPT-4 achieved the highest score of 58.2 on this metric. Following closely were ChatMusician-Base and ChatMusician, with scores of 40.2 and 39.5, respectively, surpassing GPT-3.5’s score of 31.2 and LLaMA2-7B-Base’s score of 33.3. This demonstrates the superiority of our method, which significantly enhanced the model’s music knowledge capability by around 7 percentage points compared to LLaMA2-7B-Base through continued training. Simultaneously, we observed the alignment tax (Zhao et al., 2023), where the fine-tuned ChatMusician scored approximately 0.7 points lower on this metric than the Base model.

**Music Reasoning.** Contrary to the performance in knowledge metrics, as shown in Figure 5, all systems exhibit subpar results in music reasoning metrics. The majority of systems do not significantly surpass the baseline in a zero-shot setting. Remarkably, even the most advanced system, GPT-4, only scored 25.6 on this metric. Interestingly, ChatMusician-Base achieved a score of 27.1 in music reasoning metrics, surpassing GPT-4. Furthermore, despite the alignment tax, ChatMusician still obtained a score of 26.3, outperforming GPT-4 in the zero-shot music reasoning metric.

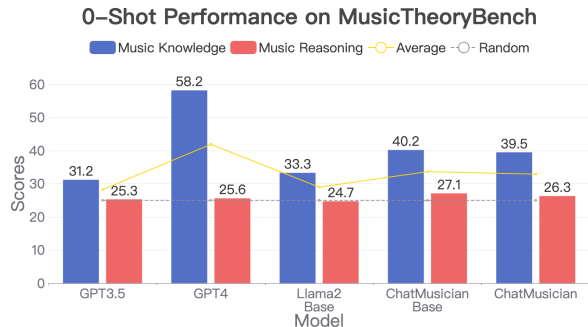


Figure 5: Zero-shot accuracy on MusicTheoryBench. We included GPT-3.5, GPT-4, LLaMA2-7B-Base, ChatMusician-Base, and ChatMusician. The blue bar represents the performance on the music knowledge metric, and the red bar represents the music reasoning metric. The dashed line corresponds to a random baseline, with a score of 25%.

Method	Mus. Knowledge	Mus. Reasoning
GPT4-0-shot	58.2	25.6
+5-shot ICL	64.1	38.0
GPT4-RolePlay	68.3	36.6
+5-shot ICL	68.8	<b>39.5</b>
GPT4-CoT	68.4	36.7
+5-shot ICL	<b>69.9</b>	34.9

Table 3: We further conducted prompt engineering on GPT-4 to check the upper limit on MusicTheoryBench. We included the techniques of chain-of-thoughts, role-play, and 5-shot in-context-learning. The highest score we achieved on music knowledge metric is 69.9, and 39.5 on music reasoning metric.

**How Far Can GPT-4 Go?** MusicTheoryBench represents the first initiative aimed at quantitatively assessing music knowledge and reasoning abilities. In pursuit of this objective, we endeavored to explore the limits of GPT-4 within our benchmark to ascertain its capabilities. GPT-4 is renowned for its robust in-context learning (ICL) and chain-of-thought (CoT) skills. Accordingly, we opted to employ prompt engineering techniques on the GPT-4 baseline to evaluate its performance on the MusicTheoryBench across various conditions, including 5-shots, CoT, and musician role-play prompts.

Table 3 displays GPT-4’s performance scores under different prompt engineering strategies. Utilizing a combination of role-play and 5-shot ICL techniques, we achieved a peak score of 39.5 in music reasoning. Meanwhile, the integration of CoT and 5-shot ICL techniques resulted in a top score of 69.9 in music knowledge. These results significantly surpass the performance of the vanilla zero-shot approach, yet they still fall short of fully

saturation of the proposed benchmark.

## 5.2 Music Generation

In this section, we demonstrate that the ABC notation format we have selected serves as an efficient means to encode and compress musical structures and repetitions in a string format. We then provide both qualitative and quantitative evidence to show that our methodology significantly enhances musicality. Moreover, it seamlessly integrates up to six conditional music generation tasks into an LLM without compromising performance.

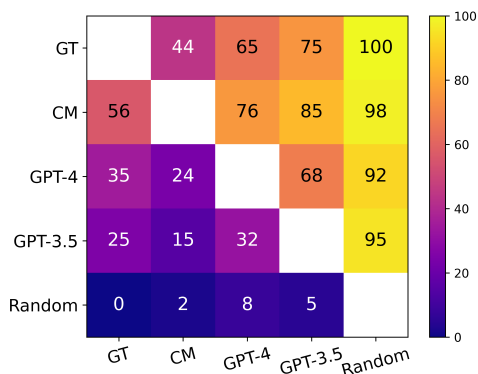


Figure 6: Results from our qualitative study where listeners judged pairs of music come from two different sources. Each row indicates the % of times listeners preferred instrumentals from that system compared to those from each system individually (N = 80). ChatMusician is denoted by CM. i.e. 76 means that listeners preferred ChatMusician over GPT-4 in 76% of cases.

### 5.2.1 Musicality

**Learning Music Repetitions.** Automatically detecting music repetition and structure remains an unresolved issue. However, by employing ABC notation, we have designed a straightforward experiment to detect phrase-level music repetition by checking the existence of repeat signs. Table 4 reports results. It appears that 76% of the generated samples from ChatMusician contain repeat signs, higher than GPT-4 and GPT-3.5. This suggests that ChatMusician is more likely to generate music with repetition and structure.

System	Repetition Det. Rate(%)
ChatMusician	<b>76.0</b>
GPT-4	70.2
GPT-3.5	32.2

Table 4: We calculate the phrase-level repetition detection rate of ABC notation strings generated by ChatMusician, GPT-4, and GPT-3.5. The higher the better.

