CONTENTS

# Foundation Models for Music: A Survey

Yinghao Ma, Anders Øland, Anton Ragni, Bleiz MacSen Del Sette, Charalampos Saitis, Chris Donahue, Chenghua Lin, Christos Plachouras, Emmanouil Benetos, Elona Shatri, Fabio Morreale, Ge Zhang, György Fazekas, Gus Xia, Huan Zhang, Ilaria Manco, Jiawen Huang, Julien Guinot, Liwei Lin, Luca Marinelli, Max W. Y. Lam, Megha Sharma, Qiuqiang Kong, Roger B. Dannenberg, Ruibin Yuan, Shangda Wu, Shih-Lun Wu, Shuqi Dai, Shun Lei, Shiyin Kang, Simon Dixon, Wenhu Chen, Wenhao Huang, Xingjian Du, Xingwei Qu, Xu Tan, Yizhi Li, Zeyue Tian, Zhiyong Wu, Zhizheng Wu, Ziyang Ma, and Ziyu Wang

*Abstract*—In recent years, foundation models (FMs) such as large language models (LLMs) and latent diffusion models (LDMs) have profoundly impacted diverse sectors, including music. This comprehensive review examines state-of-the-art (SOTA) pre-trained models and foundation models in music, spanning from representation learning, generative learning and multimodal learning. We first contextualise the significance of music in various industries and trace the evolution of AI in music. By delineating the modalities targeted by foundation models, we discover many of the music representations are underexplored in FM development. Then, emphasis is placed on the lack of versatility of previous methods on diverse music applications, along with the potential of FMs in music understanding, generation and medical application. By comprehensively exploring the details of the model pre-training paradigm, architectural choices, tokenisation, finetuning methodologies and controllability, we emphasise the important topics that should have been well explored, like instruction tuning and in-context learning, scaling law and emergent ability, as well as long-sequence modelling, etc. A dedicated section presents insights into music agents, accompanied by a thorough analysis of datasets and evaluations essential for pre-training and downstream tasks. Finally, by underscoring the vital importance of ethical considerations, we advocate that following research on FM for music should focus more on such issues as interpretability, transparency, human responsibility, and copyright issues. The paper offers insights into future challenges and trends on FMs for music, aiming to shape the trajectory of human-AI collaboration in the music realm.

*Index Terms*—Self-Supervised Learning, Foundation Model, Music Information Retrieval, Music Instruction Following, Music Generation

Yinghao Ma, Bleiz MacSen Del Sette, Charalampos Saitis, Christos Plachouras, Emmanouil Benetos, Elona Shatri, György Fazekas, Huan Zhang, Ilaria Manco, Jiawen Huang, Julien Guinot, Luca Marinelli, and Simon Dixon, are with the Centre for Digital Music, Queen Mary University of London, London E1 4NS, U.K., email: {yinghao.ma, emmanouil.benetos}@qmul.ac.uk

Anders Øland, Chris Donahue, Roger B. Dannenberg, Shih-Lun Wu, and Shuqi Dai are with the School of Computer Science, Carnegie Mellon University. 15213, PA, U.S.

Anton Ragni, University of Sheffield, Western Bank, Sheffield, S10 2TN, U.K.

Chenghua Lin, Xingwei Qu, and Yizhi Li, The University of Manchester, Oxford Rd, Manchester, M13 9PL, U.K.

Fabio Morreale, University of Auckland

Ge Zhang, Wenhao Huang, are with 01.AI, Beijing 100089, CN

Gus Xia, Liwei Lin, and Ziyu Wang, are with Music X Lab, Mohamed bin Zayed University of Artificial Intelligence & New York University Shanghai

Max W. Y. Lam, independent researcher

Megha Sharma, University of Tokyo

Qiuqiang Kong, is with the Chinese University of Hong Kong

Ruibin Yuan, Zeyue Tian, are with the Division of Emerging Interdisciplinary Areas, The Hong Kong University of Science and Technology, Hong Kong SAR

Shangda Wu, is with the Central Conservatory of Music, CN

Shun Lei, is with the Tsinghua Shenzhen International Graduate School, Tsinghua University

Shiyin Kang, is with the Skywork AI PTE. LTD., Beijing

Wenhu Chen, is with the Department of Computer Science, University of Waterloo, N2L 3G1, CA & Vector Institute, CA

Xingjian Du is with University of Rochester

Xu Tan, is with Microsoft

Zhiyong Wu, is with the Tsinghua Shenzhen International Graduate School, Tsinghua University

Zhizheng Wu, is with the Chinese University of Hong Kong, Shenzhen.

Ziyang Ma, is with Shanghai Jiao Tong University

## I. INTRODUCTION

Music is an important part of human culture, universal in its cross-cultural presence, yet taking many different forms across cultures. Its functions include emotion regulation, communication, and promoting social cohesion; it appears in art, entertainment, worship, and advertising; and represents a large industry contributing to the global economy. It presents opportunities to benefit both human society culturally and music industries economically, as well as unique technical challenges when combined with AI.

The field of computer music is at the intersection of music, computer science, electrical engineering, and artificial intelligence, drawing upon fields such as philosophy (aesthetics), psychology (perception, cognition, and production), and physics (acoustics). Computational approaches to music often employ signal processing and other techniques to extract features from audio signals, and then apply machine learning algorithms for music information retrieval (MIR) tasks or music composition.

Although natural language processing, computer vision, and speech processing have widely used foundation models (FMs), we are still only scratching the surface of AI for art, of which music is an essential component. One challenge specific to music is polyphonic signal modelling. Unlike speech and language signals, music usually has several simultaneous "speakers", and the "meaning" of what they "say" is not grounded in real-world objects or events. The occurrences of different note events are not independent, making it a challenging modelling task to capture the "language(s)" of music. Moreover, music typically has a much longer duration with a much higher sample rate compared to speech or general audio, making it harder to model the whole musical piece.

Recent advances in pre-trained language models (PLMs) significantly outperform traditional algorithms on a range of music-related computational tasks, demonstrating the potential

of modern machine learning techniques to understand and process music on an unprecedented scale [LYZ$^+$24]. However, a critical bottleneck has emerged in terms of dataset size and quality. For algorithms to be reliable, especially those deployed in complex, realistic scenarios, they need to be trained on diverse and representative datasets. The performance of these algorithms is deeply dependent on the size of the annotated dataset and the quality of its annotations, which justifies the need for large quantities of high-quality data. Unfortunately, music datasets are often size-constrained due to the limited availability of copyright-free public domain data and the high costs associated with labelling and annotation.

FMs address this problem by employing self-supervised learning (SSL) approaches for pre-training on a large amount of unlabelled music data. SSL enables the model to learn meaningful representations without the need for explicit labelling, by exploiting intrinsic structures within the data. This approach is similar to the natural human learning process. For example, when children hear different instruments played, they learn the characteristics of each unknown instrument and are able to identify the instruments in new pieces of music without necessarily knowing their names. Similarly, SSL enables machine learning models to derive general knowledge from large unlabelled datasets, thereby improving their performance on downstream tasks that lack large amounts of labelled data. As has proven successful across other domains, models trained through such approaches show promising results for music understanding and generation.

### A. What is a Foundation Model?

The term *foundation model* was coined to describe a multi-purpose machine learning model that, rather than being trained for a single specific task, serves as the basis of multiple derived models that are able to perform a wide range of tasks [BHA$^+$21]. This terminology reflects the shift from the traditional emphasis on the specifics of architectures or tasks to a focus on broadly applicable models whose emergent abilities and generalisation are unlocked by significantly increasing the number of model parameters [WBZ$^+$21], [CND$^+$22]. Contrary to terms such as *large language models* or *self-supervised learning*, which emphasise narrow aspects of AI development, *foundation model* captures the essence of the generality of these models.

The rise of foundation models has come about due to advances in computational hardware, architectural innovations in neural networks (e.g., the Transformer architecture), and an enhanced focus on minimally supervised training paradigms. Foundation models employ deep neural network architectures that are typically trained on large-scale unlabelled datasets with SSL. Following this *pre-training* phase, foundation models can be adapted for various downstream tasks via a relatively lightweight finetuning or in-context learning stage, for example, using a labelled dataset that is orders of magnitude smaller than the pre-training data.

Beginning with language models such as Google's BERT (Bidirectional Encoder Representations from Transformers [DCLT18]) and OpenAI's GPT (Generative Pre-trained Trans-

former [BMR$^+$20]) series, foundation models have demonstrated the power of SSL for training on extensive web-sourced datasets, freeing researchers from reliance on labelled data, which does not scale economically to web-scale data sizes. In addition to text analysis and text generation, such PLMs have also demonstrated their utility across various modalities, including image processing with CLIP [RKH$^+$21a], DALL-E [RPG$^+$21], and Flamingo [ADL$^+$22]; speech and audio generation with Audiobox [VSL$^+$23]; music generation with Jukebox [DJP$^+$20a], MusicLM [ADB$^+$23] and MusicGen [CKG$^+$24]; and robotic control with RT-2 [BBC$^+$23].

The release of Stable Diffusion[1] and ChatGPT[2] in 2022 marked a significant turning point for foundation models, in terms of public impact, as well as industrial and academic interest in the creation of AI-generated content (AIGC). This significant progress is primarily due to the ability to follow language instruction, emergent ability in algorithmic advances when scaling up to large language models (LLMs), and the authentic quality of Latent Diffusion Models (LDMs) [RBL$^+$21]. These methods suggest a paradigm shift in AI, as generalised frameworks can support multiple applications across disparate domains. Although developing AI with universal capability for multiple and unseen tasks has been the goal of AI researchers since its earliest days [NSS59], most AI research in the ensuing decades has focused on a single or a limited number of pre-defined task(s). In addition, access to advanced problem-solving capabilities through natural language interaction facilitates uptake by non-specialists. Although the development of foundation models demands substantial financial and computational investments plus significant human effort, the adaptation of pre-existing models for tailored needs is more cost-effective, and the release of open-source foundation models like Stable Diffusion, Llama [TLI$^+$23a], Mistral [JSM$^+$23], and MAP-NEO [ZQL$^+$24] gives access to users, developers and researchers alike to explore the possibilities of the models.

In this paper, we will discuss two types of self-supervisedly pre-trained foundation models which can perform multiple downstream tasks. The first is the single-modality pre-trained model in the waveform or symbolic domain that requires finetuning on downstream tasks. This can be some variant of PLM for music understanding such as MERT[LYZ$^+$24] or music generation such as Jukebox[DJP$^+$20a]. The second is multimodal pre-trained models that can take both natural language and music as input and have the potential to solve downstream tasks with in-context learning. This includes Latent Diffusion Models (LDMs) with multiple text inputs such as MusicLDM[CWL$^+$23a], a music encoder prepended to an LLM such as Mu-llama[LHSS24] or an LLM with multimodal tokenisers such as AnyGPT[ZDY$^+$24], Gemini 1.5[RST$^+$24] and GPT-4o.

### B. Why Foundation Models for Music?

FMs for music not only address data scarcity and reduce annotation costs, but also enhance generalisation in music information retrieval and creation. By pre-training on large

---

[1]https://github.com/CompVis/stable-diffusion
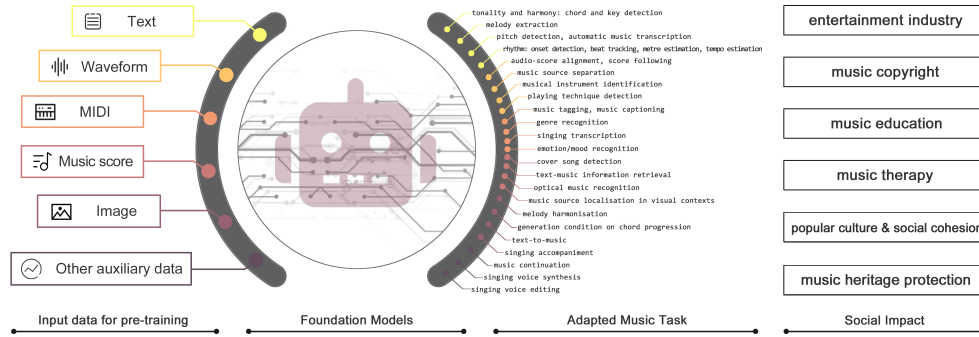[2]https://chatgpt.com/

Figure 1: The input modalities, downstream applications and social impacts of foundation models for music

music datasets, these models provide a better understanding capability of unseen structures, genres, or instruments. These algorithms can also contribute to the protection of the cultural heritage of music through world music analysis, music education, and new forms of artistic expression.

*1) Impact on Industry:* Foundation models have, or will potentially have, more robust and commercially viable applications for music than previous methods, including in creative processes, music understanding, and approaches within the entertainment industry.

In the area of **Creative Applications**, AIGC is perhaps the most obvious application of foundation models, including music, such as personalised music generation and collaborative composition with musicians. Foundation models enable the generation of music based on user-specified preferences as input such as genre, mood, tempo, and instruments. Following recent progress in LLMs and LDMs in music, many music generation startups with proven commercial impact such as SunoAI[3], TiangongAI[4] and Udio[5], have emerged. Musicians and producers can manipulate aforementioned parameters to steer the composition processes, assisting in the ideation process. Such music generation applications can enable new forms of interaction with users and musicians. Music can change based on the listener's feedback or input prompts, potentially creating more immersive and personalised listening experiences. Additionally, FMs show potential in collaborating with musicians or music editors by following their instructions more professionally and robustly.

Foundation models address several aspects of **Music Understanding**. By analysing listening habits and understanding musical preferences, FMs can offer listeners more personalised recommendations, improving the user experience on streaming platforms. FMs may also be used to better detect cover songs and identify copyright infringement, helping artists and companies protect their intellectual property more efficiently. They may also provide analyses of musical pieces to aid musicologists in understanding music structure, characteristics, and innovation.

For **Entertainment and Media**, foundation models can create adaptive soundtracks that correspond to the narratives of visual media for musicians and music editors, enhancing the impact and immersion of movies and video games.

*2) Social Impact:* Foundation models for music, with their capacity to understand, generate, and process music, can provide profound implications for culture and society. As FM has the potential to better resolve all kinds of music-related tasks, most main applications of MIR can be regarded FM territory, and therefore FMs have the potential to change the way we interact with, preserve, and understand music, raising significant ethical and heritage considerations.

Concerning **Cultural Preservation and Diversity**, foundation models can play a role in preserving world cultural and musical traditions that are at risk of being lost. By analysing diverse musical datasets, these models can identify unique characteristics of styles, compositions, and performances from cultures around the world, much similar to current LLMs' capability of understanding minor languages. Moreover, FMs can promote cultural awareness by facilitating exploration of music from different parts of the world.

For the field of **Music Anthropology**, foundation models may serve as a tool for studying the evolution of music across different nations and eras. By analysing vast amounts of music data, FMs can uncover music patterns and cultural influences. By relating this analysis to social and historical data, FMs could potentially provide insights into music's role in human societies.

Foundation models can improve access to **Music Education** by creating personalised learning experiences adapted to the learner's pace and style; For example as a virtual tutor that provides theoretical and practical knowledge, feedback, virtual accompaniment and simulated ensemble playing. This could make music education more accessible, regardless of access to traditional music education resources, encouraging a more inclusive culture of music learning, and removing barriers that have historically limited participation in music-making.

In **Music Therapy**, FMs could be tailored to produce music for therapeutic purposes, aligning with individual therapeutic goals or emotional needs, potentially offering mental health support. Likewise, in non-clinical settings, by generating music that reflects or counteracts listeners' emotional states, foundation models can play a role in mood regulation and

---

[3] https://suno.com/

[4] https://music.tiangong.cn/

[5] https://www.udio.com/

wellness practices.

