## The Music Maestro or The Musically Challenged, A Massive Music Evaluation Benchmark for Large Language Models

**Anonymous ACL submission** 

#### Abstract

Benchmark plays a pivotal role in assessing the advancements of large language models (LLMs). While numerous benchmarks have been proposed to evaluate LLMs' capabilities, there is a notable absence of a dedicated benchmark for assessing their musical abilities. To address this gap, we present ZIQI-Eval, a com-800 prehensive and large-scale music benchmark specifically designed to evaluate the musicrelated capabilities of LLMs. ZIQI-Eval encompasses a wide range of questions, covering 10 major categories and 56 subcategories, resulting in over 14,000 meticulously curated 013 data entries. By leveraging ZIQI-Eval, we conduct a comprehensive evaluation over 15 LLMs to evaluate and analyze LLMs' performance in the domain of music. Results indicate that only 017 GPT-4 is capable of effectively understanding and generating music, achieving an average accuracy rate, suggesting that there is ample room for improvement in existing LLMs. With ZIQI-Eval, we aim to provide a standardized 023 and robust evaluation framework that facilitates a comprehensive assessment of LLMs' music-024 related abilities.

#### 1 Introduction

027

034

040

041

In recent years, large language models (LLMs) have made significant advancements, revolutionizing various natural language processing tasks. These models have showcased their proficiency in tasks such as accessing and reasoning about world knowledge.

Benchmark evaluation has played a crucial role in assessing and quantifying the performance of LLMs across different domains. Traditional benchmarks tailored to particular tasks such as coding (Austin et al., 2021), reading comprehension (Li et al., 2022), and mathematical reasoning (Cobbe et al., 2021), in light of the advancements made by LLMs, are increasingly regarded as inadequate for assessing their comprehensive capabilities. Consequently, there has been a surge in the emergence of more comprehensive benchmarks (Liang et al., 2022; Srivastava et al., 2022).

043

044

045

047

048

050

051

054

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

078

079

081

However, both the traditional and comprehensive benchmarks have failed to adequately address the musical accuracy of large language models. Music is an essential part of human life and culture, and assessing LLMs' understanding and generation of music presents a unique and challenging task. This oversight emphasizes the necessity for a comprehensive evaluation framework specifically designed to capture the nuances of the musical domain.

Therefore, we present ZIQI-Eval, an extensive and comprehensive music benchmark specifically crafted to assess the music-related abilities of LLMs. ZIQI-Eval comprises a diverse range of questions, systematically organized into 10 major categories and 56 subcategories. These categories cover various aspects of music, including music theory, composition, genres, instruments, and historical context. In addition, this music benchmark actively contributes to the recognition of female music composers. By incorporating valuable content from these composers, it rectifies the gender disparity prevalent in historical literature, fostering advancement and inclusivity within the realm of music scholarship. With over 14,000 carefully crafted data entries, ZIQI-Eval provides a rich and extensive resource for evaluating LLMs' comprehension and generation of music-related content.

Utilizing ZIQI-Eval, we carry out a comprehensive experiment over 17 LLMs, comprising APIbased models and open-source models, to evaluate the performance of LLMs in the realm of music. Specifically, we fed music knowledge or the first half of a musical score, along with four options, into LLMs to assess their ability to select the correct option and provide meaningful explanations. With an average accuracy rate of just 51.9%, even the top-performing model, GPT-4, falls short in demonstrating comprehensive music understand-

139

140

141

142

143

144

145

146

147

148

149

150

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

131

ing and generation capabilities. This observation not only exposes the overlooked aspect of music in LLMs but also emphasizes the significance of ZIQI-Eval in bridging this gap and tackling the inherent challenges associated with it.

084

101

102

103

104

105

106

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

In summary, our main contributions are as follows:

- We find that existing evaluations of the capabilities of large models have overlooked their musical abilities. Therefore, we propose ZIQI-Eval benchmark, a manually curated, largescale, and comprehensive benchmark for evaluating music-related capabilities. It consists of 10 major categories and 56 subcategories, encompassing over 14,000 data entries.
  - We conduct evaluations on the music comprehension and music generation capabilities of 17 LLMs and find that almost all of them struggled to understand music effectively, let alone generate it.
  - We explore the issue of bias in LLMs' music capabilities, focusing on gender bias, racial bias, and region bias. Our research reveals that....

#### 2 Related Work

Music Comprehension Inspired by the field of natural language processing (NLP), previous studies represented music as embedding sequences for music understanding. Chuan et al. (2020) and Liang et al. (2020) partition music pieces into distinct, non-overlapping segments of fixed duration, and train embeddings for each segment.