**Human Evaluation.** A more comprehensive evaluation of music repetition and structure requires human assessment. Following (Donahue et al., 2023) and (Thickstun et al., 2023), we conduct a listening study to measure the qualitative performance of *CM* against the ground truth (*GT*) and baselines consisting of *GPT-4*, *GPT-3.5* and random note sequences (*Random*). For our study, listeners are presented with a pair of music excerpts generated from different sources, are asked to indicate which of the two pieces of music excerpts is more musical and are encouraged to pay attention to the musicality from these two aspects: how consistent the music sounds as a whole (e.g., in terms of its melodic contours, rhythmic patterns, and chord progression); and how likely the development of the music follows a clear structure (e.g. verse-chorus division, repetitions).

Results for all systems appear in Figure 6. When comparing our ChatMusician to GPT-4, listeners preferred music from our system in 76% of cases. A Wilcoxon signed-rank test of these pairwise judgments indicates that listeners preferred music from CM significantly more often than GPT-4 and GPT-3.5 ( $p = 2.7 \times 10^{-6}$  and  $p = 3.8 \times 10^{-10}$ , respectively).

A study was conducted to assess the quality of the music generated, as detailed in the Appendix E.

### 5.2.2 Controllability

**Format Correctness Evaluation.** To assess the success rate at which the output ABC notation was correctly formatted and parsed, we conducted a randomized sampling of 500 music generation prompts from the dataset. These prompts were then entered into ChatMusician, GPT-3.5, and GPT-4. To improve the parsing success rates for the GPT series, we prefixed the prompts with the directive "*Please respond in ABC notation.*". Table 5 presents the comparative success rates across the three systems. Notably, both ChatMusician and GPT-4 demonstrated success rates exceeding 90%, whereas GPT-3.5 achieved a markedly lower rate of 65.4%.

**Task-wise Metrics.** We sampled 100 prompts from each of the 5 generation tasks, and calculated average percentile scores as metrics for the 5 music generation tasks, the higher the better. Figure 7 presents the detailed score for each task of each model we have tested. See Appendix H for details.



System	Success Rate(%)
ChatMusician	<b>99.6</b>
GPT-4	94.6
GPT-3.5	65.4

Table 5: We evaluated the parsing success rates of ABC notation strings generated by ChatMusician, GPT-4, and GPT-3.5.

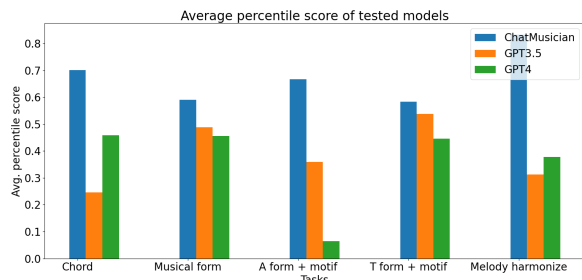


Figure 7: Here we provide the average percentile score for 5 out of 8 total musical tasks of ChatMusician, GPT-3.5 and GPT-4. Task names are abbreviations of the tasks in Table 2 (A form + motif is the abbreviation for "Alphabetic musical form and motif conditioned music generation" and T form + motif is the abbreviation for "Terminology musical form and motif conditioned music generation").

We can see that ChatMusician outperforms both GPT-3.5 and GPT-4 at all five tasks. Note that the low score of GPT-4 at task "Alphabetic musical form and motif music generation" is because most samples generated by GPT-4 of this task contain malformed ABC notation.

### 5.2.3 Compression Ratio of ABC Notation

The ABC Notation exhibits a markedly elevated compression rate in comparison to other encoding methods. For a detailed comparison, see the Appendix F. The underlying reason is straightforward. The musical score, ingeniously devised by humans, inherently encodes musical repetition. With just a repeat sign denoted as |: and :|, repeating phrases or even entire sections can be succinctly notated, corresponding to durations ranging from several seconds to minutes.

## 5.3 Language Ability

We report the MMLU score of ChatMusicians, as compared to LLaMA2-7B-Base, in Table 6. Our findings indicate that both ChatMusician and ChatMusician-Base achieve higher scores on the MMLU than the LLaMA2-7B-Base model. This suggests that incorporating our method, which infuses intrinsic music understanding and generation

System	MMLU Score(%)
ChatMusician-Base	<b>48.50</b>
ChatMusician	46.80
LLaMA2-7B-Base	46.79

Table 6: MMLU score of ChatMusicians and LLaMA2-7B-Base.

capabilities, does not compromise the general language abilities of the model. On the contrary, it appears to enhance them to a certain extent.

## 5.4 Memorization Effect of ChatMusician

We analyze the memorization abilities of ChatMusician following (Copet et al., 2023). We randomly select 500 samples from our training set and we feed the model with an instruction prompt. We compare the generated ABC notations with the ground truth. The fraction of examples where the generated and ground truth tokens are identical for the entire sequence is 0.02%. Furthermore, partial matches occur in 0.24% of the training examples, where the generated and ground truth sequences share at least 80% of their tokens.

## 6 Conclusions

In conclusion, our study introduces ChatMusician, an innovative open-source LLM capable of advanced music reasoning and composition. By leveraging a text-compatible music representation and achieving notable performance on both music and language benchmarks, ChatMusician represents a significant step forward in integrating musical creativity within language models. Our findings underscore the potential of LLMs as powerful tools for music understanding and creativity, highlighting the untapped possibilities in the fusion of music and artificial intelligence. The release of MusicPiles, MusicTheoryBench, and ChatMusician provides a valuable resource for further research in this exciting domain.

## Limitation

The current iteration of ChatMusician predominantly generates music in the style of Irish music, attributable to a significant portion of the dataset being sourced from this genre. The model exhibits hallucinations and faces limitations in supporting open-ended music generation tasks due to the lack of diversity in handcrafted music instructions.