The ability of foundation models to generate music that mimics human compositions raises important **Ethical Considerations**. The fact that the models benefit from the intellectual property of the millions of musicians and artists who created the training data leads to legal challenges and debates about the lawful use of data. Ethical discussions focus on issues of copyright, originality, and the role of AI in creative processes, ideally with interpretability and transparency. As these models become more prevalent, society must navigate between leveraging technology for innovation in music creation and respecting the rights and contributions of human artists.

The impact of foundation models for music is likely to be profound, offering new tools for the generation, analysis and interaction of music, as well as for music education and therapy. As these models are developed, it is essential to consider their ethical implications thoughtfully, ensuring that they serve to enrich human culture and promote a more fair and inclusive global society. For more information about the ethical issues of FMs in music, please refer to Section VI.

### C. Aims of the Survey

The purpose of this survey is to provide a thorough and comprehensive overview of foundation models as they relate to the domain of music, including LLMs and LDMs. While some previous overview papers have addressed FMs [BHA+21] or LLMs [ZZL+23], [HLC+24] in general and as they apply to specific areas such as vision [ZHJL24], speech [ZLL+23], [MMB+23], [LCP+23] and audio [WCL+24], [MLPW22], [LSS+23], [TTG+24], they do not comprehensively cover music-related applications of FMs. Besides, previous music surveys also exhibit limitations in providing a comprehensive overview of FMs. For example, [JLY20] fails to incorporate new advancements post-2021, particularly in LLMs and audio LDMs. Similarly, [HSF+24] focuses predominantly on digital signal processing methods, neglecting the integration of FMs into music synthesis and understanding. [HOHOB22] briefly mentions LLMs and LDMs but lacks an in-depth exploration of their applications in music understanding as well as multimodality. [ZBRR23] provides a limited discussion on music generation models, primarily focusing on commercial scenarios and overlooking critical technical details and ethical considerations.

Our survey aims to bridge this gap by reviewing the wide array of FM applications, from music understanding to generation, therapy, and the ethical implications associated with these technologies. By doing so, we seek to highlight the unique challenges and opportunities that musical data presents for FMs, including aspects such as modelling long-term temporal dependencies, and the evaluation of artistic outputs. Additionally, this survey endeavours to update the literature with recent advances in LLMs and audio LDMs that are not covered in existing surveys.

The survey provides a detailed exploration of FMs in music. Section 2 examines music modalities and representations, including psychoacoustics, audio representations, symbolic music representations, and their integration with other modalities.

We then turn to the diverse applications of FMs in music in section 3, spanning understanding, generation, and medical application. Section 4 covers FMs' technical aspects, focusing on pre-training strategies, (instructional) finetuning, model architecture, audio tokenisation, applying LLM foundation models, music agents, scaling laws and emergent abilities, along with future work. The discussion in Section 5 extends to datasets and evaluation approaches, highlighting the challenges and solutions in both acoustic and symbolic domains of music understanding and generation tasks. The last sections critically assess the ethical and social impacts along with copyright concerns of utilising FMs in music. They address potential cultural issues, including transparency and interoperability of algorithms, effect on humans, people's responsibility, and copyright issues. We suggest researchers in general ML focus on sections 2 and 3, and computer music researchers focus on FM methodology in section 4. For quick start, please refer to the GitHub repository[6].

## II. REPRESENTATIONS OF MUSIC

Music representations can be used to describe, encode, and communicate musical information mentally and computationally. These representations can capture different aspects of music, such as pitch, rhythm, harmony, dynamics, expression, timbre or structure. Humans perceive acoustic music signals in both the temporal and frequency domains, at multiple time and frequency scales, and abstract this rich sensory information to perceptual representations (experiences) of musical concepts such as pitch, rhythm, dynamics, and timbre, among others. Such a process inspires people to design a broad range of music representations with often very different focuses and characteristics. Many representations have been adopted to or exclusively devised as computer representations of music, and yet more are developed for or used as input representations to neural networks. Studying music representations in the context of foundational music models is essential because the choice of representation affects the model's effectiveness, efficiency and performance. Appropriate representations ensure that the data is relevant and informative, enhancing accuracy and reducing computational load. They may also help extract meaningful features and aid in the generalisation of new data. A suitable representation may also incorporate domain-specific knowledge, better align with the task requirements and facilitate interpretability, ultimately leading to more robust, transparent and explainable models.

In this section, we discuss manually designed computer representations of musical notation and audio content at various levels of abstraction, starting from a very brief overview of how sound and musical concepts are experienced by humans. Additionally, the representations of foundation modelsF(Ms) can also be fully learnable audio tokens other than manually designed representations; please refer to Section IV-C for more information. Finally, In the last subsection of this part, we discuss multimodal representations for music.

We can tell from the following paragraphs that the current FMs for music are developed on limited modalities, including

waveform, MIDI or ABC notations. And FMs with more comprehensive input modalities remain underexplored.

### A. Music Perception & Notation

Humans process musical information through a complex interplay of perceptual and cognitive processes. A sound wave transmits a pattern of oscillations, created by an excitation force (e.g., bowing) that stimulates a vibrating object (e.g., a viola's string), through a physical medium such as air [CSS19]. On the way from the source (e.g., bowed viola string) to our ears, the sound waves carry information about the vibrating object, the driving force, and possibly other interactions through the physical medium (e.g., the walls of a room). When the waves reach the ears, these oscillations are translated into a pattern of neural pulses that the brain processes to make sense of musical elements such as melody, harmony, rhythm, and timbre. The ear acts like a set of overlapping bandpass frequency filters [Lyo17]. Sound energy within a frequency range is integrated across frequency and time (the temporal pattern of neural pulses can also encode frequency). The notion of the *critical band* (the range of frequencies that are integrated) is important psychoacoustically because it provides a basis for explanations of the degree of masking and loudness summation, but also of frequency discrimination.

Pitch may be defined as the human perceptual correlate of acoustic frequency. Specifically, the perceptual representation of pitch involves at least two dimensions: a circular dimension of *pitch chroma* (the *pitch class*, sometimes called *tonality*) and a vertical dimension of *pitch height* [She82]. A pitch class is a set of all heights of pitches that are perceived as repeating once per octave (e.g., all C notes). Both pitch chroma and pitch height dimensions generally allow ordering pitches on a scale from low to high, but not always unambiguously (e.g., Shepard tones [She64]). Furthermore, the relationship between pitch and frequency is not linear: we can better detect differences in lower frequencies than higher frequencies [SVN37]. For example, we can tell the difference between 500 and 1500 Hz, but probably struggle to discriminate between 9 and 10 kHz, although the frequency distance remains the same. This is because the width of the critical band increases at higher frequencies (frequencies are smeared) and decreases at lower frequencies (frequencies can be resolved).

Musically, chroma is the most important part of the pitch and forms the basis of melody (distinct pitches played in sequence) and harmony (distinct pitches coinciding with one another), but also of rhythm. Rhythm is what orders the movement of melodic and harmonic patterns in time, both in the sense of how listeners perceive and how musicians perform music, consisting of *beats* (perceived pulses that mark equally spaced points in time), *tempo* (their frequency), and *meter* (perceived cyclical groupings of beats into units marked by accents; the *time signature*). Staff notation is a *symbolic* visual representation of music that allows musicians to read the pitch and rhythm of notes they are supposed to play, and forms the basis of digital systems to encode musical documents in a machine-readable format. The pitch (class) and octave (height) of a note are indicated by the vertical position of the note within, below, or above the staff. Notes have different durations or note lengths, represented by whole notes, half notes, quarter notes, eighth notes etc. Each note length has a specific duration relative to the beat defined by the meter. Roughly speaking, rhythms and melodies are made up of notes with different durations.

Timbre, or *tone colour*, is a "catch-all" term, as it has been called, that refers to everything on what a sound or musical piece "sounds like," except for pitch information, loudness, and rhythm-related features, although changes in pitch, dynamics, and timing also produce changes in timbre [SSM+19]. Timbre is primarily determined by the relative amplitudes of frequency components. Phase can also play a role, especially when certain changes in the relative phases of harmonics within a critical band can change the resulting envelope modulation [Lyo17]. Temporal characteristics of the signal (e.g., attack time, decay) and vibrato, are key contributors to timbre too [APS05]. However, changes in harmonic amplitudes affect timbre in a much more direct way, such that timbral (i.e., perceptual) distances between sounds correlate well with spectral differences, such as measured in terms of the log power outputs of a one-third-octave filterbank [Plo76]. Such an approach to timbre is not accurate or complete [Lyo17], but captures a reasonable part of the way a sound or musical piece "sounds" and, together with non-linear frequency scales, is the basis for computer representations of musical spectra.

### B. Computer Representations of Music

Inspired by the process through which humans perceive music, as mentioned in the previous subsection, we will introduce the various computational methods for music representation. Echoing how the ear transforms sound waves into neural signals, we begin by detailing the conversion of the waveform in the time domain into spectral representations using bandpass frequency filters. Subsequently, we delve into symbolic music representations at both performance-level and note-level that abstract neural signals in terms of pitch, timbre and time through our brain. Though MIDI (Musical Instrument Digital Interface) in performance-level and ABC notations in note-level are the most widely-used presentations for symbolic music LLMs, we introduce other music representations which have been used with traditional deep learning methods in this subsection as they have potential to be trained for foundation models. We then explore more abstract note-level representations which prioritise relational and structural aspects of music notation. An example is the ABC notation system, which organises music by relative beats in staff notation rather than specific timestamps, offering a more human-understandable representation suited to computational models. At last, we will discuss combined representations like note-tuples and the potential of learnable note-level tokenisers. These advanced topics pave the way for integrating these representations into foundation models, enhancing our ability to process and generate music through AI.

*1) Acoustic-level Representations of Music:* **The log Mel spectrogram** is a key audio representation technique that synthesises signal processing with psychoacoustic principles to closely mirror human auditory perception. It starts with converting audio signals from the time domain to the frequency domain using the Fast Fourier Transform (FFT). The Short-Time Fourier Transform (STFT) further analyses these signals over time, segmenting the audio into overlapping frames and applying FFT to capture temporal dynamics. Transitioning to the Mel scale involves applying Mel band-pass filters to the STFT output. These filters, aligned with the non-linear nature of human pitch perception, group frequencies into bins that are linearly spaced at lower frequencies and logarithmically at higher ones. The logarithmic transformation of this Mel-scaled output adjusts for the human ear's non-linear response to loudness and pitch, compressing the spectrogram's dynamic range and emphasising perceptually significant changes in sound energy. The log Mel spectrogram has been directly applied in tasks of music generation [HSR$^+$22], [FM22a] due to its perceptual relevance, enabling accurate modelling of musical characteristics. In the foundation model paradigm, however, log Mel spectrograms usually serve as input into audio representation models, which learns an efficient latent space [GLCG22a], [FFH24], [HXL$^+$22b] in which generation operates.

**The Mel-Frequency Cepstral Coefficients (MFCCs)** are a crucial audio feature extraction technique that encapsulates the characteristics of human speech perception. This process begins with converting audio signals from the time domain to the frequency domain via FFT. It then applies the STFT to segment the audio into overlapping frames and analyse temporal changes. The Mel scale is employed through Mel filters that mimic the human ear's logarithmic frequency perception, focusing on discriminative lower frequencies. The computation of MFCCs involves taking the logarithm of the energies in each Mel filter, followed by a Discrete Cosine Transform (DCT). Due to their effectiveness in capturing vocal tract configurations, MFCCs are widely used in speech recognition and speaker identification applications. AV-HuBERT [SHLM22], Hubert [HBT$^+$21a], and MusicHubert [MYL$^+$23] use MFCC features as re-prediction targets.

**The Constant-Q Transform (CQT)** enhances audio processing in music analysis by providing a log-frequency spectrogram aligned with the musical scale. Unlike the linear Fourier Transform, the CQT uses a logarithmic scale that reflects the exponential nature of musical pitch, simplifying the extraction of note frequencies. Its key feature is that the centre frequency to bandwidth ratio remains constant (denoted by Q), allowing variable filter lengths that optimise performance. Despite less popularity compared to the log-Mel spectrogram, CQT has proven useful for tasks like musical timbre transfer [HLA$^+$19] and music representation learning [LYZ$^+$24].

**Chroma features** focus on the twelve pitch classes of Western music, capturing the harmonic core of a piece regardless of timbre or instrumentation. These features, displayed in chromagrams, condense music into pitch class profiles that highlight the harmonic and melodic structure. Chroma features are especially effective in identifying chords, key sig-natures, and modulations due to their alignment with the equal-tempered scale. The most recent application of chromagram is melody-guided music generation in MusicGen [CKG$^+$24].

*2) Symbolic Music Format & Their Content:* In this subsection, we introduce the low-level symbolic music formats and their informational content at both performance-level and note-level. Symbolic music formats are essential for representing musical notation and facilitating digital music processing. This subsection will cover key formats, highlighting their structures, capabilities, and limitations, especially in the context of training computer music foundation models.

*a) Performance-level symbolic music repsentation:* **Musical Instrument Digital Interface (MIDI) protocol** is a standard way of saving and transferring MIDI data for playback on different systems. Unlike MusicXML or MEI, MIDI encodes musical performance data, capturing information about note pitches, durations, velocities, and other performance aspects like dynamics, articulations, and control changes. This makes MIDI particularly suited for applications requiring precise control over musical performance and real-time interaction between devices. In ML, the detailed performance data in MIDI is invaluable for training models focused on tasks such as music generation, transcription [BDDE19], Optical Music Recognition (OMR) [SF20] and style transfer [BKWW18]. Models can learn from the nuances of human performances captured in MIDI data, enabling the generation of expressive and realistic musical outputs. MIDI's widespread use and standardisation mean there is a vast amount of data available for training purposes. However, the format's focus on performance rather than notation means it lacks explicit information about musical structure, such as key signatures, time signatures, and notation details.

Note-level representations focus on the structural elements of music, such as pitch and rhythm, offering a simplified, human-readable format ideal for sharing musical scores. In contrast, performance-level formats like MIDI encompass additional expressive details of a piece's execution, including dynamics, musical techniques and duration/tempo variations. While note-level symbols capture the composition's essence, MIDI provides a detailed account of its performance, capturing both the notation and its expressive execution.

*b) Note-level symbolic music representation:* The following paragraphs introduce the diverse note representation of symbolic music which are shown in Figure 2.

**MusicXML** is a significant format in the digital music domain, encapsulating the complexities of musical notation in a universally transferable manner [Goo12]. MusicXML is rooted in the Extensible Markup Language (XML), and offers a textual format for encoding music scores (see Figure 2f), ensuring human and machine readability. This format enables a detailed representation of musical elements, from the pitch and duration of notes to their graphical representation in sheet music. By design, MusicXML does not encode all document information in a fully-featured music notation system, which limits capabilities for other non-CWMN scores.

Despite MusicXML's strength in accurately representing musical notation and the fact that many music scores are shared and published in MusicXML format, it is not com-

(a) Snippet from a musical score in staff notation.