Later, with the development of large language models (LLMs), recent research has utilized the modeling capabilities of these models to further enhance the understanding of music. MidiB-ERT (Chou et al., 2021) and MusicBERT (Zeng et al., 2021) both utilize pre-trained BERT to tackle symbolic-domain discriminative music understanding tasks. MusicBERT further designs OctupleMIDI encoding and bar-level masking strategy to enhance pre-training with symbolic music data. Gardner et al. (2023) extracts music-related information from an open-source music dataset and uses instruction-tuning to instruct their proposed model LLark to do music understanding, music captioning, and music reasoning. NG-Midiformer (Tian et al., 2023) first processes music pieces into sequences, followed by leveraging N-gram encoder to understand symbolic music.

**Music Generation** Before the proliferation of LLMs, there are some other traditional methods proposed for music generation, mainly falling into three categories: neural networks, neural audio codecs, and diffusion models.

Engel et al. (2019), Marafioti et al. (2020), Greshler et al. (2021), Yu et al. (2021), Caillon and Esling (2021) employ neural network architectures such as CNNs, RNNs, or GANs to achieve music generation. A neural audio codec typically contains an encoder and a decoder. Valenti et al. (2020) follows the typical structure. Petermann et al. (2021) additionally employs skip connections between the corresponding pair of encoder-decoder layers to promote reconstruction performance. Grachten et al. (2020) encodes the input as a distribution rather than a single value for each dimension. Some models such as Jukebox (Dhariwal et al., 2020), AudioLM (Borsos et al., 2023), and MusicLM (Agostinelli et al., 2023) further insert a vector quantizer between the encoder and the decoder to learn a discrete latent representation. A diffusion model iteratively adds Gaussian noise and then learns to reverse the diffusion process to construct desired data samples from the noise. Kong et al. (2020) proposed DiffWave, a non-autoregressive model that converts the white noise signal into structured waveform through a Markov chain. Chen et al. (2020) combines score matching and diffusion models to generate high fidelity audio samples. Yang et al. (2023), Huang et al. (2023a), and Liu et al. (2023) utilize latent diffusion approach to generate high-quality music.

Since the advent of LLMs, researchers gradually began to explore the application of LLMs in music domain. AudioGen (Kreuk et al., 2022) and Music-Gen (Copet et al., 2023) both use an autoregressive transformer-based decoder (Vaswani et al., 2017) that operates on the discrete audio tokens. Macaw-LLM (Lyu et al., 2023) incorporates visual, audio, and textual information by using an alignment module to unite multi-modal features to textual features for LLM to generate response. M<sup>2</sup>UGen (Hussain et al., 2023) exploits the potential of LLM to bridge multi-modal music comprehension and generation. It utilizes LLaMA2 model to comprehend the multi-modal contextual information of the input and perform downstream tasks such as

269

270

271

272

273

274

275

276

277

278

279

music question-answering and music generationguidance.

Benchmark Evaluations Benchmark evaluation 184 plays a crucial role in assessing the development 185 of LLMs. Previous traditional benchmarking ef-186 forts (Hendrycks et al., 2021; Sakaguchi et al., 187 188 2020) focused on evaluating certain capabilities of models in individual tasks or single-task types. 189 However, with the advancement of LLMs, these 190 benchmarks have become insufficient for comprehensive and accurate assessment of LLM ca-192 193 pabilities. Consequently, researchers have proposed more comprehensive and challenging bench-194 marks (Hendrycks et al., 2020; Li et al., 2023; 195 Huang et al., 2023b) to test whether LLMs pos-196 sess general world knowledge and reasoning abil-197 ity. Additionally, there are task-specific evalua-198 tions such as LawBench (Fei et al., 2023) and Ar-199 cMMLU (Zhang et al., 2023). However, whether in English or Chinese, there is currently a lack of 201 benchmarks for evaluating the musical abilities of LLMs, despite music being an important part of human life. Therefore, we propose ZIQI-EVAL, a benchmark for evaluating the musical abilities of LLMs, to fill the gap in benchmark evaluations of LLMs' musical capabilities.

## **3** ZIQI-Eval Benchmark

#### 3.1 Dataset Curation

210

211

212

213

214

215

216

217

218

219

222

227

228

**General Principle** This dataset integrates the renowned music literature database Répertoire International de Littérature Musicale (RILM), providing a broad research perspective and profound academic insights into the dataset. The inclusion of "The New Grove Dictionary of Music and Musicians" injects the essence of musical humanism into the dataset. Furthermore, dozens of domestic and foreign monographs, such as "Music in Western Civilization" by Paul Henry Lang, the availability of past exam materials from Baidu Wenku, and the advanced data processing capabilities of GPT-4 (Achiam et al., 2023), collectively enhance the data integrity and reliability of the model.

**Data Statistics** ZIQI-Eval dataset consists of two parts: music comprehension question bank and msuic generation question bank.

The music comprehension question bank which is presented in the form of multiple-choice questions consists of 10 major categories and 56 subcategories, encompassing 14244 data entries. It not only includes traditional classifications such as music performance, composition theory, and world ethnic music, but also covers popular music, Western music history, Chinese music history, Chinese traditional music, music aesthetics, and music education. The topics range from popular music, rock music, blues, to female music and more. Additionally, the dataset adopts a decentralized design philosophy, fully showcasing the diversity and inclusiveness of global music cultures.