## Ethics Statement

The model exhibits illusions, which, if employed in music education, could potentially mislead learners. Additionally, the model exhibits some memorization, raising concerns about the potential infringement of music copyrights if it inadvertently regurgitates private training data. We plan to develop a music plagiarism detection algorithm to identify instances of the memorization effect. Furthermore, we aim to implement further alignment strategies to mitigate the occurrence of illusions.

## Contributions and Acknowledgments

This paper is a tribute to our talented friend Anqiao Yang, for his friendship and valuable advice to this work. Thank Skywork AI PTE. LTD. funded the computing resources. The research was supported by Early Career Scheme (ECS-HKUST22201322), Theme-based Research Scheme (T45-205/21-N) from Hong Kong RGC, NSFC (No. 62206234), and Generative AI Research and Development Centre from InnoHK. Yizhi Li is a Ph.D. student fully funded by the Department of Computer Science, University of Manchester, UK. Yinghao Ma is a research student at the UKRI Centre for Doctoral Training in Artificial Intelligence and Music, supported by UK Research and Innovation [grant number EP/S022694/1]. Emmanouil Benetos is supported by a RAEng/Leverhulme Trust Research Fellowship [grant number LTRF2223-19-106].

### Core Contributors

Ruibin Yuan, [ryuanab@connect.ust.hk](mailto:ryuanab@connect.ust.hk)  
Hanfeng Lin, [linhanfeng@bjtu.edu.cn](mailto:linhanfeng@bjtu.edu.cn)  
Yi Wang, [ezemonyi@gmail.com](mailto:ezemonyi@gmail.com)  
Zeyue Tian, [ztianad@connect.ust.hk](mailto:ztianad@connect.ust.hk)  
Shangda Wu, [shangda@mail.ccom.edu.cn](mailto:shangda@mail.ccom.edu.cn)

### Contributors

Tianhao Shen  
Ge Zhang  
Yuhang Wu  
Cong Liu  
Ziya Zhou  
Ziyang Ma  
Liumeng Xue  
Ziyu Wang  
Qin Liu  
Tianyu Zheng  
Yizhi Li  
Yinghao Ma  
Yiming Liang

Xiaowei Chi  
Ruibo Liu  
Zili Wang  
Pengfei Li  
Jingcheng Wu  
Chenghua Lin  
Qifeng Liu  
Tao Jiang  
Wenhao Huang  
Wenhu Chen  
Emmanouil Benetos  
Jie Fu  
Gus Xia  
Roger Dannenberg

## Correspondence

Wei Xue, [weixue@ust.hk](mailto:weixue@ust.hk)  
Shiyin Kang, [shiyin.kang@kunlun-inc.com](mailto:shiyin.kang@kunlun-inc.com)  
Yike Guo, [yikeguo@ust.hk](mailto:yikeguo@ust.hk)

## References

- Andrea Agostinelli, Timo I Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, et al. 2023. Musiclm: Generating music from text. *arXiv preprint arXiv:2301.11325*.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. 2023. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*.
- CCARH at Stanford University. 2023. [A library of virtual musical scores in the humdrum \\*\\*kern data format](#).
- Miguel Civit, Javier Civit-Masot, Francisco Cuadrado, and Maria J Escalona. 2022. A systematic review of artificial intelligence-based music generation: Scope, applications, and future trends. *Expert Systems with Applications*, page 118190.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Mike Conover, Matt Hayes, Ankit Mathur, Jianwei Xie, Jun Wan, Sam Shah, Ali Ghodsi, Patrick Wendell, Matei Zaharia, and Reynold Xin. 2023. [Free dolly: Introducing the world's first truly open instruction-tuned llm](#).