(b) LiyPond

```
\relative c'' {
  \key es \major
  \time 3/4
  \clef treble
  b4  ...}
```

(c) **kern

```
*M3/4
*kf
*clefG2
4b
*-
```

(d) ABC

```
X: 1
T: snippet
M: 3/4
K: Eb
L: 1/4
B
```

(e) GUIDO

```
[ \\key<es \\major>
  \\meter<3/4>
  \\clef<"G">
  (b5 q)  ...]
```

```
<measure number="1">
  <attributes>
  <divisions>1</divisions>
  <key> <fifths>-3</fifths> </key>
  <time><beats>3</beats><beat-type>4</beat-type></time>
  <clef><sign>G</sign><line>2</line></clef>
  </attributes>
  <note>
  <pitch><step>B</step><octave>5</octave></pitch>
  <duration>4</duration><type>quarter</type>
  </note> ...
</measure>
```

(f) MusicXML

```
<scoreDef xml:id="scoredef-35" key.sig="3f" meter.count="3" meter.unit="4">
  <staffGrp xml:id="staffgrp-11">
    <staffDef xml:id="staffdef-59" clef.shape="G" clef.line="2" n="1" lines="5" />
  </staffGrp>
</scoreDef>
<section xml:id="section-33">
  <measure xml:id="measure-32" right="single">
    <staff xml:id="staff-25" n="1">
      <layer xml:id="layer-55" n="1">
        <note xml:id="note-02" dur="4" oct="5" pname="b" />
        ...
      </layer>
    </staff>
  </measure>
```

(g) MEI

Figure 2: Symbolic music representations for the same piece of music

monly used to train FMs for music in previous works. For one thing, this is due to its symbolic nature, like ABC notation, which, while rich in notational detail, lacks the audio information necessary for models to learn the acoustic properties of music. Foundation models, especially those geared towards music generation and understanding, require data that encompass the timbral, dynamic, and expressive aspects of music as it is heard, not just as it is written. For another, the process of encoding or decoding MusicXML can be cumbersome for AI models, which thrive on larger datasets and longer context lengths for development. The complexity and verbosity of XML-based formats may introduce additional challenges in processing and interpreting data efficiently, but may also provide more benefits on music details compared to other music notations. Therefore, while MusicXML excels in notation accuracy and interoperability between software, its application in training FMs is limited by the challenges and still underexplored.

**Music Encoding Initiative (MEI)** aims at creating a digital format to represent music notation [Rol02]. MEI utilises an XML-based schema, see Figure 2g, providing rules for documenting the intellectual and physical properties of music notation, enabling consistent search, retrieval, display, and exchange of information across platforms. MEI's modular and extensible structure supports encoding various music notation systems, including common Western notation, mensural (Renaissance), and neume (Medieval) notations. This flexibility ensures an accurate representation by preserving the unique structure and semantics of each notation system rather than merely imitating them visually. One of the primary goals of

MEI is to create a semantically rich model for music notation. It enables the encoding of traditional facsimile, critical, and performance editions. Its foundation on open standards and platform independence promotes the development of extensive and international archives of notated music, which serve as essential resources for editions, performances, analyses and research. In comparison, while MEI and MusicXML encode musical elements such as notes, staves, rests, and clefs using XML, they differ in focus. MusicXML primarily facilitates interchange between notation editors. In contrast, MEI provides a more detailed and systematic encoding of notation information, supporting various notation systems beyond standard common Western notation. Moreover, MEI can document relationships among digital page images, notational features and audio recordings, offering a more comprehensive framework for music representation. Similarly to MusicXML, MEI faces challenges in training computer music FMs. Its complexity, lack of standardisation, and focus on detailed musical notation make it less suitable for ML, which requires more streamlined and uniform data formats. MEI's emphasis on visual representation rather than audio features, large data files, and limited integration with machine learning tools further hinder its effectiveness. Foundational models benefit more from compact, standardised data representations [BHLP16].

**LilyPond** is a syntax similar to LATEX, allowing users to write musical notation in plain text files, which are then compiled into engraved scores [NN03]. The syntax is comprehensive (see Figure 2b) covering a wide range of musical symbols and formatting options, making it suitable for complex and professional-quality scores. It has been previously used for

automatic transcription of organ tablatures [SKM+21] and polyphonic music transcription [CS17]. However, its complexity and focus on visual presentation can pose challenges for model training and efficiency.

**\*\*kern** format is a symbolic music representation system designed for encoding and analysing Western music notation. Developed as part of the Humdrum Toolkit [Hur02], \*\*kern aims to provide a flexible and comprehensive way to represent musical scores in plain text. The format encodes musical information such as pitch, rhythm, meter, and articulation using a straightforward syntax, see Figure 2c, facilitating both human readability and computational processing. \*\*kern is highly extensible, allowing users to include additional musical parameters and metadata as needed. Its standardised text-based structure makes it suitable for training models in symbolic music analysis, pattern recognition, and music generation [RSIM21]. Models can leverage \*\*kern to explore musical structures, identify stylistic features, and generate compositions that adhere to specific musical conventions. The flexibility in encoding can lead to variations that require careful handling to ensure consistent and effective model training. Additionally, while \*\*kern captures a broad range of musical elements, it may lack some of the expressive nuances found in performance-level data, potentially limiting its application in models focused on expressive performance analysis.

**ABC notation** represents a concise and computer-friendly approach to musical notation, employing a simple alphabetic system alongside ASCII characters to depict musical elements. It primarily utilises the letters a–g and A–G to denote notes, with 'z' marking rests (see Figure 2d). Additional notation specifies musical nuances such as sharps, flats, octave shifts, note lengths, keys, and ornamentation. This notation system, which integrates elements from Helmholtz pitch notation, was designed to simplify the sharing of music online and provide a straightforward language for software development, paralleling the simplicity of tablature and solfège. The structure of ABC notation is divided into a header and a body. The header contains metadata such as reference numbers for organising multiple pieces, titles, time signatures, default note lengths, and keys. Conversely, the body focuses on the musical content, encoding each note and rest into tokens that reflect pitch, duration, and rhythmic placement, delineated by bar lines. Like other symbolic music notations, ABC notation is compatible with natural language notations, making it an ideal candidate for developing text-to-symbolic music models, such as chatMusician [YLW+24a]. This compatibility allows for the seamless integration of ABC notation with AI models designed to convert textual descriptions into musical compositions, offering a bridge between linguistic descriptions and musical notation. By leveraging the simplicity and clarity of ABC notation, developers can create models that not only generate music from text but also understand and manipulate music in its symbolic form.

**GUIDO Music Notation (GMN)** is a text-based format designed to represent musical scores in a human-readable and machine-processable way [HHRK98], [HHRK01]. It aims to provide a robust and flexible method for encoding music notation, focusing on simplicity and ease of use. The format utilises a straightforward syntax to represent pitches, durations, dynamics, articulations, and other musical elements. Its hierarchical structure allows for the clear organisation of musical information, making it suitable for a wide range of musical styles and notational complexities. GUIDO's design prioritises both readability and computational efficiency, facilitating its use in various digital music applications. The basic concept behind the GUIDO design is to represent simple musical ideas in a simple manner, reserving complex representations for complex notions. GMN's simplicity and structured syntax (see Figure 2e) make it an option for training models in symbolic music analysis, composition, and pattern recognition. ML models can utilise GUIDO to analyse musical structures, generate new compositions, and study stylistic patterns across different genres. The human-readable format also aids in debugging and understanding the encoded data, which is advantageous during model development and evaluation. However, the trade-off for GUIDO's simplicity is its potential lack of detail compared to more complex formats like MusicXML or MEI. Although GUIDO captures essential musical information, it may not encompass all the expressive nuances required for detailed performance analysis or highly intricate compositions. Additionally, variations in encoding practices can introduce inconsistencies, necessitating careful preprocessing for effective model training. This format has yet to be explored in the realm of ML/DL.

**Ontologies for music representation:** In the past two decades, along with the development of Web Ontologies, several music-related ontologies were also developed. However, they either need more expressiveness or are too complex for machine learning models. Ontologies like the Music Theory Ontology (MTO) [RDRM18] and the Music Score Ontology (MusicOWL) [JdSBTK17] attempt to model music-theoretical concepts and notation, respectively. MTO is limited in expressiveness for polyphonic and multi-voice music, while MusicOWL focuses on semantic representation and only partially captures engraving rules and visual formatting. The Music Note Ontology [PG21] focuses on high-level concepts and musical notes but falls short of fully integrating visual score characteristics. While these ontologies provide a rich semantic framework, their complexity and focus on semantics over visual layout limit them for use in machine learning models that need to align and understand image data (scores) with symbolic data (notation).

*3) Transformed Symbolic Music Representations for Foundational Models:* Unlike audio, symbolic music files in the previous section often contain varied informational content and do not provide a straightforward way of processing through deep-learning-based models. In this section, we introduce different mid-level representations that have been applied in literature (a visualisation can be found in Figure 3), with an emphasis on foundation models.

**Tokenised sequence representation:** Tokenised sequence is currently the most popular way of inputting symbolic music. It can be extracted from MIDI or MusicXML, designed with flexible encodings and emphasise specific musical information. Early designs like [DHYY18] encode music as sequences comprising pitch, duration, chord, and bar position, with
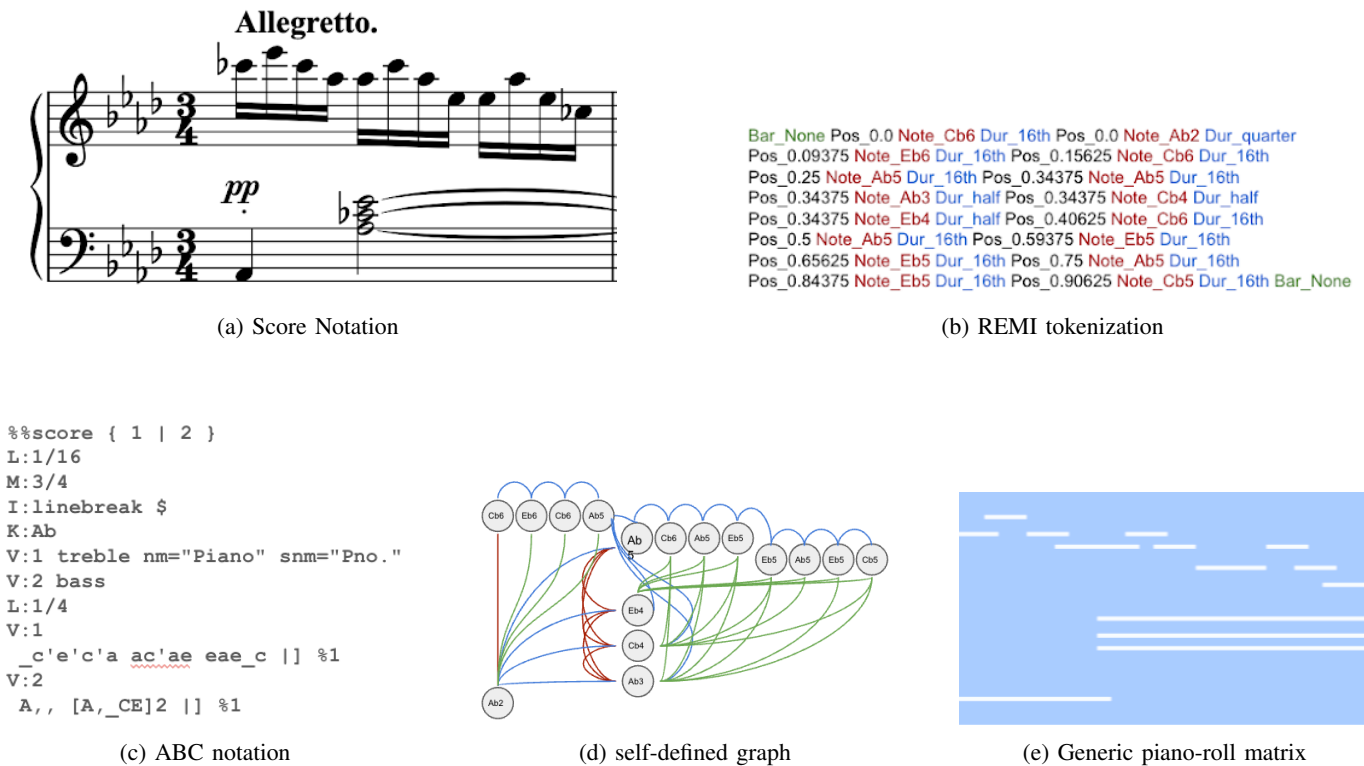
(a) Score Notation

(b) REMI tokenization



(c) ABC notation

(d) self-defined graph

(e) Generic piano-roll matrix

Figure 3: Excerpt of Schubert's *Impromptu Op. 90 No.4* and its input visualisations

chords represented by a 12-dimensional binary vector and bar positions indicating a note's relative location within a bar. Other approaches [DVVD20] include representing a lead sheet with sequences of pitch, rhythm, and chord, emphasising monophonic texture, or using dual sequences for pitch and rhythm with unique symbols for note continuation and rests. DeepBach [HPN17] distinguishes itself by using multiple lists to denote Bach's chorales across four parts, incorporating symbols for pitch maintenance and quantising time to accommodate rhythm and fermatas. Anticipation-RNN [HN17] adapts by using real note names or operationalising scores into time series to simplify the learning of polyphonic music with complex rhythms. Some models [JJDZ20], [CWBKD20] introduce special symbols or encoding strategies for note continuation, varying from shared "hold" symbols to part-specific indications, thereby balancing computational efficiency with the need to capture detailed musical structures.

**Note tuples** is a subcategory of note sequences that groups tokens in a more structural way. For instance, in the BachBot project [LGJS17], each note in a Bach chorale is encoded as a tuple containing pitch and tie information, with each time frame comprising four tuples to represent the four voices. This method organises notes by descending pitch and uses special delimiters and symbols to indicate frame separations, fermatas, and the start and end of scores. Following a similar methodology, BandNet [ZCYC19] employs a Z-scan for score scanning. Furthermore, NoteTuple in [HHIE18] includes additional attributes like velocity, duration, and time offset for piano performance modelling, offering a succinct representation that avoids the interleaving of notes. This tuple-based

representation, also used in studies like BachProp [CG18] and PianoTree VAE [WZZ+20a], simplifies musical information encoding, making it a preferred choice for many researchers.

More recent MIDI tokenization designs emphasise capturing structures and reducing sequence length: The REMI (REvamped MIDI-derived events) tokenization [HY20a] enhances MIDI representation by focusing on rhythm modeling and structural clarity. REMI introduces Note Duration events to directly represent rhythms, replacing traditional Noteoff events. It uses a combination of Bar and Position markers to delineate musical structure, providing a clear metric grid that reflects the actual layout of music, which is employed in works like [GSK+22] and [CCC+24]. The Compound Word [HLYY21a] employs embedding pooling to shorten the sequence length: note tokens (Pitch, Velocity, and Duration) are independently converted to embeddings, then merged into a single one. Similarly, Octuple [ZTW+21] uses embedding pooling to make sure each pooled embedding represents a single note. MuMIDI [RHT+20a] is designed for multitrack tasks. It contains a *Track* token that precedes note tokens and includes a "built-in" and learned positional encoding. The python package MIDITok [FBC+21] provides a unified implementation for these methods.

**ABC tokens**: ABC format, as introduced in the previous section as a file type, is by itself a natural sequential representation that's compatible with language models without additional conversion steps. ChatMusician [YLW+24a] is among the first to leverage the language models' ability to understand and generate symbolic music, with the text-compatible ABC music as input. [QBM+24] further developed a Synchronised

Multi-Track ABC Notation (SMT-ABC) that aims to preserve coherence across multiple musical tracks by concatenating elements from different tracks and then enclosing them by newly introduced symbols < | > indicating the barline. Large-scale cross-modal pre-training also demonstrated ABC notation's capability in semantic search and zero-shot classification [WYTS23a].