The music generation question bank consists of 200 questions, testing the ability of music continuation. Considering the difficulty in the evaluation of the generated music, the music generation questions are also presented in the form of multiplechoice questions.

We conduct a comprehensive evaluation of LLMs' music capabilities across the entire dataset. It is worth mentioning that this music dataset has made positive contributions in highlighting female music composers. By including relevant content about female composers, it addresses the gender imbalance in historical literature and promotes progress and inclusivity in the music academic community. This initiative not only reflects the model's profound recognition of gender equality issues but also demonstrates its efforts in advancing the diversification of the music field.

## 3.2 Evaluation Criteria

The evaluation is divided into two parts: music comprehension evaluation and music generation evaluation. The music comprehension evaluation aims to assess the LLMs' music comprehension abilities, specifically their understanding of music harmony, melody, and rhythm. The music generation evaluation, on the other hand, seeks to evaluate the LLMs' capacities for music generation, namely their ability to generate music across diverse styles and genres.

**Music Comprehension Evaluation** We turn the music-related knowledge into the question stem and provide them with options to LLMs, making LLMs to choose the right answer. For example, as shown in Figure 2, take "What is the milestone representative work of Impressionistic orchestral music?" as the stem, "The Sea", "Prelude to the Afternoon of a Faun", "Pelléas et Mélisande", and "Clair de Lune" as the options, we examine whether LLMs can select the right answer "Prelude to the Afternoon of a Faun".



Figure 1: ZIQI-Eval task overview.

Music Comprehension Test         题目:印象主义管弦音乐里程碑式的代表作为。         Question: What is the milestone representative work of Impressionistic orchestral music?         A.《大海》 "The Sea"         B.《牧神午后》 "Prelude to the Afternoon of a Faun"         C.《俱里亚斯与梅丽桑德》 "Pelléas et Mélisande"         D.《月色满庭台》 "Clair de Lune"
Music Generation Test           题目: 请根据输入的旋律选择最匹配的旋律续写片段:           Question: Please choose the most fitting continuation for the given melody based on the input:           X:3           M:2/4           L:1/8           R:Country Dance           N:"Allegro"           "Allegro"D/E/[FFF]F2 dF] {F}EDEF[DA,DE]           FFFF]F2 dF]EDFF[D3;]           :AABB[ccdd]ccBB[A2;]
[:f/g/]aaaa a2 fd[g2ec defg  aaaala2 fd[g2 cc ds]] B. F2F2G2FD   GFD2D2F2-   F8   DFF2GFG2   D2FGD2D2-   D8 C. [:{^fg}a3g ^f2f2[g2 {b}ag/a/ b2 z2]{e=f}g3f e2e2]f2 {a}gf/g/ a2 z2  {de}18e d2^c2]d2 {f}ed/c/ f2 {a}gf/g/a2 bg f2e2 d2 d^d d' 22 z2!D.C.!:[] D. [: d2g f/g/af   e/f/ge d/e/fd   B2f fdB   BeB fdB   d2g f/g/af   e/f/ge d/e/fd   A2e ecA   ABA ecA :]]

Figure 2: Examples of music comprehension and music generation test.

Music Generation Evaluation Given that most LLMs can only accept textual inputs, we utilize ABC notation to convert the musical scores of audio into a textual format, which serves as the input for LLMs. We partition the sheet music written in ABC notation into two segments. The initial segment serves as the question, while the subsequent segment presents four alternative options, also in ABC notation, for the potential continuation of the composition. Then we make LLMs discern the most likely continuation fragment, assessing their music continuation ability. For instance, as shown in Figure 2, we split the original score from "AABBlccddlccBBIA3:!", and test whether LLMs

287

290

have the ability to choose the most fitting option.

295

297

298

299

300

301

302

303

304

306

307

308

309

## Experiments

#### 4.1 Setup

**Baselines** We comprehensively assess 17 LLMs, including API-based models and open-source models. The API-based models contain GPT-4 (gpt-4-1106-preview) (Achiam et al., 2023), GPT-3.5-Turbo (OpenAI, 2022), Claude-instant-1 (An-thropic, 2022), and ERNIE-Bot (Baidu, 2023) series. The open-source models contain Aquila-7B (WUDAO, 2023), Bloomz-7.1B (Muennighoff et al., 2022), ChatGLM2-6B (THUDM, 2023), Mixtral (Jiang et al., 2024), Qwen-7B-Chat (Bai et al., 2023), XuanYuan-70B (Zhang and Yang, 2023), and Yi-6B (01-ai, 2023).