- Jade Copet, Felix Kreuk, Itai Gat, Tal Remez, David Kant, Gabriel Synnaeve, Yossi Adi, and Alexandre Défossez. 2023. Simple and controllable music generation. *arXiv preprint arXiv:2306.05284*.
- Shuqi Dai, Huiran Yu, and Roger B Dannenberg. 2022. What is missing in deep music generation? a study of repetition and structure in popular music. *arXiv preprint arXiv:2209.00182*.
- Prafulla Dhariwal, Heewoo Jun, Christine Payne, Jong Wook Kim, Alec Radford, and Ilya Sutskever. 2020. Jukebox: A generative model for music. *arXiv preprint arXiv:2005.00341*.
- Chris Donahue, Antoine Caillon, Adam Roberts, Ethan Manilow, Philippe Esling, Andrea Agostinelli, Mauro Verzetti, Ian Simon, Olivier Pietquin, Neil Zeghidour, et al. 2023. Singsong: Generating musical accompaniments from singing. *arXiv preprint arXiv:2301.12662*.
- Alexandre Défossez, Jade Copet, Gabriel Synnaeve, and Yossi Adi. 2022. [High fidelity neural audio compression](#).
- Xidong Feng, Yicheng Luo, Ziyang Wang, Hongrui Tang, Mengyue Yang, Kun Shao, David Mguni, Yali Du, and Jun Wang. 2023. Chessgpt: Bridging policy learning and language modeling. *arXiv preprint arXiv:2306.09200*.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. 2020. The pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. In *International Conference on Learning Representations*.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
- Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Noam Shazeer, Ian Simon, Curtis Hawthorne, Andrew M Dai, Matthew D Hoffman, Monica Dinulescu, and Douglas Eck. 2018. Music transformer. *arXiv preprint arXiv:1809.04281*.
- Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Ian Simon, Curtis Hawthorne, Noam Shazeer, Andrew M. Dai, Matthew D. Hoffman, Monica Dinulescu, and Douglas Eck. 2019. [Music transformer: Generating music with long-term structure](#). In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.
- Yu-Siang Huang and Yi-Hsuan Yang. 2020a. [Pop music transformer: Beat-based modeling and generation of expressive pop piano compositions](#). In *MM '20: The 28th ACM International Conference on Multimedia, Virtual Event / Seattle, WA, USA, October 12-16, 2020*, pages 1180–1188. ACM.
- Yu-Siang Huang and Yi-Hsuan Yang. 2020b. [Pop music transformer: Generating music with rhythm and harmony](#). *CoRR*, abs/2002.00212.
- Harsh Jhamtani and Taylor Berg-Kirkpatrick. 2019. Modeling self-repetition in music generation using generative adversarial networks. In *Machine Learning for Music Discovery Workshop, ICML*.
- Matthew Kenney. 2023. [arxiv-math-instruct-50](#).
- Guohao Li, Hasan Abed Al Kader Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023. [Camel: Communicative agents for "mind" exploration of large scale language model society](#).
- Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2022. Emergent world representations: Exploring a sequence model trained on a synthetic task. *arXiv preprint arXiv:2210.13382*.
- LinkSoul-AI. 2023. [LinkSoul/instruction\\_merge\\_set](#). [https://huggingface.co/datasets/LinkSoul/instruction\\_merge\\_set](https://huggingface.co/datasets/LinkSoul/instruction_merge_set).
- Peiling Lu, Xin Xu, Chenfei Kang, Botao Yu, Chengyi Xing, Xu Tan, and Jiang Bian. 2023. [Musecoco: Generating symbolic music from text](#). *arXiv preprint arXiv:2306.00110*.
- Haipeng Luo, Qingfeng Sun, Can Xu, Pu Zhao, Jianguang Lou, Chongyang Tao, Xiubo Geng, Qingwei Lin, Shifeng Chen, and Dongmei Zhang. 2023. [Wizardmath: Empowering mathematical reasoning for large language models via reinforced evol-instruct](#). *arXiv preprint arXiv:2308.09583*.
- Elizabeth Hellmuth Margulis and Rhimmon Simchy-Gross. 2016. Repetition enhances the musicality of randomly generated tone sequences. *Music Perception: An Interdisciplinary Journal*, 33(4):509–514.
- Nobuo Masataka. 2007. Music, evolution and language. *Developmental science*, 10(1):35–39.
- Nobuo Masataka. 2009. The origins of language and the evolution of music: A comparative perspective. *Physics of Life Reviews*, 6(1):11–22.
- Sageev Oore, Ian Simon, Sander Dieleman, Douglas Eck, and Karen Simonyan. 2018. [This time with feeling: Learning expressive musical performance](#). *CoRR*, abs/1808.03715.
- Christine Payne. 2019. [Musenet](#). OpenAI Blog.
- Christine Payne. 2022. [Musenet](#). <https://openai.com/research/musenet>.

- Guilherme Penedo, Quentin Malartic, Daniel Hesslow, Ruxandra Cojocaru, Alessandro Cappelli, Hamza Alobeidli, Baptiste Pannier, Ebtesam Almazrouei, and Julien Launay. 2023. The refinedweb dataset for falcon llm: outperforming curated corpora with web data, and web data only. *arXiv preprint arXiv:2306.01116*.
- Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. 2023. Instruction tuning with gpt-4. *arXiv preprint arXiv:2304.03277*.
- Maria Chiara Pino, Marco Giancola, and Simonetta D’Amico. 2023. The association between music and language in children: A state-of-the-art review. *Children*, 10(5):801.
- Jeff Rasley, Samyam Rajbhandari, Olatunji Ruwase, and Yuxiong He. 2020. Deepspeed: System optimizations enable training deep learning models with over 100 billion parameters. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 3505–3506.
- Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*.
- Bob Sturm, Joao Felipe Santos, and Iryna Korshunova. 2015. Folk music style modelling by recurrent neural networks with long short term memory units. In *16th international society for music information retrieval conference*.
- John Thickstun, David Hall, Chris Donahue, and Percy Liang. 2023. Anticipatory music transformer. *arXiv preprint arXiv:2306.08620*.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajwal Bhargava, Shruti Bhosale, et al. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- Guan Wang, Sijie Cheng, Xianyuan Zhan, Xiangang Li, Sen Song, and Yang Liu. 2023a. Openchat: Advancing open-source language models with mixed-quality data. *arXiv preprint arXiv:2309.11235*.
- Yizhong Wang, Hamish Ivison, Pradeep Dasigi, Jack Hessel, Tushar Khot, Khyathi Raghavi Chandu, David Wadden, Kelsey MacMillan, Noah A Smith, Iz Beltagy, et al. 2023b. How far can camels go? exploring the state of instruction tuning on open resources. *arXiv preprint arXiv:2306.04751*.
- Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hananeh Hajishirzi. 2022. Self-instruct: Aligning language model with self generated instructions. *arXiv preprint arXiv:2212.10560*.
- Wikipedia contributors. 2023. [Wikipedia database](#).
- Shangda Wu, Xiaobing Li, and Maosong Sun. 2023. Chord-conditioned melody harmonization with controllable harmonicity. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE.
- Shangda Wu and Maosong Sun. 2022. Exploring the efficacy of pre-trained checkpoints in text-to-music generation task. *arXiv preprint arXiv:2211.11216*.
- Shangda Wu and Maosong Sun. 2023. Tunesformer: Forming tunes with control codes. *arXiv preprint arXiv:2301.02884*.
- Fei Yu, Hongbo Zhang, and Benyou Wang. 2023. Nature language reasoning, a survey. *arXiv preprint arXiv:2303.14725*.
- Xiang Yue, Xingwei Qu, Ge Zhang, Yao Fu, Wenhao Huang, Huan Sun, Yu Su, and Wenhua Chen. 2023. Mammoth: Building math generalist models through hybrid instruction tuning. *arXiv preprint arXiv:2309.05653*.
- Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. 2023. A survey of large language models. *arXiv preprint arXiv:2303.18223*.
- Le Zhuo, Zhaokai Wang, Baisen Wang, Yue Liao, Chenxi Bao, Stanley Peng, Songhao Han, Aixi Zhang, Fei Fang, and Si Liu. 2023. Video background music generation: Dataset, method and evaluation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15637–15647.