**Piano roll:** The piano roll serves as a historical symbolic representation of music, dating back to the era of player pianos. These self-playing instruments used piano rolls—continuous paper rolls with punched perforations—to automatically perform music. The perforations on the roll, representing note control data, trigger the playing of notes as they pass over a tracker bar. Player pianos were capable of capturing and reproducing the performances of renowned pianists, encoding not only the pitch and duration of notes but also the dynamics of the performance. In modern music technology, the piano roll has evolved into a geometric visualisation that is used for music analysis and generation. This representation plots time on the horizontal axis and pitch on the vertical axis, with each note depicted as an axis-parallel rectangle that encodes onset time, pitch, and duration. Such a two-dimensional representation is particularly compatible with diffusion-based models, and it has been applied to tasks like transcription [CSU+23] and generation [MJXZ23], [LQZ+24a].

**Notes graph:** Graphs emerge as a natural representation of symbolic music since music exhibits innate structures like voices and chords that can be formed into a graph with musical heuristics. [JKKN19] introduced a novel approach to representing music scores as graphs, where each note forms a node and various musical relationships between notes are depicted as edges. This graph-based representation utilises six primary types of edges—next, rest, set (onset), sustain, voice, and slur—to capture the intricate connections within the score. The model distinguishes forward and backward directions, along with a unique self-connection for each note, culminating in a total of 12 edge types for a comprehensive and detailed musical score representation. Similarly, [KFW23] also created a heterogeneous graph from score notes and tackles the voice separation problem as graph link prediction in multi-trajectory tracking. [ZKD+23] further examined the potential of applying score or performance graphs with various edge designs on music understanding problems and compared them with other representations counterparts.

**Tonnetz matrix** offers a distinctive two-dimensional approach to representing music in computational models, setting it apart from conventional piano roll representations. Inspired by music theory, it organises music into sequences of 2D matrices, where each node corresponds to one of the 12 pitch classes. These nodes are arranged to reflect harmonic relationships, with horizontal lines following the circle-of-fifths and triangles representing major and minor triads. The network's expandable nature allows pitch classes to appear multiple times, facilitating the depiction of complex harmonic structures. According to [CH18], the use of the Tonnetz in music generation leads to outputs with greater pitch stability and more repetitive patterns, showcasing its potential for innovative music representation and generation in computer models.

## C. Multimodal Music Representations

Music is a multidimensional artistic medium with various facets, each of which gives musical compositions a unique viewpoint. Apart from the auditory and symbolic aspects that are immediately associated with music, two important modalities in the field of multimodal learning for music have been extensively investigated:

**Text**: Words are essential to creating and comprehending music, including lyrics, metadata, and user comments. They provide a deeper understanding of the themes and storylines found in music by bridging the gap between musical expression and language context.

**Visual**: Visual components such as album art, frames from music videos, and associated visuals not only enhance the whole listening experience but also have an aesthetic impact.

The subsequent sections will delve into multimodal representations of music, specifically focusing on the interaction between music tokens and their representations in both textual and visual domains.

*1) Interactions with Textual Modality:* In this section, we will examine music-text representation learning. Traditional NLP methods, such as TF-IDF, will not be covered. Instead, we will briefly introduce contemporary deep learning approaches, including using pre-trained embedding from encoder models like the encoder of T5 [RSR+20a], using pre-trained language models as decoders like Llama2 [TMS+23], and training from scratch.

*a) Pre-trained Embeddings from Natural Language Encoders:* There are three types of pre-trained text encoders used as text conditioning in previous work: pre-trained encoders such as the encoder of T5, instructed-based encoders like FLAN-T5 [CHL+24], and pre-trained text-audio encoders like CLAP [WCZ+23a].

Audiogen [KSP+22] is an LLM conditioned on T5 embeddings to generate learned quantised audio tokens autoregressively. Noise2Music [HPW+23b], and Moûsai [SJS23a] generate high-quality music clips with a Latent Diffusion Model (LDM) conditioned on pre-trained T5 embeddings of text prompts. ERNIE-Music utilises ERNIE-M [OWP+21] to encode multi-lingual inputs to produce music waveforms directly from free-form text [ZPW+23] with a LDM architecture.

MusTango [MGG+23] leverages FLAN-T5, showcasing advancements in enhancing music generation controllability and employing a high-fidelity Latent Diffusion Model (LDM) for text-guided universal music generation [LCY+23a]. Similarly, MusicGen [CKG+23a] integrates FLAN-T5 with T5 and CLAP to generate conditioning text embeddings for an audio LLM.

MusicLDM [CWL+23b] is a diffusion model that has been proven powerful in text-conditioned music generation. It uses the pre-trained model CLAP during training to condition on audio embeddings. At inference time, CLAP is used to extract text embeddings to condition generation. Effective text-to-music generation is made possible by this method, which

guarantees varied, high-quality outputs that complement the original content. MusicLM [ADB+23] and Make-an-audio [HHY+23] also utilise CLAP embeddings as a text condition for music generation.

*b) Pre-trained Embeddings from Natural Language Decoders:* Significant progress has been achieved in LLM-integrated frameworks for multimodal music generation and understanding. MusiLingo [DML+24] combines Llama [TLI+23b] and a pre-trained acoustic music representation to perform music captioning and instruction following, such as generating music-related Q&A pairs. Similarly, LLark [GDSB23a] combines Llama2 [TMS+23] and multiple generative audio encoders, providing good performance in music instruction, specifically in zero-shot situations.

M$^2$UGen [HLSS23b] understands and produces music through the integration of images and videos, utilising the pre-trained Llama 2 model to achieve a comprehensive understanding of music. Using pre-trained audio representation models, Authors of [LHSS23a] improve text-to-music creation by adding music-related question responses and captioning.

In addition to the models mentioned previously, pre-trained natural language decoders such as GPT-4 can act as agents to make complex AI music tools more accessible to a larger audience by streamlining a range of music-related tasks. For more information, please refer to Section IV-E.

*c) Trained jointly from Scratch:* ChatMusician [YLW+24b] interprets music as a language by adding ABC notations to Llama vocabulary to train a tokeniser and a GPT model based on it. It is proficient at creating well-organised compositions from ABC notation through understanding subtle musical aspects.

SongComposer is trained based on a tokeniser as a tuple of symbolic song notes and lyrics [DLD+24a], which transforms textual descriptions into musical compositions, improving the coherence between melody and lyrics.

By training models from scratch, researchers can directly align the tokenisation and modelling processes with the unique structure and semantics of musical data.

*2) Interactions with Visual Modality:*

*a) Pre-trained Embedding from Visual Encoder:* Pre-trained visual models like the Vision Transformer (ViT) [DBK+20] and Video Vision Transformer (ViViT) [ADH+21] have revolutionised tasks that integrate visual and musical elements. Utilised by systems such as M$^2$UGen [HLSS23a], these models adapt NLP-focused transformer architectures to process image patches and spatio-temporal visual tokens.

Both V2Meow [SLH+23] and VidMuse [TLY+24] leverage the multimodal CLIP model [RKH+21b] as their visual encoders to effectively bridge the gap between visual and audio modalities. In V2Meow, the visual encoder extracts high-level features from video frames, which are then used to condition the auto-regressive music generation model. Similarly, VidMuse employs CLIP's visual encodings to capture both local and global visual cues. By integrating these cues through long- and short-term modelling, VidMuse produces music tracks that are not only acoustically rich but also semantically synchronised with the video content.

*b) Trained jointly from Scratch:* AnyGPT [ZDY+24] uses a unified approach to process multiple modalities, including speech, text, images, and music with specialised tokenisers. For example, music uses the EnCodec model [DCSA22] to convert audio into discrete tokens, while images use the SEED tokeniser [GGZ+23] for visual data. These tokens are processed by the language model and reconstructed using detokenisers, ensuring the perceptual integrity of the original modalities and enabling effective content generation across diverse inputs and outputs.

MM-Diffusion [RMY+23] introduces a joint audio-video generation framework that employs a sequential multimodal U-Net to generate aligned audio-video pairs from Gaussian noise. The integration of audio and video generation in a single framework ensures that the audio and visual elements are harmoniously aligned, resulting in more immersive and coherent multimedia content.

## III. APPLICATIONS

In this section, we will introduce the application of foundation models on music, including music understanding, music generation and musical application. Traditional approaches also have such kinds of applications but lack versatility. Foundation models (FMs) like the product of sunoAI and Q-wen-audio have or potentially have better and more generalised performance due to the self-supervised learning paradigm and a larger number of parameters.

### A. Music Understanding

Music information retrieval (MIR) refers to the field of research focusing on retrieving information from music data [Dow03]. The field's origins stem from the wider discipline of information retrieval (IR) and information science, although its scope and methodologies have shifted such as the community can be said to more broadly focus on Music Information *Research* [SMB+13]. The focal point of the MIR community is the ISMIR[7] society and conference (standing for *International Society for Music Information Retrieval*), first established in 2000 as an international symposium before converting into a conference in 2002. The original scope of the research community during the inception of the conference was on retrieving information from symbolic music representations and managing bibliographic music collections (with a large part of stakeholders being music libraries), although over the years audio signals have become the predominant modality used in MIR research. The term MIR is sometimes used interchangeably with *Music Informatics* or *Music Information Processing* [Mül15], and there are also strong links between MIR and the signal processing sub-field of *Music Signal Processing* (showcased in conferences such as IEEE ICASSP[8]) and more recently with the sub-field of AI focusing on *Music AI* [Mir21].

To the authors' knowledge, there is no systematic and commonly agreed taxonomy of the different problems and

---

[7]https://ismir.net/
[8]https://ieeexplore.ieee.org/xpl/conhome/1000002/all-proceedings

tasks that have been of interest to the MIR community over the years, which could subsequently form the basis for developing music foundation models. Certain task categorisations and groupings have been proposed in various MIR survey papers and textbooks [SMB+13], [Mül15], [SGU+14], and a fairly comprehensive list of popular MIR tasks can be found as part of the MIREX public evaluation challenge[9] that ran annually between 2005 and 2021. The below paragraphs will provide a brief overview of common MIR tasks organised based on different musical dimensions (e.g. rhythm, harmony, timbre), noting that the majority of MIR research activity uses audio as its primary modality. Increasingly, the MIR community also includes tasks beyond analysis, which include music generation or manipulation; these will be covered in the subsequent subsection III-B; in addition, tasks related specifically to multimodal MIR and singing voice will also be presented in a separate subsection III-A2 and III-A3 respectively.

*1) Traditional Music Information Retrieval Tasks:* While coming up with a single MIR taxonomy is challenging and out of scope for this work, one suggested task categorisation could be as follows:

- Tasks related to tonality and harmony: mode, chord, and key detection
- Tasks related to melody and pitch: melody estimation, pitch & multi-pitch detection, note tracking, automatic music transcription
- Tasks related to rhythm: onset detection, beat & downbeat tracking, metre estimation, tempo estimation
- Temporal alignment tasks: score following, audio-to-score alignment, score alignment
- Source separation tasks: musical instrument source separation, harmonic-percussive source separation
- Timbre-related tasks: musical instrument identification, playing technique detection
- Clip-level classification tasks: music tagging, genre recognition, emotion/mood recognition
- Content-based audio retrieval: audio identification, audio matching, cover song detection
- Temporal segmentation tasks: music detection, music structure segmentation, time boundary identification
- Tasks with visual score input: optical music recognition and subtasks, including staff line identification, music symbol identification
- Performance-related understanding: technique identification, performer identification, performance assessment, difficulty estimation

Given the extensive task list above, there is no single unified resource for providing task definitions and benchmark methodologies, although there exist textbooks that cover a subset of the above tasks (e.g. [Mül15], [Ler22]), with ISMIR conference proceedings[10] being useful for providing additional technical details. Review and overview papers do exist for specific tasks (which might or might not be up to date with related literature depending on the publication year), including on onset detection [BDA+05], music transcription [BDDE19], chord

estimation [MSRNDB14], melody extraction [SGER14], score following [OLS03], music source separation [CFL+18], music classification [FLTZ10], audio identification [CBKH05], cover song detection [SGH10], optical music recognition [CZJP20], plus general roadmaps for the community [SMB+13].

Attempting to summarise methodological trends, most MIR tasks above focus on audio as the input modality, with a smaller subset of tasks using machine-readable scores or visual representations as input. The vast majority of MIR methods focus on Western tonal music, with a much smaller subset of research focusing on non-Western cultures or folk/traditional music. The varying temporal granularity between different tasks is a key differentiator, with MIR ranging from clip-level classification tasks such as audio tagging (or even higher level tasks such as setlist identification) to tasks which require prediction of musical cues at a fine temporal level (e.g. pitch detection, onset detection). These varying task requirements in terms of temporal granularity, combined with the diversity of music cultures in a global scale, might potentially prove an obstacle for the creation of a general-purpose music foundation model, as will be discussed in the following sections. On the other hand, the clear connections that exist between different tasks (e.g. onset detection is directly linked with beat tracking) can lead to the creation of music representations and MIR models that can potentially address multiple tasks, possibly at the expense of musical diversity.

More recently, large music audio models based on pre-training and self-supervised learning have been proposed to tackle a range of MIR tasks, therefore leading to the emergence of foundation models for music. Technical details of those models will be presented in Section IV, and a brief summary for completeness will be provided here. In what might be one of the first attempts towards establishing a common testbed for multiple MIR tasks (thus aiming towards the creation of music understanding foundation models), MARBLE [YML+24] introduced a benchmark with a hierarchical taxonomy for 18 MIR tasks across acoustic, performance, score, and high-level description levels.

Specific models that might be classed as foundation models for music or large music audio pre-training models include JukeMIR [CDL21], which explores learnt representations from Jukebox [DJP+20a] for four MIR tasks; MULE [MKO+22], which is a self-supervised model that combines a SlowFast component with a variant of the well-known ResNet architecture, pre-trained on the MusicNet dataset and applied on music classification tasks; Music2Vec [LYZ+22], which is a self-supervised model using a masked prediction strategy with student and teacher models and applied to a variety of MIR tasks not limited to classification; MERT [LYZ+24] which combines two teacher models using a masked language modelling acoustic pre-training approach, applied to 14 MIR tasks; and finally MusicFM [WHL24a], which also applies a self-supervised learning approach taking advantage of random quantisation and applied to 5 MIR tasks. Fig. 4 shows results of music audio pre-trained models across a range of MIR tasks, as reported through the MARBLE benchmark [YML+24].

*2) Multimodal Music Understanding Tasks:* multimodal information is abundant in the real world, where agents un-

---

[9]https://www.music-ir.org/mirex/wiki/MIREX_HOME
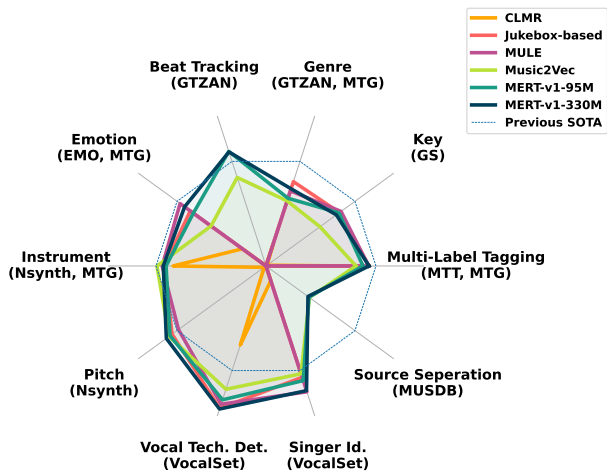[10]https://ismir.net/conferences/

Figure 4: Comparison of various music audio self-supervised models evaluated on a range of different tasks, as reported through the MARBLE benchmark [YML+24]. Figure reprinted with permission.

derstand and interact through various channels such as vision, language, speech, and touch. Therefore, AI systems need to effectively process and integrate data from these diverse sources to improve their ability to comprehend and engage with their environment. This section reviews extensive research on multimodal music understanding, which connects the domain of music with other modalities like vision and language. It highlights key advancements and methodologies from recent studies, showcasing how these developments enhance our understanding and interaction with music through multiple forms of information. Subsequently, we organise these methodologies into: 1. Audio-Visual Joint Representation Learning and Synthesis; 2. Audio-Language Joint Representation Learning and Synthesis; 3. multimodal Extensions of Large Language Models. For the multimodal application to music generation, typically text2music generation, please refer to sectionIII-B.