**Metrics** We use a regular expression R, namely 310 r'[ABCD]', to match the answer and consider 311 the first uppercase letter  $\in \{A', B', C', D'\}$ 312 matched as the response. We define Accuracy 313 (Acc.) as the proportion of correctly answered 314 questions among all questions. Precision is the 315 proportion of correctly answered questions among 316 the questions predicted as A/B/C/D. Recall is the 317 proportion of correctly answered questions among 318 the total number of questions that should be an-319 swered as A/B/C/D. In this case, the total number of questions that should be answered as A/B/C/D 321 is actually the total number of questions, so the re-322 call metric is equivalent to the accuracy metric. F1 323 score is the weighted harmonic mean of precision and recall. The specific formulas for these metrics

Models	Music Comprehension Evaluation			Music Generation Evaluation		
	Precision	Recall (Acc.)	F1	Precision	Recall (Acc.)	F1
GPT-4	62.85	100.00	77.19	53.50	97.00	68.96
GPT-3.5-Turbo	-	-	-	30.50	96.00	46.29
Claude-instant-1.2	45.86	71.43	55.86	25.00	99.50	39.96
ERNIE-Bot	49.96	66.69	57.13	29.00	96.50	44.60
ERNIE-Bot-Speed	31.18	74.57	43.97	42.00	100.00	59.15
ERNIE-Bot-Turbo	47.88	94.07	63.46	25.50	100.00	40.64
ERNIE-Bot-8k	53.17	99.16	69.22	26.50	88.00	40.73
Aquila-7B	29.06	62.40	39.65	9.00	40.00	14.69
Bloomz-7.1B	31.97	91.06	47.33	19.00	64.50	29.35
ChatGLM2-6B	39.82	60.52	48.04	15.50	62.50	24.84
Mixtral-8x7B	43.39	99.56	60.44	31.00	100.00	47.33
Qwen-14B	30.04	17.98	22.49	23.66	15.50	18.73
XuanYuan-70B	37.70	46.70	41.72	21.00	89.00	33.98
Yi-6B	60.00	11.06	18.68	0.00	0.00	0.00
Yi-34B	32.24	16.76	22.06	12.12	2.00	3.43

Table 1: Main results(%) of the Music Comprehension Evaluation and Music Generation Evaluation in ZIQI-Eval. Segment 1: API-based models; Segment 2: Open-source models.

326	are as follows:

327

332

333

334

335

339

340

341

343

345

347

349

329 Precisio

$$\hat{y} = R\left(\tilde{X}\right)$$

$$Precision = \frac{\sum_{i=1}^{N} \mathbb{I}\left(\hat{y}_{i} = y_{i}\right)}{V}$$

$$\sum_{i=1}^{N} \mathbb{I}\left(\hat{y}_{i} = y_{i}\right)$$

 $\tilde{X} = G(X)$ 

$$Recall(Acc.) = \frac{\sum_{i=1}^{n} \mathbb{I}(y_i = y)}{N}$$
$$F1 = \frac{2 * Precision * Recall}{Precision + Recall}$$

where G is the LLM generation process,  $\tilde{X}$  is the generated string,  $R(\cdot)$  is applying the regular expression for answer retrieval,  $\hat{y}$  is the predicted answer, V is the number of questions predicted as A/B/C/D, N is the total number of the questions, and  $\mathbb{I}(\cdot)$  is the indicator function.

#### 4.2 Results

Table 1 presents the main results of ZIQI-Eval. Based on the results, we can find that:

I. Overall, the performance of all LLMs on the ZIQI-Eval benchmark is poor. In both music comprehension test and music generation test, the majority of LLMs have not surpassed the passing threshold of 60. Their accuracy rates generally hover between 30 and 50, performing only marginally better than random selection. Even the top-performing model, GPT-4, achieved accuracy rates of only 77.19 and 68.96 in the respective tests. This glaring discrepancy highlights the inadequate consideration given to music accuracy within current LLM models and underscores the formidable challenges posed by the ZIQI-Eval benchmark. 350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

**II. API-based models perform better than open-source models.** In the evaluation of music comprehension test, API-based models generally exhibit higher accuracy compared to open-source models. The accuracy of API-based models is basically distributed between 50 and 70, while opensource models mostly range between 30 and 50. Only specific open-source models like ChatGLM3 and Mixtral can achieve an accuracy higher than 50.

In the evaluation of music generation questions, API-based models consistently outperform opensource models with significantly higher accuracy. The highest accuracy achieved by an API-based model is 68.96, surpassing the highest accuracy of in open-source models.

**III.** The music capabilities of LLMs are dependent but not solely on parameter size. There is a certain degree of relationship between the musical ability and parameter size of LLM models within the same series, while the musical ability of LLM models from different series is not strongly correlated with parameter size.