## A Introduction to ABC Notation

ABC notation is a text-based system for music notation, particularly popular for notating folk and traditional tunes. It was designed to be easily read by humans and to be simple to type without special musical fonts or software. It offers unique advantages when interfacing with deep learning models:

- **Data Efficiency:** ABC notation compactly represents musical information in a text format, making it highly efficient in terms of storage and transmission. This compactness is advantageous when training deep learning models as it minimizes data overhead.
- **Easy Preprocessing:** Being a structured text format, ABC notation can be easily tokenized and converted into numerical sequences or embeddings, a crucial step in preparing data for neural networks.
- **Scalability:** The simplicity of ABC notation allows for rapid collection and annotation of large datasets. Deep learning models, especially neural networks, benefit immensely from vast datasets, enabling better training and generalization.
- **Generative Models:** ABC notation's text-based nature makes it an excellent candidate for generative models such as LLMs, which have shown proficiency in generating coherent sequences in text-based domains.
- **Interpretability:** The outputs generated by deep learning models using ABC notation are human-readable, allowing for immediate feedback and iterative refinement. This is particularly useful in tasks such as music generation where understanding and tweaking the generated output is crucial.
- **Integration with Other Modalities:** ABC notation can be easily integrated with other data modalities in multi-modal deep learning systems, offering a comprehensive representation for music-related tasks.
- **Community Support:** The vast number of available tunes and compositions in ABC notation means that there is a rich dataset readily available. Deep learning models can leverage this to learn diverse musical structures and styles.

An ABC file consists of a series of headers followed by the music notation. The headers provide metadata about the tune, like its title, composer, rhythm, and more. The music notation section defines the melody.

Headers usually begin with a single letter followed by a colon. Some common headers include:

```
X: Reference Number
L: Default Note Length
Q: Tempo
M: Time Signature
K: Key Signature
```

The music is represented using letters, numbers, and symbols:

- Notes are denoted by the letters a-g (for notes in the octave above middle C) and A-G (for the octave below).
- Note duration is given by appending a number. For instance, A2 indicates a note twice the default length.
- Sharps, naturals, and flats are shown with ^, =, and \_ . For example, ^F is an F sharp.
- Chords are grouped using square brackets, e.g. [ceg] for a C major chord.
- Bars are marked by the | symbol.
- Tuplets, like triplets, are notated using special syntax, e.g., (3abc for a triplet of a, b, and c.
- Various decorations and ornamentations have unique symbols.

Here is a basic tune in ABC notation:

```
X:1
L:1/8
M:3/4
K:D
de | "D" f3 g a/gf | "A" e4 AB | "C" =c3 d"G/B"
B/AG | "A" A4 fg |
"D" a3 g fd | "A" e4 AB | "C" =c3 d e/fg | "D" f4
::
d/edcA | "Bm" B2 A2 F2 | "F#m" A3 B d/edcA | "G"
B2 A2 F2 |
"A" (E4 E)A | "Bm" B3 A FD | "F#m" C3 D (F/G/A) |
"G" B3 e"A" dc | "D" d4 :|
```

This represents a waltz set in D major. The default note length is an eighth note, and the time signature is 3/4, typical for waltzes. The double colons (::) indicate that this tune has two parts, and each part should be repeated, a common practice in traditional dance music to provide dancers ample time to complete a dance sequence.

## **B The examination scope of Music Theory Bench**

### 1. Pitch, Note Value, and Notation System

- Sound, Musical Tone Characteristics
- Musical Tonality System, Tone Row, and Scale Degrees
- Grouping of Notes
- Staff, Clef and Stave
- Division of Note Values
- Semitone and Whole Tone
- Temperament
- Harmonic Series

### 2. Rhythm, Beat and Note Value Combinations

- Rhythm and Beat
- Even Rhythmic Division and Irregular or Special Rhythmic Division
- Time Signature and Types of Time Signatures
- Syncopation
- Note Value Combination

### 3. Interval

- Definition and Classification of Intervals
- Degrees and Intervals
- Diatonic and Chromatic Intervals
- Single Intervals and Compound Intervals
- Inversion of Intervals
- Consonant Intervals and Dissonant Intervals
- Enharmonic Intervals

### 4. Triad

- Definition of Triad
- Types of Triads
- Inversions of Triads

### 5. Seventh chord

- Definition and Types of Seventh Chords
- Positions and Inversions of Seventh Chord
- Arrangements of Seventh Chords
- Enharmonic Chord
- Consonance of Chords
- Ninth Chord, Eleventh Chord and Thirteenth Chord

### 6. Modal Scales

- Key Name, Key Signature and Scale Degrees
- Major Scale
- Minor Scale
- Medieval Modes
- Ethnic Scales

### 7. Relationship between Keys

- Relationship between Major and Minor Keys.
- Modal Interchange
- Relative Major and Minor
- Tone Equal Temperament
- Relative Keys

### 8. Western Modes and Tonality

- Mode and Key Signature
- Natural Major and Minor Scales
- Harmonic Major and Minor Scales
- Melodic Major and Minor Scales
- Tonal Chromaticism in Modal Analysis

### 9. Ethnic Modal Scales and Tonality

- Pentatonic Scale
- Hexatonic Scale
- Heptatonic Scale

### 10. Transposition

- Western Transposition
- Ethnic Transposition

### 11. Tonal Analysis of Chord Progressions

### 12. Intervals and Chords in a Mode.

- Intervals in a Mode
- Resolution of Intervals
- Triads in a Mode
- Seventh Chords in a Mode
- Resolution of Dominant Seventh Chords

### 13. Transposition in Notation

- Interval Transposition

### 14. Chromatic Scale

- Major Chromatic Scale
- Minor Chromatic Scale
- Dynamic Markings and Terminology