*a) Audio-Visual Joint Representation Learning:* **Music source localisation in visual context.** PixelPlayer[ZGR+18] learns to isolate musical sources and spatially localise them within the visual context from unlabelled video sequences. Through a joint audio-visual learning approach, PixelPlayer can identify objects and their sounds, localise objects in images or video, and separate audio elements emitted by each object.

**Video music retrieval & recommendation** [LK19] deploys a dual-stream network for music retrieval from videos. This network learns to match music and videos by understanding their cross-modal distances, and it is trained end-to-end with music-video pairs. Meanwhile, it creates an emotion-based latent space, pre-trained through audio and video emotion tagging. [MAB+19] addresses the requirement for cross-modal retrieval in music, allowing users to query one modality (audio, for example) and retrieve related information from other modalities (sheet music, album covers, etc.) . VM-NET employs a specialised loss function to pair videos with music, preserving the unique characteristics of each modal-

ity [HIY18]. Improving upon this, a study enhances retrieval quality by shifting from manual feature extraction to learned feature representations [PRP21]. MRCMV, a framework for selecting background music for detailed virtual-content videos, leverages SlowFast and VGGish networks to extract video and audio features. It further employs self-attention and cross-attention modules to understand intra- and inter-modal relationships, with a fusion gate determining the contribution of each feature type [LSZ+21]. [SVRS22] capture the extended temporal context within music and video modalities, offering a sophisticated approach to content pairing. [GSL23] utilises a dual-path cross-modal network to integrate content and emotional information to enhance music recommendations for videos. Addressing the challenges of label noise and the inadequate capture of key video segments, the Saliency-based Self-training for Video-Music Retrieval (SSVMR) framework has been proposed [CZL+23]. ViML takes a trimodal approach, recommending music by leveraging integrated representations from both video and text inputs, showcasing the potential of combining multiple modalities [MSSR23].

**Visual & audio classification** MAViL [HSX+24] combines Masked Autoencoders (MAE) and contrastive learning to develop a self-supervised learning method for handling the heterogeneous audio and video modalities. MAViL reconstructs heavily masked inputs, learning audio-video representations and aligning them in a unified latent space. Moreover, MAViL introduces a pre-training task that reconstructs unified, context-aware audio-video representations instead of simple single-modality input reconstruction.

*b) Audio-Language Joint Representation Learning:* **Music captioning**, which involves creating descriptive text for music, represents a unique blend of linguistic expression and musical perception. This emerging field aims to enhance music accessibility and understanding through natural language. The curation and augmentation of datasets are crucial in this field. To tackle the data scarcity issue in music captioning, [DCLN23a] proposes a method that employs LLMs to generate pseudo captions from tags. [MWD+23b] introduces the Song Describer Dataset (SDD), a collaborative collection featuring more than one thousand natural language descriptions for 706 musical recordings. To advance music captioning, novel model architectures have been developed. [MBQF21] propose a multimodal encoder that combines a CNN for audio and an LSTM for jointly representing audio and text. This method effectively captures high-level semantics and summarises information across varying levels of input granularity. [ZJXD22] integrates a CNN-SA audio encoder with a pre-trained language model to improve song lyrics interpretation. ALCAP [HHL+23] utilises an alignment module through cross-modal contrastive learning between music and lyrics, enhancing the latent representation of music and lyrics. [LDJ23] finds that subjective annotator insight is crucial for crafting deep learning models tailored to music captioning. Based on a multi-layer cross-attention mechanism, [DYL+24] employs a joint attention model to improve the interaction between musical and textual information across various semantic layers.

**Music instruction following** Mu-llama [LHSS23a] and

MusiLingo [DML+24] integrated MERT[LYZ+24], a pre-trained music encoder, and Llama, a pre-trained natural language decoder. Mu-llama employs a Llama-adapter on the upper transformer layers for modality alignment [ZHL+23a] while MusiLingo utilises a simple linear projection on the initial transformer layer, yielding enhanced performance. Given the diverse range of tasks within the Music Understanding category, recent works have explored the application of Multimodal Instruction Tuning to enhance the performance of Music Language Models. LLark [GDSB23b] and JMLA [DYL+24] investigate how this approach can be effectively applied to the training of Music LLMs, aiming to improve their versatility and effectiveness across a wide spectrum of music-related tasks. Both models employ a similar architectural approach, utilising cross-attention modules to connect audio encoders with pre-trained LLMs. This integration allows for the effective processing of both audio and textual information. To create robust instruction-tuning datasets, these studies leverage traditional music tagging datasets, employing data augmentation techniques to expand them into comprehensive datasets suitable for fine-tuning Multimodal LLMs. By incorporating this approach, these models demonstrate the potential for more adaptive and comprehensive music understanding capabilities, effectively bridging the gap between audio processing and natural language understanding in the music domain.

**Text-music information retrieval** is an evolving field that integrates text and visual data to enhance access to and understanding of music. This interdisciplinary approach enriches interactions between sensory modalities, allowing for more nuanced music analysis and retrieval. MuLaP proposes a multimodal pre-training framework for MIR that integrates audio and text features in an early fusion manner [MBQF22b]. In contrast, other researchers have adopted a dual-tower architecture focusing on contrastive learning. This method aligns textual descriptions with their corresponding musical representations by leveraging a large corpus of music-text pairs for pre-training [WYTS23b], [MBQF22a], [HJL+22]. This dual-tower approach allows for distinct processing of the modalities before their interaction, potentially leading to more effective cross-modal alignment.

**Text-music understanding & generation** AudioGPT [HLY+23] is a multimodal model that processes and generates audio from text and audio inputs. It follows four steps: It converts speech to text, interprets the task to understand user intent, uses GPT to match tasks with audio models based on features like prosody and timbre, and then the audio models generate a response for the user. AudioPaLM [RAN+23] uses a dual-stream architecture with text [ADF+23] and speech [BMV+23] language models. It integrates speech and text into a multimodal generative model and creates a joint vocabulary, enabling the model to process and generate both speech and text.

*c) Audio-Visual-Language Joint Representation Learning:* Several studies leverage Large Language Models (LLMs) like GPT and LLaMA for multimodal applications, extending their capabilities to perform extensive music understanding tasks.

AnyGPT [ZDY+24] converts multimodal input into seman-

tic tokens which are capable of being subsequently converted back to their original forms by de-tokenisers. This method preserves the perceptual integrity of the original modalities while enabling the LLM to handle understanding and generation tasks autoregressively. M$^2$UGen [HLSS23a] is a framework designed for multimodal music understanding and generation. M$^2$UGen employs different modality encoders to interpret various inputs: the Vision Transformer (ViT) [DBK+20] for images, ViViT [ADH+21] for videos, and the MERT [LYZ+24] model for music. After processing the inputs into feature representations, modality-specific adaptors refine them. Subsequently, the LLaMA2 [TLI+23b] model interprets these modality signals for various downstream tasks. Specifically, M$^2$UGen then uses AudioLDM2 [LTY+23a] for music generation. NExT-GPT [WFQ+23], a model constructed on the M$^2$UGen framework, can process and generate content across text, images, videos, and audio. This model employs specific encoders and projection layers to convert inputs from diverse modalities into a unified representation space. It utilises a large language model Vicuna [CLL+23] for interpreting these multimodal inputs. NExT-GPT utilises pre-trained diffusion models, such as AudioLDM [LCY+23c], for audio generation.

*3) Vocal Music Understanding:* Vocal music understanding is similar to speech understanding and/or speech processing, and many techniques for vocal music or singing voice are inspired by the algorithms for speech processing, such as speaker identification for singer identification and speech recognition for lyrics transcription. The algorithms for speech are typically optimised for more uniform pitch, rhythm, and dynamic articulation. Different from speech signals, the singing voice has a wider pitch range, greater variability in rhythm and pronunciation, and involves complex vocal techniques. Therefore, keeping the differences in mind is essential when applying the algorithms to singing voice.

In the past few years, SSL models for speech processing have proven effective in various vocal understanding tasks. There has also been an exploration into training SSL features specifically for the singing voice. This section introduces representative work in singer identification, lyrics transcription, singing transcription, vocal source separation, and lyrics interpretation.

*a) Singer Identification (Singer ID):* This task involves recognising the identity of a singer[11]. It is often formulated as a classification problem within a dataset containing a fixed number of singers. Many approaches extract Mel-Frequency Cepstral Coefficients (MFCCs) to capture voice characteristics, paired with Gaussian Mixture Models (GMMs) [Zha03], [CLG11]. Embeddings designed for speaker recognition and verification, such as i-vector [PART13] and x-vector [SGS+18], have also been used in singer ID [Kru14], [ZWCX22]. With the advancement of neural networks, it has become possible to extract latent features from general audio representations, such as spectrograms [LN19], [ZQY+21], using CNNs and CRNNs. Recent methods propose using singing source separation to suppress background music [SDL19] or data augmentation [HCF+20]. These approaches have

---

[11]We do not differentiate it from similar tasks such as artist classification

improved accuracy and F1 scores on the Artist20 [Ell07]. However, these methods have been tested on a small, fixed number of singers, and their ability to scale up and generalise to unseen singers remains unknown.

There have been several attempts to approach singer ID through representation learning. [LN19] proposed a triplet network that learns a joint embedding for monophonic and mixture singing. [ZQY⁺21] transformed the classification problem into a query searching task by replacing the softmax layer with a k-nearest neighbours (KNN) layer. [TLR23] trained SSL features on source-separated singing voices and compared them with speech SSL features in terms of classification accuracy.

*b) Automatic Lyrics Transcription (ALT) & alignment::* These are two closely related tasks aimed at recognising and locating lyrics within singing. Lyrics transcription focuses on identifying the linguistic content of the singing voice, while lyrics alignment involves retrieving the timestamps of each text unit of lyrics, typically words. Both depend on an acoustic model that captures the relationship between audio signals and phonetic information. Previous work has utilised the same training pipelines for both tasks to develop an effective acoustic model.

Early approaches to automatic lyrics transcription and alignment utilised a Gaussian Mixture Model-Hidden Markov Model (GMM-HMM) design or a factorised time-delay neural network model, often combined with a language model [DB19]. Building on the traditional modular system, some researchers proposed enhancements in architecture [DAD20] and multi-domain training [DAD21]. [SDE19] introduced end-to-end training using the connectionist temporal classification (CTC) loss function. Subsequent works explored various aspects of the task using end-to-end methods, including multilingual application [VHM⁺20], [HB24], training domain [XGL22], joint training [GGL23], [HBE22], etc.

Given the great success of pre-trained self-supervised learning (SSL) features in speech recognition, researchers have begun applying these features to automatic lyrics transcription (ALT). [OGW22] were the first to finetune wav2vec 2.0 features on source-separated singing, achieving significant improvements. Following this, [GYL23] applied self-training with noisy student augmentation. Whisper, a robust speech recognition model released by OpenAI [RKX⁺23a], has also proven effective for singing voices. Efforts have been made to create datasets using Whisper [ZYP⁺23] and to adapt it for low-resource language ALT [WLL⁺23].

[DSE23] trained a similarity model between audio and text using contrastive loss and applied it to lyrics alignment. There have also been attempts to apply SSL features trained on music signals to ALT, such as in the MARBLE benchmark [YML⁺23], although these efforts have had limited success.

*c) Singing Transcription: :* Singing transcription aims at estimating the musical notes of a sung melody. Vocal F0 estimation, or melody extraction, is a closely related task that focuses on extracting the fundamental frequency of the singing voice and often plays a role in typical cascading systems for singing transcription. Early approaches for singing transcription consist of an F0 estimation module

and a quantisation module with segmentation to obtain note-level transcription. In [MTBB15], [KG16], the singing voice is first processed through pitch extraction algorithms, followed by note segmentation and quantisation modules that generate symbolic notations. To avoid error propagation by modules, end-to-end approaches using encoder-decoder models have been proposed [NNGY19], [NNF⁺19]. [WJ23] proposed using both Connectionist Temporal Classification (CTC) loss and cross-entropy loss to facilitate learning from weakly labelled data. To leverage information such as beat and phoneme, these elements are incorporated either through multi-task learning [NNGY19], [GOZ⁺24] or conditioning mechanisms [YSN23]. Inspired by object detection and sound event detection work, MusicYOLO [WTY⁺23] retrieves note objects from a macro perspective through an object detection model.

Different training methods for singing transcription have also been investigated. [NNF⁺19] proposed a loss function on the attention weights using semi-supervised alignment labels. [KLK⁺22] introduced self-training with the noisy student framework [XLHL20], while [HS21] utilised virtual adversarial training [MMKI19], both leveraging unlabelled singing data.

*d) Vocal Source Separation::* The singing voice is one of the four primary stems targeted in source separation, with applications such as karaoke and as a crucial component in various vocal understanding tasks. Therefore, researchers have focused on improving singing voice separation. To enhance the performance of vocal source separation (VSS) and related singing analysis tasks, multi-task learning approaches have been proposed. These include joint vocal activity detection [SED18b], pitch estimation [JBEW19], and lyrics transcription [GGL23]. Additionally, some methods incorporate lyrics information into the model to improve intelligibility [MP20], [JCL20]. The literature has also explored generative models for VSS, including approaches using GANs [FLJ18] and normalising flows [ZDJ⁺22]. Furthermore, semi-supervised training frameworks have been proposed to enhance SVS performance [SED18a] and [WGI⁺21].

Several previous works have focused on training self-supervised or unsupervised learning features for VSS. [MDS20] and [IMDS20] conducted representation learning via a denoising objective using autoencoders. [YWW⁺23] introduced deep unfolding to learn a latent representation of vocal signals for VSS. Additionally, MERT [LYZ⁺24], a general-purpose representation for acoustic music understanding, has been successfully applied to source separation.

Inspired by denoising diffusion models, [PMS22] have proposed unconditional signal modelling to gradually convert a mixture into a singing voice or accompaniment. [PMSS23] introduced the use of a cold diffusion process as a feature followed by unsupervised clustering. More challenging tasks, such as the separation of multiple singing voices, have also been investigated [PCC⁺20], [YPRF23].

Recently, methods based on Differentiable Digital Signal Processing (DDSP) are also being proposed. Unlike the above methods, DDSP-based models can easily be trained on unlabelled music data and adapted to separate homogeneous sources like lead and backing voices. [SRK⁺23] proposed

a unified framework that can be used for both music stems and homogeneous mixtures via cascaded pitch extractors and deep neural networks. [RCT24] extends this method on leading and backing voice separation to a fully differentiable one via HCQT and voice assignment,

*e) Lyrics Interpretation::* Lyrics interpretation involves the text summarisation of song lyrics, aiding the quick understanding of the singing voice and the management of music based on its linguistic content. This task can be approached using text summarisation methods applied solely to the text [SS18], [FCGG19]. Incorporating audio provides an additional dimension, enhancing the understanding of the song. [ZJXD22] proposed a multimodal generative model with representation fusion for lyrics interpretation, leveraging both text and audio information.