The ChatGLM series, Qwen series, and Yi series LLMs consistently show improvements in both music comprehension and generation accuracy. Contrary to expectations, the model with
significantly different parameter sizes, ChatGLM26B and XuanYuan-70B, exhibits higher accuracy
in music comprehension for the ChatGLM2-6B
model, surpassing XuanYuan-70B by 6.32. Even
among models with similar parameter sizes, there
can be considerable differences in performance.
For example, the Yi-6B model achieves a music comprehension accuracy of only 28.28, while
ChatGLM2-6B achieves an accuracy of 48.04, resulting in a significant difference of 19.76 between
the two accuracy rates.

IV. The instruction-following abilities of LLMs are not directly linked to their music capabilities. The recall scores of LLMs are strongly correlated with their instruction-following abilities. However, a strong instruction-following capability does not necessarily indicate strong musical capabilities in LLMs. Some LLMs may score highly in terms of recall, but they struggle to effectively comprehend and generate music. Claude-instant-1 serves as a clear example where the subjective recall score reaches 99.5, but the precision is only equivalent to random selection.

V. The music generation capabilities of LLMs are in need of improvement. Even though some LLMs demonstrate a decent understanding of music, their music generation capabilities still have room for improvement. In general, the accuracy for music generation test in LLMs are lower compared to music comprehension test. The difference can be quite significant, such as ERNIE-Bot-8k, where the score for music comprehension test is higher by 28.49 compared to music generation test.

## 5 Analysis

392

394

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

499

423

424

425

426

427

In addition to the overall evaluation of LLMs on the dataset, we are also interested in the models' accuracy for specific categories.

# 5.1 Does LLM show any bias towards questions related to women?

We compare the accuracy of LLMs in the female music theme with the average accuracy obtained by LLMs in the female music theme to analyze whether there is bias in LLMs towards female music. We categorize LLMs into three groups: LLMs without gender bias (above the average accuracy), LLMs with no significant bias (deviating within a range of  $\pm 1.0\%$  from the average accuracy), and

Models	Female	Black	Region		
	1 0111010	Bluen	European	Other	
GPT-4	39.42	6.95	86.95	39.37	
Claude-instant-1	39.38	48.00	51.32	40.00	
Aquila-7B	43.85	46.29	28.69	34.00	
Bloomz-7.1B	53.17	68.83	46.91	49.56	
ChatGLM3-6B	27.27	33.33	11.43	11.61	
Mixtral-8x7B	46.90	67.86	43.39	48.72	
Mistral-7B	17.14	33.33	4.00	4.16	
Qwen-14B	32.51	30.30	21.57	34.46	
Yi-6B	46.90	37.31	46.02	37.45	
Yi-34B	38.26	31.58	37.37	40.77	
Average	38.48	40.38	37.77	34.01	

Table 2: Results(%) of Female Music Accuracy and Black African Music Accuracy. *Female* stands for Female Music Accuracy, *Black* stands for Black African Music Accuracy, and *Region* stands for the accuracy of LLMs regarding World Ethnic Music.

LLMs with gender bias (below the average accuracy).

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

Because we do not fine-tune LLMs, the results reflect the inherent biases of the LLMs themselves. According to the results of Table 2, 40% of the models have no gender bias, 30% of the models are neutral or have no significant bias, and 30% of the models have gender bias, as shown in Figure 3(c). LLMs with accuracy lower than the average accuracy tend to overlook relevant content related to female music themes. Mistral-7B and ChatGLM3-6B, in particular, have significantly lower scores than the average, indicating a notable gender bias issue in these two models. Overall, LLMs exhibit minimal gender bias.

## 5.2 Does LLM exhibit bias toward different races?

We calculate the accuracy of LLMs for the subtopic of Black African music, using the same partitioning method as for determining gender bias, to assess whether there is racial bias in LLMs. According to the results of Table 2, 40% of the models have no racial bias, 10% of the models are neutral or have no significant bias, and 50% of the models have racial bias, as shown in Figure 3(c). The accuracy rates of ERNIE-Bot-Speed and Aquila-7B are below the mean by x and y respectively, indicating a significant racial bias in these two models. Overall, LLMs exhibit minimal racial bias.



Figure 3: Performance of LLMs on gender bias and racial bias.

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

457

## 5.3 Does LLM display bias in terms of region?

We seek to investigate whether LLMs are influenced by Eurocentrism, which positions Europe as the cultural and knowledge center, potentially leading to lower evaluations or neglect of contributions from non-central regions and resulting in biases against these regions. To assess the presence of region bias, we computed the accuracy for the European Music subtheme within World Ethic Music, and the average accuracy for other subthemes within World Ethic Music. Among the LLMs, 30% exhibited higher accuracy rates in European Music compared to other regional music, while 50% of LLMs demonstrated higher accuracy rates in European Music than the average accuracy rate within European Music. These findings suggest that LLMs are influenced by Eurocentrism and exhibit bias towards non-central regions. Most LLMs show similar accuracy between European music and other regions. Surprisingly, GPT-4 exhibits a significantly higher accuracy in European music compared to other regions, with a difference of 47.58, demonstrating a clear bias and regional inclination.