- Tempo Markings and Terminology
- Inversions and Voicings
- Augmented Sixth Chords
- Neapolitan and Borrowed Chords

#### 15. Form and Structure

- Phrases, Periods and Sentences
- Binary, Ternary and Rondo Forms
- Sonata-Allegro Form
- Theme and Variations
- Fugue and Other Contrapuntal Forms

#### 16. Counterpoint

- Species Counterpoint
- The Rules of Voice-Leading
- Imitative Counterpoint (Canon, Fugue)
- Imitative Counterpoint

#### 17. Melody

- Melodic Construction and Development
- Motivic Development
- Sequences

#### 18. Twentieth-Century Techniques

- Atonality and Serialism
- Twelve-Tone Technique
- Set Theory
- Minimalism
- Microtonality

#### 19. Musical Styles and Genres

- Historical Overview from Medieval to Contemporary
- Characteristics of Different Musical Periods (e.g., Baroque, Classical, Romantic)

#### 20. Analysis Techniques

- Roman Numeral Analysis
- Schenkerian Analysis
- Graphic Analysis
- Neo-Riemannian Theory

#### 21. Orchestration and Instrumentation.

- Characteristics of Orchestral Instruments
- Basics of Writing for Different Instruments
- Full Orchestral Scoring

#### 22. Acoustics and the Science of Sound

- Overtones and Harmonics
- The Harmonic Series
- Timbre and Its Characteristics

### C Examples used in 5-shot evaluation in MusicTheoryBench

As shown in Figure 8, we present our held-out examples with prompt used in 5-shot evaluation.

### D Music Instruction Dataset Curation

We used the Irishman dataset as the basis of our music SFT data. The original dataset contains two fields: control code and ABC notation. Control code is the instruction to the generative model on the overall structure of the generated symbolic music. Here we provide a control code sample for a better explanation:

S:2 B:5 E:5 B:6 S:2

S:2 represents that there are 2 sections in this music sample, each section would be clearly marked by segmentation marks in ABC notation. B:5 represents that there are 5 bars in the first section, and B:6 represents that the second section contains 6 bars. E:5 between the two B sections represents the edit distance similarity between two music sections, in this sample: 0.5.

For the  $n^{\text{th}}$  B section, there exists  $n - 1$  number of E sections before it, in which the  $m^{\text{th}}$  E section represents the edit distance similarity between the  $m^{\text{th}}$  B section and the  $n^{\text{th}}$  B section.

#### D.0.1 Musical form analysis algorithm

For each E section before a B section, we can build a list of similarity levels for the current B section. In each of these lists, we use the following standards:

Similarity greater than or equal to 8 represents two sections that can be seen as identical sections, notated as  $s$ . The similarity between 6 and 8 represents a section that can be seen as a variation of the previous sections, notated as  $v$ . Similarity under 6 represents two different sections, notated as  $d$ ). Give the following example of control code to algorithm 2:

S:4 B:1 E:1 B:8 E:3 E:7 B:1 E:1 E:4 E:1 B:8

we would get this similarity level list  $a = [[d], [d, v], [d, d, d]]$

Read the following questions from the four options (A, B, C and D) given in each question. Choose the best option.

Which of the following chord progressions best describes the above example?

L:1/4

M:4/4

K:E

[G,B,E] [A,CE] [F,B,D] [F,A,C] |] %1

A. ii

6

/

4 - V - vi

6 - iii

B. I

6 - IV - V6

/

4 - ii

C. IV - V6

/

4 - I - ii

D. iii

6 - V - I

6

/

4 - IV

Answer: B

Which of the following best describes the seventh chord in the above example?

L:1/4

M:4/4

K:D

[FGBd]4 |] %1

A. Major seventh in third inversion

B. Dominant seventh in second inversion

C. Major/minor seventh in third inversion

D. Minor seventh in second inversion

Answer: A

Which of the following is the name of the note in the above example?

L:1/4

M:4/4

K:Cb

D,4 |] %1

A. B-flat

B. D

C. B

D. D-flat

Answer: D

The chord in the above example can be best described as which of the following?

L:1/4

M:4/4

K:F#

[EGB]4 |] %1

A. viio

B. V

C. ii

D. iv

Answer: A

[Actual question here]

Figure 8: 5-shot examples and prompt used in MusicTheoryBench benchmark.



Then we create a string to represent the alphabet musical form and put character *A* at its beginning, walk through each sub-list in the similarity level list, and mark the index of the first appeared *s* and *v*. If  $s > v$ , we will append the same alphabet at the index of *s*. If  $v > s$ , we will append the alphabet at the index of *v* with an added prime sign.

In the example above, we would get its alphabet musical form as *ABB'C*.

Using this alphabetic musical form, we can produce musical forms represented by terms. We gathered some commonly used musical form terms and put them into three categories: traditional musical forms from music theory, including *Only One Section*, *Binary*, *Ternary*, *Variational*, extended musical forms, including *American Popular*, *Verse/Chorus*, *Verse/Chorus/Bridge*, *Verse/Chorus/Verse/Bridge*, *Through Composed*, and compound musical forms, including *Compound Binary*, *Compound Ternary*.

### D.0.2 Motif extraction algorithm

The motif extraction algorithm starts by separating the sample into each section with section length information provided in the control code, then processes the token sequence *s* of each section with algorithm 1.