## B. Music Generation

*1) Symbolic Music Generation:* In the area of symbolic music generation, whether it involves generating score-level symbolic music (e.g., ABC notation) or performance-level music (e.g., MIDI), the approach generally falls into two categories: generating from scratch or based on given conditions like chords, tracks, textual descriptions, or other musical properties.

*a) From scratch.:* Some early works explored **monophonic/melody music generation** from scratch as well as polyphonic music. In early research stage, using RNNs to generate melodies was the most common approach. Folk-RNN [SSBTK16] trained LSTM networks on 23k music pieces represented in ABC notation in two approaches: one predicts characters based on the previous 50 characters, and the other predicts tokens based on all previous tokens of a transcription. Anticipation-RNN [HN17] generates melodies through a RNN-based generative model while allowing user-defined positional constraints. Besides RNNs, some works utilise GAN-based models, such as SeqGAN [YZWY17], and VAE-based models, such as the hierarchical model proposed by Roberts et al. [RER+18a], to generate music melodies.

For **polyphonic music generation**, MuseGAN [DHYY18] was the first to propose models to generate multi-track symbolic music using the jamming model, the composer model, and the hybrid model. This method generates music with five tracks—bass, drums, guitar, piano, and strings in MIDI format. To generate more natural multi-track symbolic music, DMB-GAN [GYY19] employs a self-attention mechanism to extract both spatial and temporal features, building dual GANs for each branch to create a more harmonious structure.

With the advent of recent generative models, more advanced methods have been researched for generating symbolic music. Multitrack Music Transformer (MMT) [DCD+23] introduced a Transformer-based multitrack music representation that accommodates a diverse array of instruments while keeping the sequence length short, which aims to effectively address the limitations on the number of instruments supported by previous models. Mittal et al.[MEHS21] proposed a method to train diffusion models using symbolic music data, leveraging a pre-trained VAE to map the discrete domain to the continuous

latent domain. Building on the Transformer model, Muse-former [YLW+22] introduced a fine-coarse-grained attention mechanism to handle the challenge of long music sequences. In this approach, fine-grained attention captures structural information, while coarse-grained attention gathers contextual information, improving the model's ability to generate coherent and complex musical pieces.

*b) Based on condition.:* Generating music from scratch can often result in music that follows a certain pattern, limiting the creativity of music generation. By incorporating various types of input conditions, such as chords, melody tracks, lyrics, and text descriptions, not only can users interact with the music generation process more dynamically, but they also gain higher and more fine-grained control over the output. Consequently, there is a growing group of research focused on generating symbolic music under different input conditions.

**Conditioned on chord sequences**, MIDINet [YCY17] utilises a generative adversarial network to create music with multiple MIDI channels. This network is capable of generating music by adhering to a specified chord sequence or by conditioning on the melody established in preceding bars. Choi et al. [CPH+21] introduced a Transformer model conditioned on chords for generating K-POP melodies. This model generates rhythm and pitch components while adhering to a specified chord progression. MelodyDiffusion [LS23a] introduces a transformer-based diffusion model specifically designed for chord-conditioned melody generation with discrete musical data. Unlike traditional U-nets, this model leverages transformers to capture long-range dependencies via attention mechanisms and parallel processing. Besides, DeepChoir [WLS23] can generate a four-part chorale based on given melody and chord progression.

**Melody harmonisation** refers to crafting a chord progression that complements a given melody. This progression must harmonise with the melody while also aligning with its rhythmic pattern. Lim et al. [LRL17] propose a technique for generating chord sequences from symbolic melodies by leveraging bidirectional long short-term memory (BLSTM) networks. These networks are trained on a database of lead sheets to achieve the desired chord generation. AutoHarmoniser [WYW+24] is a system designed for melody harmonisation with controllable harmonic density and rhythm. It features an extensive vocabulary of 1,462 chord types, enabling it to generate chord progressions with varying harmonic density for a given melody. Not limited to melody, Multi-Track Music Machine (MMM) [EP20], which is based on the Transformer model, treats each musical note as a time-ordered sequence. In this method, notes from different tracks are interleaved into a single sequence. GETMusic [LTL+23b] employs a different strategy by representing musical notes as tokens and arranging them in a 2D grid, with tracks stacked vertically and time progressing horizontally. This diffusion-based model randomly designates each track of a music piece as either the target or the source during training.

**Generating symbolic music from text descriptions** with the advancement of pre-trained models and LLMs in recent years has been more and more popular. Wu et al.[WS22] initiated the first investigation into generating complete and

semantically coherent symbolic music scores from textual descriptions. They explore the efficacy of leveraging publicly available NLP checkpoints for the task of text-to-music generation. Music Composition Copilot (Musecoco)[LXK+23a] is a two-stage framework for music generation. In the first stage, ChatGPT synthesises and refines the input text into musical attributes such as instrument, time signature, tempo and pitch, etc. In the second stage, these attributes are used to generate the symbolic music.

ChatMusician [YLW+24a] proposes a novel approach that incorporates music as a secondary language in Large Language Models (LLMs). By utilising ABC notation, this method effectively merges music and text, allowing for internal music composition and analysis without the dependency on external multimodal frameworks. In contrast to ChatMusician, SongComposer [DLD+24b] employs MIDI for representing symbolic music and introduces a unique tuple structure. This structure formats lyrics alongside three specific note attributes: pitch, duration, and rest duration, ensuring accurate musical symbol interpretation and precise alignment of lyrics with the melody. MuPT [QBM+24], a generative pre-trained transformer for symbolic music, design the Synchronised Multi-Track ABC Notation (SMT-ABC Notation) to maintain coherence across multiple tracks. This allows the model to handle up sequences containing up to 8192 tokens, enhancing its capability to generate complex musical pieces. Apart from directly generating music from text descriptions, there are also efforts focused on using LLM agents to create symbolic music. Examples include Musicagent [YSL+23], ComposerX [DYY+24], and ByteComposer [LLD24].

**Music generation conditioned on video.** Music, integral to video production, has sparked researchers' interest in generating music that is conditioned on video content.

CMT [DJL+21], Video2Music [KPH23], and Video background music generation [ZWW+23a] all focus on creating background music with general video content. CMT initially establishes rhythmic connections between video and background music, offering local manipulation of these rhythmic elements and comprehensive control over music genre and instrument choices. Video2Music [KPH23] and Video background music generation [ZWW+23a] extract both low-level and high-level semantic features for generating music. Diff-BGM [LQZ+24b] generate background music using a diffusion-based method.

Different from the broader scope of music generation for general video content, some works specifically focus on generating music tailored for dance videos. Research explores translating human motions from dance videos into musical notation [SLS21], enhancing rhythmic sound generation. Multi-Instrumentalist Net innovates by generating instrumental music directly from videos of musicians without relying on supervised learning [SLS20b]. Similarly, Dance2Music [AP21] utilises dance movements as a foundation for generating music [AP21], illustrating the potential of dance-centric approaches in music creation. Efforts extend to synthesising music from silent videos of musical performances, converting body movements into MIDI sequences for realistic music synthesis [GHC+20a] and generating music from silent piano performances [SLS20a]. Advances include deep learning frameworks for transcribing piano music from visual data [KWMZ20]. To generate background music corresponding with the body movements of musicians in video clips, Foley Music [GHC+20b] utilises a Graph-Transformer architecture, which includes a Graph Convolutional Network (GCN) encoder and a Transformer decoder. It learns to map the relationship between human body key points detected in videos and MIDI events.

*2) Language Model for Acoustic Music Generation:* The state-of-the-art (SOTA) generative models various fields are predominantly transformer-based language models (LMs) [VSP+17a]. These LMs have revolutionised natural language processing [RNS+18], [BMR+20], [TDFH+22], [TLI+23b], demonstrating extraordinary capability in capturing contextual dependencies and relationships within sequences efficiently. In light of this, the application of LMs has extended beyond natural language processing to various tasks, including image generation [ERO21], [YXK+22], [YLK+22] ,video generation [GHY+22], [YLG+23]. Furthermore, this breakthrough in capturing sequence relationships has been extensively studied in speech and audio generation and achieves competitive performance.

**Application in speech synthesis.** GSLM [LKH+21] and pGSLM [KLP+22] proposed to quantise speech representations derived from the self-supervised learning models for speech (such as CPC [vdOLV18], wav2vec 2.0 [BZMA20b] and HuBERT [HBT+21b]) and employed a generative language model for both conditional and unconditional speech generation. Similarly, approaches that combine LMs with discrete representations have also been widely applied in speech resynthesis [PAC+21], speech emotion conversion [KPC+22], spoken dialog system [NKC+23] and speech-to-speech translation [LCW+22], [PCW+22]. Recently, Encodec [DCSA22], SoundStream [ZLO+21], and their derivative works [KSL+23], [YLH+23] proposed representing audio signals as multiple streams of discrete tokens by utilising residual vector quantizer (RVQ). This approach allows for the generation of high-quality audio from quantised tokens. Subsequently, several studies combined such discrete representation of audio with LMs for text-to-speech synthesis [WCW+23], [ZZW+23], [KVB+23], [BSV+23].

**Unconditional music generation.** To achieve music generation, ADAS [DvdOS18] has already attempted to use hierarchical VQ-VAEs [VDOV+17a], [RVdOV19] to learn discrete representations of music samples at various temporal resolutions before RVQ was proposed. They combined these with LMs to generate music with high temporal coherence, but the audio quality remains limited. With the advent of RVQ, Perceiver AR [HJC+22] sought to model the flattened multi-stream discrete token sequences of SoundStream [ZLO+21] directly with LMs, achieving high-quality piano music generation. However, due to the excessive length of the flattened sequences, these methods face challenges in maintaining long-term temporal coherence. AudioLM [BMV+23] addressed this issue by using semantic, coarse acoustic and fine acoustic tokens to represent audio and employing multiple cascaded LMs to generate different token sequences. This allows AudioLM

to generate coherent, high-quality piano music uncondition-ally. Besides the GPT style models, DrumGAN [NLR20] is a progressive generative adversarial network (GAN) model designed for the synthesis of drum sounds. It leverages con-ditional inputs based on perceptual features, enabling intuitive and musically relevant control over the generation process, thereby enhancing both the quality and the expressiveness of the synthesized drum sounds.

**Text to music.** Compared to unconditional music gen-eration, recent works have introduced descriptive text as a condition to achieve controllable music generation, also known as text-to-music generation. AudioGen [KSP+23] introduced textual information through T5 [RSR+20b] encoder, enabling text-to-audio generation. Inspired by CLAP [WCZ+23a], Mu-sicLM [HJL+22] further introduces descriptive text as a con-dition through MuLan [HJL+22], which is trained with con-trastive loss, and is capable of modelling a large variety of long music sequences beyond piano music SingSong [DCR+23] followed a similar modelling approach but for producing instrumental music that harmoniously aligns with the pro-vided vocals. Additionally, some works have explored visual-conditioned music generation. In generating background music for videos, models introduce controllable features [DJL+21] and explore mapping visual arts to music [ZZS+22]. A Transformer-based model aligns music with video content [KPH23], showcasing diverse applications. Innovative meth-ods focus on generating music audio with style control from silent videos [SLH+23] and complex music samples from dance videos [ZOW+22], highlighting the integration of visual content into music generation.

**Long high-quality music generation.** However, the cas-caded language models also present computational challenges. In the context of modelling the token sequences of neural audio codecs [DCSA22], [BSV+23], [KSL+23], [YLH+23], efficiently generating long, high-quality music segments re-mains an unresolved issue. SoundStorm [BSV+23] utilises a non-autoregressive decoding scheme [CZJ+22] to significantly accelerate the AudioLM [BMV+23], enabling the acoustic LM to complete decoding within 27 forward passes. Similarly, VAMPNET [GSKP23] uses the Descript Audio Codec (DAC) [KSL+23] as the audio tokeniser and incorporates parallel iterative decoding to achieve music generation. Parallel to this work, MusicGen [CKG+24] introduces efficient token interleaving patterns that eliminate the need for cascading multiple LMs, achieving the generation of high-quality music samples with a single-stage LM. Similarly focused on produc-ing high-quality music samples with a single-stage model, but in contrast to MusicGen which conditions on text descriptions, VidMuse [TLY+24] generates music based on video input.

**Lyrics to singing**. All aforementioned methods only take the instrumental music generation into consideration, generat-ing music with singing remains a significant challenge. To this end, Jukebox [DJP+20b] can be seen as the first and only attempt from published literature so far to simultaneously generate music with singing from lyrics using a single LM. Recently, while the industry has seen the emergence of song

generation tools like Suno[12] and Udio[13], neither has disclosed their methodologies. For more information, please refer to sectionIII-B4

**Music & general audio generation.** Furthermore, some research attempts to design universal audio generation methods rather than solely focusing on music generation. WavJourney [LZL+23] proposed to utilise Large Language Models (LLMs) to connect various audio generation models across different tasks, enabling the generation of comprehensive audio content (covering speech, music, and sound effects) from textual story narratives. Meanwhile, SpeechGPT [ZLZ+23] and UniAudio [YTT+23] strive to implement universal models based on textual instructions, utilising a single LM to execute various tasks with correct output on different instructions.

*3) Audio Diffusion Model:* In addition to the language model-based approaches, diffusion models [SDWMG15], [HJA20], [KSPH21], as a competitive class of generative models, have recently delivered impressive results in various domains of generative modelling. These models, through a series of diffusion and reverse diffusion processes, effectively formulate the mapping between data and latent distributions, enabling the generation of highly realistic and quality outputs. Diffusion models have demonstrated exceptional capabilities in tasks such as text-to-image generation [DN21], [RBL+21], [HSC+22], image super-resolution [SHC+22], [LYC+22] and image inpainting [SSDK+21], showcasing their ability to generate high-quality and diverse samples for images. Fur-thermore, advancements such as DDIM [SME21] and progres-sive distillation [SH22] have introduced sampling acceleration techniques that significantly reduce the long sampling time, a known drawback of traditional diffusion models, while maintaining the quality of generation. This improvement in efficiency makes diffusion models more practical for a wider range of applications.

**Application in speech & general audio.** The adaptability and versatility of diffusion models have impressively extended beyond image generation to a variety of modalities, including speech and audio. In speech, diffusion models are widely applied in vocoders [CZZ+21a], [KPH+21], [gLKS+22], [LWSY22], [HLW+22], to transform the mel-spectrograms into high-quality speech waveforms. WaveGrad [CZZ+21a] and DiffWave [KPH+21] based on the diffusion process to incrementally generate speech waveforms by gradually adding noise and reversing this noise. FastDiff [HLW+22] enhances the speed of diffusion by decreasing the number of sampling steps while maintaining high generation quality. Additionally, some works [JKC+21], [PVG+21], [CZZ+21b], [KKY22a], [KKY22b], [Ano24] also apply diffusion models to acoustic models. The most significant contribution of diffusion models in this field is the ability to generate high-fidelity speech, which is a critical advancement for realistic and naturalistic speech generation.