From Figure 3(b), it is evident that LLMs demonstrate similar tendencies towards both gender bias and racial bias, displaying a trend where both ends (with bias and without bias) are relatively higher, while the middle (neutral) is lower. Some LLMs have accuracy rates significantly lower than the mean, such as Aquila-7B with an accuracy rate lower than the mean by 20%, suggesting that LLMs still have a long way to go in eliminating biases. It is worth noting that LLMs with a propensity for gender bias are likely to exhibit racial bias as well, as evidenced by models such as Aquila-7B, ChatGLM2-6B, and Llama-7B. Consequently, it is imperative for future developments in LLMs to address biases comprehensively, not limited to gender and racial biases.

## 6 Futher Analysis

#### 6.1 Phenomenon Analysis of LLMs

To further explore the subjective capabilities of LLMs in the realm of music, we conducted an indepth analysis of the responses provided by each model. Our findings categorize the existing LLMs into three distinct types: 497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

I. Lack of melodic understanding: This type includes LLMs that demonstrate a complete lack of comprehension regarding musical notation. When faced with questions that require the continuation of a melody after a format transformation, these models predominantly resort to evasion, often responding with statements like "Unable to determine, need more information." They fail even to understand the format of the input melody. ChatGLM2-6B and Aquila-7B are prototypical examples of this type, characterized by a high frequency of evasive responses, resulting in a significantly low efficacy in their replies. A notable phenomenon is their tendency to "guess" by consistently selecting option A, leading to most responses without any analytical explanation. For instance, in the responses from ChatGLM2-6B, option A was chosen up to 60%. Besides a preference for option A, Aquila-7B also shows a partiality towards option D.

**II. Limited appreciation, misaligned with human preferences:** A representative model in this type is ERNIE-Bot-8K. This model provides highly interpretable analyses for each option of every question, offering seemingly logical explanations concerning melody, rhythm, and pitch. However, the model's performance, with accuracy barely exceeding that of random selection, underscores the challenge of encapsulating the subjective essence of music appreciation through algorithmic processes. This discrepancy not only highlights the limitations

631

632

633

634

585

535of current AI models in understanding complex,536subjective domains but also underscores the need537for more sophisticated approaches that can better538capture the intricacies of human preferences.

**III. Relatively good appreciation skills:** GPT-4 stands out as a typical example of this type. Its responses consider aspects such as melodic coherence, stylistic similarity, and the seamless integration of musical structures, aligning to a certain extent with human preferences. Further analysis of the questions GPT-4 answered incorrectly revealed a strong inclination towards musical continuity. In many instances, it was observed that GPT-4 prioritized coherence, which led to the selection of incorrect options.

## 6.2 Analysis of GPT-4

541

542

545

546

547

548

549

550

552

554

555

556

557

558

560

562

563

564

567 568

572

574

575

577

578

580

581

582

Taking GPT-4 as a case study, we have gained further insights into the performance of LLMs in the realm of music. The performance of GPT-4 in the domains of women's music and world ethnic folk music indicates a commendable understanding of specific musical areas, reflecting GPT-4's focus on diversity and inclusivity. Women's music and world ethnic folk music, each representing unique cultural and social perspectives, have shown through GPT-4's relatively higher scores the extensive coverage of different cultures and musical traditions.

GPT-4 has demonstrated exceptional performance in the realm of popular music, achieving scores close to 90. This may be due to the abundant and accessible resources in popular music, including lyrics, genres, and artist information. The popularity and media coverage of pop music may also have facilitated the model's learning efficiency in this field.

It has also scored highly in Western music history and music performance, showcasing its capability in processing music history and practical music-making. The higher scores in Western music history over all other regions suggest a certain degree of geographical bias.

In the area of music aesthetics, GPT-4 scored low, revealing a significant weakness. This may be attributed to the complexity and subjectivity of music aesthetics, which might surpass the model's ability to learn from existing textual materials, indicating that there is room for improvement in the model's perception, evaluation, and theoretical analysis of music. Through analysis, we identified that GPT-4 tends to make errors in several distinct categories, primarily falling into three types:

**Matching Errors:** This category encompasses questions related to musical knowledge, specifically matching-type queries, such as identifying the first Hungarian national opera or the composer of "The Song of the Red Flag". GPT-4's responses often affirmatively stated incorrect options, indicating inaccuracies within its knowledge base for specific factual information.

**Comprehension Errors:** These errors involve understanding specific musical terminologies and the relationships between certain concepts. Questions like "What function of art does edutainment refer to?" or "What role do work songs play in labor as a genre of folk music?" exemplify where GPT-4 misinterprets multiple word meanings, leading to a misunderstanding of the intended concept. This suggests a need for improvement in GPT-4's understanding and reasoning within the musical domain.