---

#### Algorithm 1 ABC Notation Motif Extraction

---

**Input:**  $s^{(0)} \dots s^{(n)}$   
**for**  $x = 0, 1, \dots, n$  **do**  
    **if**  $s^{(x)}$  is a bar, chord, annotation, decoration symbols, decoration characters, embellish symbols **then** Drop  $s^{(x)}$  from *s*  
    **end if**  
**end for**  
Suppose *m* is the new token sequence length. Create an empty token frequency tuple list *a*.  
**for**  $y = 0, 1, \dots, m$  **do**  
    **for**  $z = 1, 2, \dots, 8$  **do**  
        **if**  $s^{(y)}, \dots, s^{(y+z)} \notin a$  **then**  
            Add  $(s^{(y)}, \dots, s^{(y+z)}, 1)$  to *a*  
        **else** Get tuple  $(s^{(y)}, \dots, s^{(y+z)}, k)$  from *a* and update it to  $(s^{(y)}, \dots, s^{(y+z)}, k + 1)$   
        **end if**  
    **end for**  
**end for**  
**end for**  
From *a* get the tuple *b* with the largest  $b[1]$  value and  $len(b[0])$  value.  
**return**  $b[0]$  as the motif

---



---

#### Algorithm 2 Control Code Based Musical Form Analysis

---

**Input:**  $s^{(0)} \dots s^{(n)}$   
Create musical form string  $m = "A"$ , two empty list *a* and *b*, current alphabet  $n = 'A'$   
**for**  $x = 1, 2, \dots, n$  **do**  
    **if**  $s^{(x)}[0] = "B"$  **then**  
        Append *b* to *a*, create a new empty list *b*  
    **else**  
        **if**  $s^{(x)}[-1] \geq 8$  **then**  
            Append "s" to *b*  
        **else if**  $8 > s^{(x)}[-1] \geq 6$  **then**  
            Append "v" to *b*  
        **else** Append "d" to *b*  
        **end if**  
    **end if**  
**end for**  
**for** sub list *c* in list *a* **do**  
     $p_v = c.index("v")$   
     $p_s = c.index("s")$   
    **if**  $p_v > p_s$  **then**  
        Append  $m[p_v] + 'v'$  to *m*  
    **else if**  $p_v < p_s$  **then**  
        Append  $m[p_s]$  to *m*  
    **else** Append current alphabet *n* to *m* and move *n* to the next alphabet  
    **end if**  
**end for**  
**return** *m* as the musical form

---

## E Quality Study

Figure 9 presents an example of a generated musical score. The upper area displays the ABC notation, while the lower area illustrates the corresponding staff notation. Notably, repeat signs are marked in blue. Within the phrases, there is also evidence of repetition and motifs, as well as variations that echo these motifs marked in color bars.

## F Compression Ratio of ABC Notation

We sampled a set of 1,000 songs from our training corpus to evaluate the compression ratio of different music representations. As ABC notation can be converted to MIDI or rendered into WAV, we then represent these songs using widely adopted music representations (Huang and Yang, 2020b; Oore et al., 2018; Défossez et al., 2022) such as ABC notation, MIDI-like, REMI, and audio discrete audio tokens. We show that the sequence length represented by ABC strings is the shortest,

```

X:1
L:1/8
M:2/4
K:F
F/G/ |: "F" BA"C7" GG | "F" FA"C7" G2 | "F" F>G"C7" AB |
"Am" cA"C7" GF/G/ | "F" BA"C7" GG | "F" FA"C7" G2 | "F" F>G"Bb" Bd |
1"C7" cE"F" FF/G/ :|2"C7" cE"F" F z |: "F" f3 (c/d/)(d/e/) |
"Gm" (e/f/)(f/g/) g>ec | "C7" e/d/ d/c/c/B/ B/A/A/G/ | "F" GA/B/ c/d/e/f/ | f3 (c/d/)(d/e/) |
"Gm" (e/f/)(f/g/) g>ec | "C7" e/d/ d/c/c/B/ B/A/A/G/ | "F" FA/c/ f z :|

```

Figure 9: ABC notation and corresponding staff notation of a generated music. Repetition symbols are marked blue in both notations and demonstrate a clear phrase-level repetition. Red and yellow rectangles mark clear motif-level repetition in both sections. Green rectangles mark variation notes following the motif of the first section.

significantly less than other representations.

As shown in Table 7, ABC notation reaches 288.21 average tokens per song, and 5.16 average tokens per second. This is around 38% of MIDI-based representations. This suggests that using ABC notation not only facilitates compatibility with text but also reduces training costs and learning complexity.

Format	Tokenizer	Tok./Song	Tok./Sec.
ABC	LLaMA Tokenizer	288.21	5.16
MIDI	REMI(Huang and Yang, 2020b)	753.41	12.84
MIDI	MIDI-like(Oore et al., 2018)	728.60	12.42
WAV	EnCodec(Défossez et al., 2022)	12577.46	200.00

Table 7: The average number of tokens per song (Tok./Song) and tokens per second (Tok./Sec) on 1000 songs with different encoding methods. ABC notation achieves the best compression ratio.

**How Does ABC Notation Achieve Such a High Compression Ratio?** The underlying reason is straightforward. The musical score, ingeniously devised by humans, inherently encodes musical repetition. With just a repeat sign denoted as |: and :|, repeating phrases or even entire sections can be succinctly notated, corresponding to durations ranging from several seconds to minutes.

## G Settings of Data Mixture

To consider the limited computing power and explore the impact of data mixtures, we downsampled our data to a size of 52k or 78k, with different mixture proportions. This allows for the experiment to be completed in approximately one day. All the settings are in Table 8. Table 1 contains the categorization of data domains. Music verbal refers to a combination of music knowledge and music summary. Empirically, we find that setting 18 gives a balanced performance among music understanding, music generation, and language understanding abilities. Subsequently, we scaled up setting 18 to 1.1M samples and denoted it as setting 21. Setting 21 is the reported ChatMusician system in the main paper. Here we report the results of the evaluation of these models in different domains, as shown in Table 9.