**Text-to-music generation.** This capability is particularly salient in audio generation and music generation, as these models excel in modeling complex temporal dynamics and

---

[12]https://suno.com/
[13]https://www.udio.com/

generating high-quality audio. Diffsound [YYW+23] proposed to quantise the mel spectrogram by VQ-VAE [VDOV+17a] and apply a discrete-diffusion-model-based non-autoregressive decoder [AJH+21], [GCB+22] instead of the traditional autoregressive (AR) decoder for generating quantised tokens from textual CLIP [RKH+21b] embeddings, effectively enhancing generative performance. However, the introduction of VQ-VAE, while beneficial for quantisation, has compromised the generation quality. AudioLDM [LCY+23c] and Make-An-Audio 1 & 2 [HHY+23], [HRH+23] have advanced text-to-audio generation by employing latent diffusion models (LDMs)[RBL+21]. By utilising a latent space derived from contrastive language-audio pre-training (CLAP) [WCZ+23a] to learn continuous audio representations, these models achieve superior generation quality and computational efficiency compared to Diffsound. Notably, both AudioLDM and Make-An-Audio 2 have extended their capabilities to generate text-conditional music by incorporating music data into training processes. TANGO [GMMP23a] employs an instruction-tuned LLM Flan-T5 [CHL+22] in place of CLAP, providing the model with advantages in the domain of audio and music generation from text input. Furthermore, AudioLDM 2 [LTY+23a] introduces a novel approach by using the self-supervised learning model AudioMAE [HXL+22a] to extract a general representation of audio, called "language of audio" (LOA), achieving unified generation of speech, music, and sound effects.

Compared to these works that generate audio, including some music, some works specifically focus on music generation, exploring generating long-term music with high-quality and high-fidelity. Riffusion [FM22b] represents an attempt to apply diffusion models in music generation, directly fine-tuning Stable Diffusion [RBL+21] on mel-spectrograms of music pieces from a paired text-music dataset to generate 5-second music clips. Building on this, MusicLDM [CWL+23a] further incorporates the framework of AudioLDM [LCY+23c] and refined the CLAP model [WCZ+23a] on music dataset, enabling the model to generate a more diverse and richer variety of music based on the provided text input. However, MusicLDM can only generate music with a 16 kHz sampling rate, while most standard music productions are 44.1 kHz, which falls short of the requirements for producing high-quality, high-fidelity music.

**Long high-quality music generation.** To address this limitation, Moûsai [SJS23b] not only utilises text-conditioned LDMs to learn and generate the reduced latent representations, but also employs a novel diffusion autoencoder to directly compress and generate audio. By cascading two LDMs, this model can handle the long-term structure of music and generate high-quality stereo music with a 48kHz sampling rate condition on a given textual description. Simultaneously, Noise2Music [HPW+23a] proposed an alternative approach to use cascaded diffusion models, employing multiple diffusion models to gradually increase the sampling rate the generated music, which show SOTA generative performance with high-fidelity. However, due to the employment of multiple cascaded diffusion models, the success of these two approaches incurs substantial computational costs, which would be a serious impediment to their practicalities.

MeLoDy [LTL+23a] introduces a novel solution by combining the advantages of language models and diffusion models. It employs a semantic LM based on MuLan [HJL+22] to capture the semantic structures of music, and proposes a dual-path diffusion (DPD) model to simultaneously model both coarse and fine acoustic information conditions on the semantic LM. This strategy enables the efficient generation of competitively high-quality music. JEN-1 [LCY+23b] introduces an omnidirectional diffusion model that integrates both autoregressive and non-autoregressive modes, enhancing sequence generation while improving sequence dependency. Furthermore, the model has been extended to music continuation and music inpainting.

**Visual-music generation.** Some studies have explored acoustic music generation conditioned on visual information. V2Meow [SLH+23] generates music through video input and style control via text input. M$^2$UGen [HLSS23a] employs LLMs as a bridge between visual and audio features, integrating a video module and a music generation module. Vid-Muse [TLY+24] uses a long-short-term modeling strategy to generate music conditioned on input videos autoregressively. MELFUSION [CNJ+24] introduces a model that leverages cues from textual descriptions and corresponding images to synthesise music. Another line of research in multimodal musical understanding leverages denoising diffusion models to learn the recovery of latent features from noise. DiffMAViL [NJR+23] integrates denoising diffusion processes with the MAViL [HSX+24] in a masked encoder-decoder framework to learn multimodal representations. It processes audio and video by masking pixels, adding noise, and using diffusion. This helps it pick up detailed features and rebuild the original audio-visual content. Through coupled architectures and bidirectional diffusion processes, models such as EasyGen [ZLL+24] and MM-Diffusion [RMY+23] combine audio and visual modalities, improving the quality of joint generation. With an emphasis on audio-video synchronisation, MM-Diffusion's unified architecture includes multimodal attention techniques and a connected U-Net [RFB15]. Easy-Gen uses BiDiffuser, which combines LLMs and diffusion models for efficiency in a variety of generative applications. Differently, Seeing and Hearing [XHT+24] build aligners across text, audio, and video modalities to jointly generate audio and video. Conditional Discrete Contrastive Diffusion (CDCD) [ZWO+23] utilises contrastive learning within diffusion frameworks to align music with specific images, enhancing cross-modal content generation and creative processes.

**Beyond text condition.** Beyond text-to-music generation, many works focus on other tasks of music generation based on diffusion models. AUDIT [WJT+23] and InstructME [HDS+23a] acheieved zero-shot music editing by feeding text instructions into LDMs. Music ControlNet [WDWB24] proposes a method similar to the ControlNet [ZRA23] approach in the image domain, to offers multiple precise, time-varying controls over the generated music.

*4) Vocal Music Generation:* The human voice is the most natural musical "instrument" possessed by humans. Singing voices, rich with emotive lyrics and diverse melodies, vividly

convey explicit human feelings and are an essential element of music. Vocal music generation focuses on producing singing voices under various conditions, broadly classified into the following key tasks: **Singing Voice Synthesis (SVS)**, which generates voices from input lyrics and music scores; and **Singing Voice Conversion (SVC)**, which adapts the singing voice from one person to match the timbre of another using a timbre prompt.

*a) Singing Voice Synthesis (SVS)::* SVS can be viewed as an extension of Text-to-Speech (TTS), incorporating not only text input (lyrics) but also musical notation to mediate the speech generation process.

Like TTS, the popular framework in SVS usually involves a two-stage process. Initially, an "acoustic model" converts lyrics and musical notation into an intermediate acoustic-level representation, usually in the form of a mel spectrogram. The "vocoder" then renders these acoustic representations into audio waveforms in a second stage. Given that singing voices exhibit a wider tonal range and greater dynamics than speaking voices, SVS presents a more complex challenge than TTS. Considerable research efforts have focused on developing effective acoustic models capable of producing expressive acoustic characteristics. In [GYR+21], autoregressive (AR)-based Bytesing is developed, which adopts a Tacotron[WSRS+17]-like structure and can generate high-fidelity speech, but at a slow speed. Later research showed that non-autoregressive (NAR) models such as FastSpeech 1&2 [RRT+19], [RHT+20b] can also generate high-quality speech for TTS but at a much faster speed. Therefore, FastSpeech-based SVS methods have also been developed, such as XiaoiceSing 1&2 [LWL+20], [WZH22]. To further improve the quality of synthesised singing voice, methods based on generative adversarial networks (GAN) are proposed, including HifiSinger [CTL+20], WGANSing [CBBG19], XiaoiceSing2 [WZH22] and SingGAN [HCC+22].

For the vocoder of SVS, being data-driven, a significant challenge with conventional TTS vocoders like MelGAN [KKDB+19] and HiFi-GAN [KKB20a] is their tendency to generate fragmented tones at unseen pitches. However, the ability to accurately render wide-pitch singing is a crucial requirement for SVS systems. To address this issue, neural source filter models [WTY19], [YWT23], [WHY+22] have been developed, which treat the singing voice as a filtered version of the stimuli source signal. By incorporating the parametric digital signal processing, such vocoders could generalise well to unseen pitch ranges. It has been argued that the two stages of SVS, i.e., acoustic modeling and vocoding, are separately optimised, the cascade combination may lead to a sub-optimal overall system. End-to-end (E2E) SVS systems, such as VISinger [ZCX+22] and UniSinger[HCH+23] are proposed to produce more realistic and expressive audios by leveraging the learned hidden representations.

Large speech language models have been employed to enhance TTS systems [DGCY22], [SDvRJ22], [LKL+22], [SL24]. These approaches either train their own SSL feature or use pre-trained features as guidance for their acoustic model. Similar techniques have been applied to SVS. Similar techniques have been applied to SVS. TokSing [WST+24] proposes token blending across models and layers to capture better singing semantics and acoustics. StyleSinger [ZHL+24] captures the timbre and emotion with the pre-trained wav2vec 2.0 feature from the reference recording. However, due to the lack of singing-dedicated pre-trained models, these approaches rely on pre-trained speech models, which may lead to sub-optimal performance.

Beyond the above single-purpose methods, NANSY++ [CYLK23] is a unified voice analysis and synthesis framework built with normalizing flow, that is able to perform 4 tasks including SVS, SVC and TTS. It is widely recognised that diffusion models can produce excellent performance in terms of generation quality of images [DN21], [RBL+22], videos [HNM+22], [HCS+22], and audio [LCY+23c], [LTY+23a], [YYW+23], therefore, Diffsinger [LLR+22a] and HiddenSinger [HLL23] are proposed to apply the diffusion to singing voice generation. However, since diffusion sampling typically involves tens to hundreds of steps, acceleration has become a critical issue. In [YXT+23], CoMoSpeech is proposed to apply the consistency model [SDCS23] to TTS and SVS with only one step of sampling to achieve fast speed.

Besides generating high-fidelity singing voices, many other researchers have been focusing on modelling the "singing technique", which is crucial to vocal performances. Typical works include [LCL22] for intensity and breath modelling, [KKJK23], [NKKK15], [SSZ+22] for vibrato control, [LCKL20], [NKKK15] for timbre modelling and [SSZ+22] that also enables controlling of the singing emotions.

*b) Singing Voice Conversion (SVC)::* SVC, analogous to voice conversion for singing, involves modifying various attributes of voices, most commonly changing the identity of the singer/speaker [SYKL21]. AutoVC [QZC+19] is an autoencoder designed for speech conversion conditioned on speaker embeddings, where tuning the bottleneck dimension allows the encoder to focus on speaker-independent features. This model was adapted for SVC in [Ner20], where AutoVC applies only on the harmonic spectral envelope produced by WORLD [MYO16]. More recent work further enhanced the AutoVC model with a latent regressor loss for improved identity embedding [OD23]. Other approaches based on disentanglement use Variational Autoencoders (VAE), phonetic posteriorgrams (PPG), and adversarial training. For example, [LHAH20] achieves many-to-many SVC by training a Gaussian mixture variational autoencoder (GMVAE) on non-parallel data. Another method [LTY+21] condenses phonetic information into PPGs while encoding other content information with a separate encoder. In [LZSL20], the authors adapted VAW-GAN [HHW+17] to the SVC task. The adversarial losses in [LTY+21] and [LZSL20] enhance the robustness.

Similar to SVS, self-supervised features have been applied to SVC for extracting rich and meaningful features. [JWS+23] explores the use of wav2vec 2.0 [BZMA20a] and HuBERT [HBT+21a] together with an $f_0$ harmonic generation module. Additionally, [WLT+22] demonstrates that HuBERT features outperform Mel-spectrogram and PPG-like features when using a contrastive predictive coding module. [HWY+23] replaces the traditional acoustic features input with WavLM [CWC+22a] to extract content features and for reconstruction.

[NJW+23] utilises Whisper's [RKX+23b] encoder to extract bottleneck features as content representation.

The introduction of diffusion decoders has significantly enhanced conversion quality in SVC. DiffSVC [LCSM21] is the first SVC method to use a diffusion decoder, though it is limited to any-to-one conversion. CoMoSVC [LYX+24] extends this capability to many-to-many conversion using a consistency model, achieving fast and high-quality inference. The so-vits-svc [Tea24] is an impactful open-source project providing state-of-the-art singing voice conversion tools based on VITS [KKS21]. It supports training with SSL features including HuBERT, Whisper PPG, and WavLM. Shallow diffusion [LLR+22b] is provided as a post-processing method.

*c) Lyrics Generation:* The LOAF-M2L model [OMW23] presents a melody-to-lyrics generation method. Through a hybrid training strategy that combines unsupervised learning with supervised training objectives based on musicological insights, LOAF-M2L improves the structure and syllable alignment of lyrics with respect to the melody, and it maintains a high level of text fluency. Significant improvements in both objective metrics and subjective assessments were validated.

### C. Music Therapy & Medical Applications

In the last three decades, music therapy research has seen a growing number of publications investigating the role of music in non-pharmacological care. Music has shown positive effects in supporting a number of different healthcare conditions and challenges such as depression [THZY20], stress and anxiety [dPS+22], pain analgesia [Lee16], dementia [MCMP20], and Alzheimer disease [MK22].

Although the use of music foundation models for healthcare is currently under-explored, we argue here that it could potentially boost the domain due to the generalisation capability acquired by unlabelled data and the enormous integrated domain knowledge learned by foundation models in other modalities. In other words, the use of foundation models can help both improve our understanding of the therapeutic effects of music while helping us generate better musical content to support music therapy interventions. By either implicitly learning interdisciplinary tasks by updating parameters during pre-training and fine-tuning, or explicitly specifying mapping relationships in a hand-crafted pipeline, foundation models can be utilised as a powerful plug-and-play component in music healthcare applications for a variety of tasks like developing bespoke technologies, improving accessibility, or building musical databases for specific interventions [ASV+21].

The existence of universal traits in the effects of music on human cognition has been challenged in recent years [SS15], and foundation models could help instead foster multimodal approaches that take into account a plurality of aspects in the care process [RBR24]. The rich semantics embedded in representations during pre-training [LYZ+22], [LYZ+24], [LHSS24] can help identify relationships between different modalities of treatment and foster personalised treatment [PR24], thus directing music therapy interventions towards new paradigms in the music generation process. Foundation models can also support the creation of interventions rooted in the biopsychosocial model of illness [Eng77], that proposes a multifaceted view of health conditions including biological, social and psychological factors such as personality traits, stress or socioeconomical status, that can improve the outcomes of music therapy interventions, in particular for chronic illnesses [WH17].

The relationship between multimodal biosignals, emotions and music has been explored in recent works [RGC+21] with the objective to further guide generative models and recommendation systems. Studies share a similar purpose when identifying audio features or metrics to select music tracks in playlists to improve for example sleep quality [ZWW10] or alleviate dementia [NdSC+24], looking to find optimal machine learning methods and suitable music representations. Multimodal datasets involving for example kinematic measures or biosignals such as electro-encephalogram (EEG) can further help explore the embodied relationship music has with our body and psychology, particularly when related to perceived and induced emotions [dSdLT+21], and affective computing applications involving recommendation systems have already been proposed in recent literature [GZJ+24], [JETT+23].

A natural step forward from recommendation systems would be to potentially leverage automatic music generation to meet the diversified needs of music creation in healthcare applications. For instance, generative foundation models could be introduced to enhance the production of a diversity of existing automatic improvisation techniques [XD17] or during therapy sessions involving improvisations [Wig04]. Music generation can also be used to evaluate and analyse interventions in different scenarios [LLLC22] or complement specific treatments [CPZ+20]. Fine grain controls, such as phrase or bar level generation, can be leveraged to guide the music generation process when designing a music therapy intervention aimed at eliciting specific functions or emotions dynamically [WMZ+23]. Music generation can be a powerful aid also in broader healthcare and social contexts, such as supporting people on the autism spectrum [JP23] or movement rehabilitation [FCZ+11], [BCA18], and already showing comparable results to traditional methods for binaural beats music therapy [YHH+24], an auditory illusion technique created by playing slightly different frequencies in each ear that the user will perceive as a single, pulsating tone.