**Reasoning Errors:** In instances where GPT-4 correctly understands the question and possesses the relevant knowledge background, errors occur during the reasoning or calculation process, resulting in incorrect conclusions. An example can be seen in questions involving the calculation of musical intervals, where GPT-4 confuses semitones and whole tones. This indicates a gap in GPT-4's ability to perform downstream tasks that require precise logical deductions.

## 7 Conclusion and Future Work

Our research sheds light on the oversight of existing evaluations in recognizing the musical abilities of large models. To address this gap, we introduce ZIQI-Eval, a comprehensive benchmark that encompasses 10 major categories and 56 subcategories, comprising over 14,000 data entries. Notably, this benchmark also actively contributes to the acknowledgment of female music composers, rectifying the gender disparity and promoting inclusivity. We conduct a comprehensive experiment involving 15 LLMs, including both API-based and open-source models, to assess their performance in the domain of music. The results indicate that there is significant scope for enhancing the musical capabilities of existing LLMs. We intend to create a multimodal benchmark to evaluate the musical expertise of LLMs in the future.

### Limitations

635

647

651

662

664

674

675

676

677

678

679

683

684

636Our research to date has been exclusively focused637on objective questions, without delving into the638study of subjective questions. One limitation of our639current music benchmark is the absence of multi-640modal data. While the benchmark may excel in641evaluating and comparing the quality and creativ-642ity of musical compositions based on audio data643alone, it fails to incorporate other essential aspects644of the music experience, such as visual elements or645textual information.

#### References

01-ai. 2023. 01-ai. 01-ai Blog.

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
  - 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*.
- Anthropic. 2022. Introducing claude. Anthropic Blog.
- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021.
  Program synthesis with large language models. arXiv preprint arXiv:2108.07732.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. 2023. Qwen technical report. *arXiv preprint arXiv:2309.16609*.
- Baidu. 2023. Wenxin yiyan. Baidu Blog.
- Zalán Borsos, Raphaël Marinier, Damien Vincent, Eugene Kharitonov, Olivier Pietquin, Matt Sharifi, Dominik Roblek, Olivier Teboul, David Grangier, Marco Tagliasacchi, et al. 2023. Audiolm: a language modeling approach to audio generation. *IEEE*.
- Antoine Caillon and Philippe Esling. 2021. Rave: A variational autoencoder for fast and highquality neural audio synthesis. *arXiv preprint arXiv:2111.05011*.
- Nanxin Chen, Yu Zhang, Heiga Zen, Ron J Weiss, Mohammad Norouzi, and William Chan. 2020. Wavegrad: Estimating gradients for waveform generation. *arXiv preprint arXiv:2009.00713*.
- Yi-Hui Chou, I Chen, Chin-Jui Chang, Joann Ching, Yi-Hsuan Yang, et al. 2021. Midibert-piano: largescale pre-training for symbolic music understanding. arXiv preprint arXiv:2107.05223.

Ching-Hua Chuan, Kat Agres, and Dorien Herremans. 2020. From context to concept: exploring semantic relationships in music with word2vec. *Neural Computing and Applications*, 32:1023–1036. 686

687

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

721

723

724

725

726

727

728

729

730

731

732

733

735

736

737

- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- 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*.
- 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*.
- Jesse Engel, Kumar Krishna Agrawal, Shuo Chen, Ishaan Gulrajani, Chris Donahue, and Adam Roberts. 2019. Gansynth: Adversarial neural audio synthesis. *arXiv preprint arXiv:1902.08710.*
- Zhiwei Fei, Xiaoyu Shen, Dawei Zhu, Fengzhe Zhou, Zhuo Han, Songyang Zhang, Kai Chen, Zongwen Shen, and Jidong Ge. 2023. Lawbench: Benchmarking legal knowledge of large language models. *arXiv preprint arXiv:2309.16289*.
- Josh Gardner, Simon Durand, Daniel Stoller, and Rachel M Bittner. 2023. Llark: A multimodal foundation model for music. *arXiv preprint arXiv:2310.07160*.
- Maarten Grachten, Stefan Lattner, and Emmanuel Deruty. 2020. Bassnet: A variational gated autoencoder for conditional generation of bass guitar tracks with learned interactive control. *Applied Sciences*, 10(18):6627.
- Gal Greshler, Tamar Shaham, and Tomer Michaeli. 2021. Catch-a-waveform: Learning to generate audio from a single short example. *NeurIPS*, 34:20916–20928.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. 2021. Measuring mathematical problem solving with the math dataset. *NeurIPS*.
- Rongjie Huang, Jiawei Huang, Dongchao Yang, Yi Ren, Luping Liu, Mingze Li, Zhenhui Ye, Jinglin Liu, Xiang Yin, and Zhou Zhao. 2023a. Make-an-audio: Text-to-audio generation with prompt-enhanced diffusion models. *arXiv preprint arXiv:2301.12661*.