## H Details of Average Percentile Score Metric

For each task, we first calculate an initial score. For chord conditioned music generation task, the initial score is calculated by taking the edit distance between the chords in the generated music and the chords in the prompt. For the musical form

ID	Setting	# Samples	Epochs
1	general + math + code	52k	10
2	music verbal	52k	10
3	music (verbal + score)	52k	10
4	general + music (verbal + score) + code + math	78k	10
5	music (verbal + score) : general = 1:2	78k	10
6	music (verbal + score) : general = 2:1	78k	10
7	general + math + music (verbal + score)	78k	10
8	general + code + music (verbal + score)	78k	10
9	general(exclude linksoul) + music (verbal + score) + code + math	78k	10
10	music verbal : general + math + code = 1:2	78k	10
11	music verbal : general + math + code = 2:1	78k	10
12	music verbal : general + math + code (en) = 1:2	78k	10
13	music verbal : general + math + code (en) = 2:1	78k	10
14	music verbal : irishman = 5:1	52k	10
15	music verbal : irishman = 1:1	52k	10
16	music verbal : synthetic music chat = 5:1	52k	10
17	music verbal : general(en) = 1:1	52k	10
18	music verbal : music score = 2:1	78k	10
19	music verbal + math : music score = 2:1	78k	10
20	music verbal + code : music score = 2:1	78k	10
21	music verbal : music score = 2:1	1.1M	2
22	music verbal : bach = 2:1	78k	10
23	music verbal : music score(half bach) = 2:1	78k	10
24	music verbal : music score(bach repeat 10) = 2:1	78k	10

Table 8: Settings of Data Mixture in Supervised Finetuning.

conditioned music generation task, the initial score is calculated by taking the difference between the set of musical forms calculated from generated music and the set of musical forms in the prompt. For the alphabetic/terminology musical form and motif-conditioned music generation task, the initial score is calculated by taking both the difference between the set of musical forms calculated from generated music and the set of musical forms in the prompt and the longest common sub-sequence of motif calculated from generated music and motif in the prompt. For the melody harmonization task, the initial score is calculated by taking the edit distance between the melody in the generated music and the melody in the prompt.

Since we have different initial score calculation methods for each task, we normalize the score to the same scale by taking the percentile of initial scores under each task. A percentile value represents that the initial value of a sample is larger than how much percentage of all the initial values in this task. For example, a percentile value of 0.6 in chord conditioned music generation task means that the initial score of the sample is larger than

60% of all the initial scores in chord conditioned music generation task. Finally, we take the average value of the percentile for each task of each model and produce the average percentile score of each task for tested models at Figure 10.

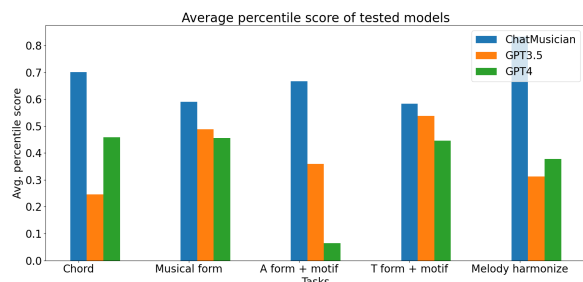


Figure 10: Here we provide the average percentile score for 5 out of 8 total musical tasks of ChatMusician, GPT-3.5 and GPT-4. Task names are abbreviations of the tasks in Table 2 (A form + motif is the abbreviation for "Alphabetic musical form and motif conditioned music generation" and T form + motif is the abbreviation for "Terminology musical form and motif conditioned music generation").

<b>ID</b>	<b>avg-music-knowledge</b>	<b>avg-music-reasoning</b>	<b>avg-music</b>	<b>gsm8k</b>	<b>openai-human-eval</b>	<b>mmlu-avg</b>	<b>abc-gen-success-rate</b>
1	40.00	27.55	33.78	28.13	10.98	48.05	37.00%
2	42.53	30.20	36.37	27.67	14.63	46.16	40.80%
3	30.63	25.10	27.87	24.64	9.76	45.14	29.40%
4	41.49	25.71	33.60	31.16	13.41	47.94	98.80%
5	32.94	24.08	28.51	14.56	7.32	45.16	17.40%
6	32.64	26.94	29.79	11.83	10.98	44.00	41.40%
7	33.68	27.14	30.41	16.76	9.76	45.24	16.40%
8	35.32	25.10	30.21	17.74	9.76	46.05	24.00%
9	41.78	26.12	33.95	29.19	13.41	47.57	96.00%
10	41.41	28.57	34.99	30.02	13.41	48.06	26.40%
11	40.37	23.47	31.92	30.55	11.59	47.67	23.00%
12	41.04	29.79	34.94	30.25	10.98	47.64	29.60%
13	41.79	27.14	34.46	28.35	7.93	48.01	8.40%
14	40.89	28.98	34.94	26.76	12.80	47.87	98.60%
15	42.08	26.73	35.22	29.75	12.80	47.70	95.00%
16	40.30	26.73	33.52	28.73	14.02	47.80	99.60%
17	43.12	25.51	34.31	23.96	9.76	46.47	27.00%
18	37.55	24.49	31.02	26.00	12.20	47.17	99.60%
19	40.59	28.37	34.48	33.28	12.20	47.87	92.20%
20	39.63	30.00	34.81	27.67	14.02	47.23	94.40%
21	39.33	24.08	31.71	30.78	12.20	46.92	99.60%
22	38.81	26.12	32.46	24.34	15.24	46.65	86.40%
23	39.48	26.12	32.80	24.94	10.98	46.74	95.40%
24	39.40	29.39	34.39	-	-	-	92.80%
LLaMa2	37.47	26.73	32.10	16.38	12.20	46.79	-
ChatMusician-Base	41.04	26.73	33.89	29.87	14.02	48.17	-
sft-ablation-base-pt	39.33	27.14	33.24	30.93	13.41	48.50	-

Table 9: Evaluation on SFT ablation models. This table presents the results of the evaluation of all supervised finetuning ablation models, including those in the MusicBench, GSM8k, OpenAI-HumanEval, MMLU and success rate of ABC Notation Generation.