- 738 739 740 741 742 744 745 746 747 748 749 751 752 753 754 755 756 758 759 760 761 762
- 766 767 768 769 770 771 772
- 774 775 776 777 778 779 780
- 781 782 783 784
- 784 785 786 787

790

- Yuzhen Huang, Yuzhuo Bai, Zhihao Zhu, Junlei Zhang, Jinghan Zhang, Tangjun Su, Junteng Liu, Chuancheng Lv, Yikai Zhang, Jiayi Lei, Yao Fu, Maosong Sun, and Junxian He. 2023b. C-eval: A multi-level multi-discipline chinese evaluation suite for foundation models. In *NeurIPS*.
- Atin Sakkeer Hussain, Shansong Liu, Chenshuo Sun, and Ying Shan. 2023. M<sup>2</sup> UGen: Multi-modal music understanding and generation with the power of large language models. *arXiv preprint arXiv:2311.11255*.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088*.
- Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. 2020. Diffwave: A versatile diffusion model for audio synthesis. *arXiv preprint arXiv:2009.09761*.
- Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre Défossez, Jade Copet, Devi Parikh, Yaniv Taigman, and Yossi Adi. 2022. Audiogen: Textually guided audio generation. *arXiv preprint arXiv:2209.15352*.
- Haonan Li, Martin Tomko, Maria Vasardani, and Timothy Baldwin. 2022. MultiSpanQA: A dataset for multi-span question answering. In ACL, pages 1250– 1260.
- Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai Zhao, Yeyun Gong, Nan Duan, and Timothy Baldwin. 2023. Cmmlu: Measuring massive multitask language understanding in chinese. arXiv preprint arXiv:2306.09212.
- Hongru Liang, Wenqiang Lei, Paul Yaozhu Chan, Zhenglu Yang, Maosong Sun, and Tat-Seng Chua.
  2020. Pirhdy: Learning pitch-, rhythm-, and dynamics-aware embeddings for symbolic music. In *MM*, pages 574–582.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. 2022. Holistic evaluation of language models. *arXiv preprint arXiv:2211.09110*.
- Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. 2023. Audioldm: Text-to-audio generation with latent diffusion models. *arXiv preprint arXiv:2301.12503*.
- Chenyang Lyu, Minghao Wu, Longyue Wang, Xinting Huang, Bingshuai Liu, Zefeng Du, Shuming Shi, and Zhaopeng Tu. 2023. Macaw-Ilm: Multi-modal language modeling with image, audio, video, and text integration. *arXiv preprint arXiv:2306.09093*.

Andres Marafioti, Piotr Majdak, Nicki Holighaus, and Nathanaël Perraudin. 2020. Gacela: A generative adversarial context encoder for long audio inpainting of music. *IEEE*, 15(1):120–131. 791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, et al. 2022. Crosslingual generalization through multitask finetuning. *arXiv preprint arXiv:2211.01786*.
- OpenAI. 2022. Chatgpt: Optimizing language models for dialogue. *OpenAI Blog*.
- Darius Petermann, Seungkwon Beack, and Minje Kim. 2021. Harp-net: Hyper-autoencoded reconstruction propagation for scalable neural audio coding. In *WASPAA*, pages 316–320.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2020. Winogrande: An adversarial winograd schema challenge at scale. *AAAI*, 34(05):8732–8740.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2022. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *arXiv preprint arXiv:2206.04615*.
- THUDM. 2023. ChatGLM2. https://github.com/ THUDM/ChatGLM2-6B.
- Jinhao Tian, Zuchao Li, Jiajia Li, and Ping Wang. 2023. N-gram unsupervised compoundation and feature injection for better symbolic music understanding. *arXiv preprint arXiv:2312.08931*.
- Andrea Valenti, Antonio Carta, and Davide Bacciu. 2020. Learning style-aware symbolic music representations by adversarial autoencoders. *arXiv preprint arXiv:2001.05494*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *NeurIPS*, 30.

WUDAO. 2023. Aquila. Github repository.

- Dongchao Yang, Jianwei Yu, Helin Wang, Wen Wang, Chao Weng, Yuexian Zou, and Dong Yu. 2023. Diffsound: Discrete diffusion model for text-to-sound generation. *IEEE*.
- Yi Yu, Abhishek Srivastava, and Simon Canales. 2021. Conditional lstm-gan for melody generation from lyrics. *TOMM*, 17(1):1–20.
- Mingliang Zeng, Xu Tan, Rui Wang, Zeqian Ju, Tao Qin, and Tie-Yan Liu. 2021. MusicBERT: Symbolic music understanding with large-scale pre-training. In *ACL*, pages 791–800.

844	Shitou Zhang, Zuchao Li, Xingshen Liu, Liming Yang,
845	and Ping Wang. 2023. Arcmmlu: A library and
846	information science benchmark for large language
847	models. arXiv preprint arXiv:2311.18658.

Xuanyu Zhang and Qing Yang. 2023. Xuanyuan 2.0: A
large chinese financial chat model with hundreds of
billions parameters. In CIKM, pages 4435–4439.