MIR techniques have recently been applied to sonification (i.e., the translation of information from a given domain to the audio domain [KWB+10]), and can potentially improve a range of tasks including cancer diagnosis [WRK+19], sleep stage classification [MMR+20], recognising significant status changes during therapies [FMG+19], and other biosignal-related tasks that could benefit not only from analysis techniques but, potentially, from data augmentation or the production of ex-novo synthetic data, similar to analogue applications in protein sequencing [YQMB19]. Analysis techniques can also be applied a-posterior to allow for a reacher and more fine-grained analysis of the outcomes of a Music Therapy intervention [ZWM17]. Sonification has been seen as well to support movement and rehabilitation [NBG+16], and generative foundation models can further improve the ability of wearable and integrated systems to interact with the body in

both its physical and psychological facets.

The introduction of AI tools in Music Therapy should prompt also to accelerate the critical discussion around how music is currently conceptualised and applied in healthcare interventions [MCA24], using AI also to challenge the Western perspective that currently monopolises music therapy interventions [KKK24]. Participation and collaborative design are a key asset to ensure both therapists and patients are directly involved into the development of AI-driven Music Therapy interventions [SYZ+24], [DSCS23], mitigating some of the ethical and social issues that are discussed later in this work.

## IV. TECHNICAL DETAILS OF FOUNDATION MODELS

In this section, we introduce key technical aspects of developing foundation models, drawing from research in relevant fields. We start by detailing the main pre-training paradigms for music PLMs and multimodal foundation models (FMs), including LLMs and LDMs in Section IV-A. We then cover domain adaptation techniques such as finetuning on downstream tasks, instruction tuning and in-context learning in Section IV-B, and outline the design of audio tokenisers and model architectures in Sections IV-C and IV-D respectively. We conclude with a discussion of interpretability and controllability (Section IV-E), LLMs for music agents (Section IV-F) and scaling laws (Section IV-G), before looking at future work in Section IV-H. We note that many of the techniques developed for foundation modelling in other fields, such as instruction tuning and in-context learning, have not yet been fully explored in music. Similarly, research on scaling laws and emergent abilities of music foundation models is still in its infancy. In addition to this, there remain several open problems specific to music, such as long-sequence modelling (Section IV-H2), that require foundation model innovations in pre-training strategies or model architecture methodologies etc. We touch upon these open issues throughout this section.

### A. Model Design & Pre-training Strategies

Foundation models are pre-trained in a self-supervised fashion on large-scale datasets, avoiding or minimising the need for labelled data. As different pre-training strategies and model designs can lead to different capabilities and performance on different downstream tasks [YML+24], [RZP+20], [WKT+21], essential design considerations regarding tokenisation, architectures, training protocols, and other methodologies vary across foundation models, studies, and implementations. This subsection covers pre-training strategies for music foundation models, which we categorise into **Contrastive Learning and Clustering**, **Generative Pre-training**, and **Masked Modelling**. While many architectures and paradigms have been adapted to deal with music data, we observe that limited work has investigated the integration of music domain knowledge in the pre-training paradigm and the technique of instruction tuning remains largely unexplored.

### 1) Contrastive Learning & Clustering:

*a) Definition of contrastive learning:* Contrastive learning is a machine learning paradigm that typically provides impressive results in self-supervised learning. Models learn meaningful semantic representations by maximising the similarity between similar samples and minimising the similarity between dissimilar samples. In doing so, it learns high-level discriminative semantic information. Maximising similarity agreement was implemented prior to contrastive learning through Siamese networks [HCL06], which suffer from representation collapse without intervention. Prominent solutions that emerged to combat this trivial solution include Bootstrap Your Own Latent (BYOL [GSA+20]), Barlow twins [ZJM+21] (which both have applications in audio representation learning [NTO+21], [ACSS23]), and triplet networks [HA15]. Contrastive learning is a generalisation of the triplet framework to N pairs, which has seen great success in general representation learning in the fields of Computer Vision and NLP, but also in environmental audio, speech, and music.

**Formalization of contrastive learning** Concretely, contrastive learning teaches general representations of data by contrasting representations of instances and maximising the agreement between positive representations. At the same time, within a set of samples, representations of negative instances are pushed away from positive pairs. Formally, we consider a general form of the contrastive learning framework.

Let $\{z_i\}_{i \in [1...N]}$ be a set of embedding representations in the latent space of points $\{x_i\}$ in the data space, encoded by a model $E : x_i \mapsto z_i$ that maps the data space to the latent space. For each index $i$, let $p(i)$ be the index of the positive sample corresponding to index $i$. This positive sample is determined by a positive sampling process $\mathcal{P}$ that is task and design-dependent. The Info-NCE loss $\mathcal{L}_{infoNCE}$, introduced in [Soh16] is used as the objective function for a given pair

$$\mathcal{L}_i^{infoNCE} = -\log \frac{\exp(sim(z_i, z_{p(i)})/\tau)}{\sum_{j \in N(i)} \exp(sim(z_j, z_i)/\tau)} \quad (1)$$

Here $sim$ is a similarity function between representations and $\tau$ is a temperature hyperparameter. $N(i)$ is the set of negative indices within the batch for index $i$. In this formulation, the network learns to maximise the similarity between positive samples while minimising that of mutually negative samples. In doing so, the encoder captures high-level semantic information about samples necessary to perform the discriminative task.

**Positive instance sampling in contrastive learning** The notion of "positive sampling process" has been intentionally left vague previously, as it has varied significantly from study to study from the first implementation in SimCLR [CKNH20]. Contrastive learning, rather than a specific method, is a set of instance-discriminative approaches, which can be generalised across modalities, learning setups, architectures, and granularities. The canonical approach which was introduced in SimCLR is to implement a stochastic augmentation chain that generates two augmented versions of the same sample as positives and treats other augmented samples as negatives. However, the following work has explored alternatives such as masking inputs [HJL+23], generating a different version

Table I: Music Foundation Model Trained with Contrastive Learning & Instance Contrastive Learning.

| Model | Modality | Application | Training Paradigm | Music Tokeniser | Architecture |
|---|---|---|---|---|---|
| COLA | Audio (Speech, Sound & Music) | Understanding | Contrastive Learning | spectrum | CNN Encoder |
| MULE | Audio (Music) | Understanding | Contrastive Learning | spectrum | CNN Encoder |
| CLAP | Audio (Sound), Text | Understanding | Contrastive Learning | spectrum | Transformer Encoder |
| MusCALL | Audio (Music), Text | Understanding | Contrastive Learning | spectrum | CNN Encoder |
| MuLan | Audio (Music), Text | Understanding | Contrastive Learning | Spectrum | CNN Encoder & Transformer Encoder |
| CLAMP | Symbolic (MIDI), Text | Understanding | Contrastive Learning | MIDI | Transformer Encoder |
| Wav2CLIP | Audio (Sound), Text, Image | Understanding | Contrastive Learning | spectrum | CNN Encoder |
| AudioCLIP | Audio (Sound), Text, Image | Understanding | Contrastive Learning | spectrum | CNN Encoder |
| vq-wav2vec | Audio (Speech) | Understanding | MLM (Clustering via CL.) | 1-D CNN | CNN Encoder |
| wav2vec 2.0 | Audio (Speech) | Understanding | MLM (Clustering via CL.) | 1-D CNN | Transformer Encoder |
| HuBERT | Audio (Speech) | Understanding | MLM (Clustering via CL.) | 1-D CNN | Transformer Encoder |
| BEST-RQ | Audio (Speech) | Understanding | MLM (Clustering via CL.) | Spectrum | Transformer Encoder |
| musicHuBERT | Audio (Music) | Understanding | MLM (Clustering via CL.) | 1-D CNN | Transformer Encoder |
| MERT | Audio (Music) | Understanding | MLM (Clustering via CL.) | 1-D CNN | Transformer Encoder |
| MusicFM | Audio (Music) | Understanding | MLM (Clustering via CL.) | Spectrum, BEST-RQ | Conformer Encoder |



(a) Contrastive learning, which relies on an instance-discriminative framework to learn general instance-level or context-level representations.

(b) Generative Pre-training, which we categorise into three main families for music foundation models : VAEs, Diffusion models, and Autoregressive Predictive Coding (APC).

(c) Masked Modeling, where representations of music are learned through reconstructing masked inputs. Reconstruction targets vary by method.
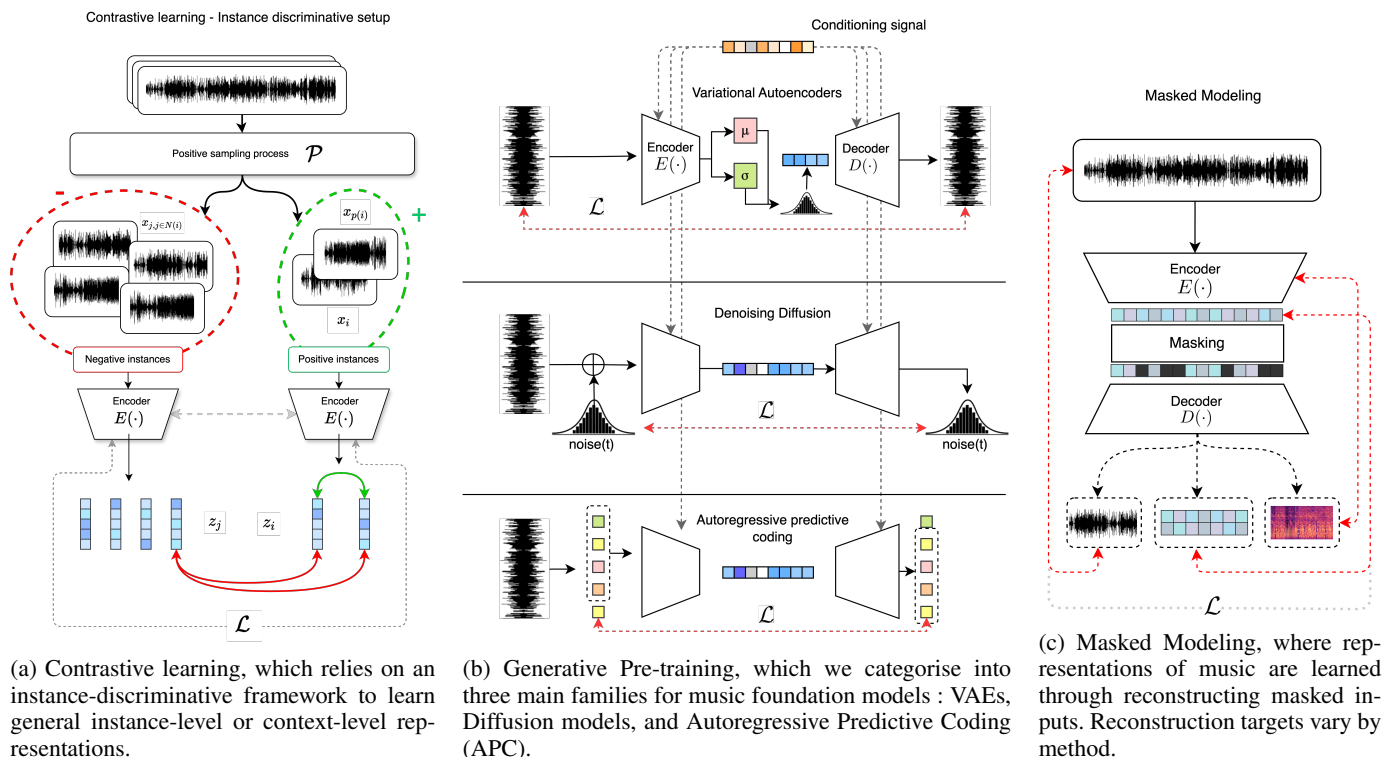
Figure 5: A broad taxonomy of pre-training strategies for Music Foundation Models. We categorise these strategies into Contrastive Learning (5a), Generative Pre-training (5b), and Masked Modelling (5c).

of the anchor through generative models [JPTI21], or using classifier-guided sampling [GS23], [YWZ+22], [HKW+22]. Some approaches consider a student-teacher setup in which positive pairs are determined by encoded representations from two encoders of the same sample [AAW+22], [CZH+21a].

In multimodal contrastive learning, a key development of the field which has enabled numerous foundation models, positive samples are determined from multimodal pairs [RKH+21b], [WCZ+23a], [HJL+22], [MBQF22a]. Further, we make the distinction between **instance contrastive learning** and **context contrastive learning**: In addition to parameterising the global similarity of data samples (context contrastive learning), contrastive learning can also parameterise the similarity between granular parts (instances) of a data sample, i.e. individual tokens of a sequential representation. [BZMA20b], [LYZ+24], [HBT+21b], [vdOLV18], [GLCG22b], [AAW+22]

Apart from data augmentation, the more general implementation of a positive sampling strategy rejoins a key consideration of contrastive learning: When teaching a model a similarity metric, how does one determine the best similarity metric to learn? Different positive sampling strategies lead to different similarity metrics, for which the model might need to capture different semantic information about the underlying data. Studies have explored the effect of different augmentation chains and sampling strategies on the learned representations, and there exists a general understanding that contrastive learning is largely dependent on the choice of positive sampling methods,

which is a key design consideration for contrastive training pipelines.

**Batch size considerations for contrastive training** Another key consideration is the batch size. Larger batches generally provide better downstream performance for contrastive foundation models in optimal settings [LKHS20], [BIS⁺23] because they provide more negative samples per batch. However, the batch formulation of the InfoNCE loss prevents gradient accumulation, which inherently bottlenecks contrastive learning approaches to computational requirements. As models scale up, training pure contrastive foundation models becomes more prohibitive. One framework that alleviates these computational costs is MoCo [HFW⁺20], which reframes contrastive learning as a key-query problem and introduces a gradient-free momentum encoder and key queue. Later implementations of the framework have since discarded such devices [CFGH20], [CXH21]. If the similarity evaluation procedure is done during preprocessing before pre-training, such as the clustering of audio features in HuBERT[HBT⁺21b], then the restriction of batch size can be mitigated by gradient accumulation.

*b) context contrastive learning:* **Examples in audio & music** Contrastive learning has been used for large-scale training approaches in environmental audio and speech representation learning. Notable examples of context contrastive learning for learning global audio representations include COLA [SGZ21], CLAR [ATM21] and CL-SER [FOM⁺21]. The most recent framework leverages a Normaliser-Free (NF) SlowFast (SF) Convolutional network and is currently SOTA across multiple tasks for context contrastive learning for speech and environmental audio [WLW⁺22].

In music, contrastive learning was first implemented with a small sampleCNN [LPKN17] network following the SimCLR approach on raw waveforms with Contrastive Learning of Musical Representations [SB21], which was later adapted with a Tailed U-Net architecture [VB⁺22]. S3T implements a MoCoV2 setup with Swin transformer encoders for music representation learning [ZZZ⁺22]. Arguably, the only unimodal contrastive foundation model for music is MULE, which was trained on 1.7M tracks from a private distribution catalog [MKO⁺22], both in a self-supervised and a supervised fashion. MULE uses the same SF-NFNet as [WLW⁺22] and approaches current SOTA on multiple tasks [YML⁺24]. As in other domains, key considerations for contrastive learning for music is the positive sampling strategy, including augmentation strategies [SB21], [MDH⁺24], relative position of the samples within the track [CJCC22], and use of generative or multi-source strategies [CPM⁺24], [GZM23].

**Multimodal contrastive learning** Multimodal contrastive foundation models are arguably the most widely used large-scale context contrastive learning approaches. The core principle of maximising similarity between positive samples remains the same, but in multimodal contrastive learning (MMCL), positive pairs originate from different modalities (e.g. text and images [RKH⁺21b], text and audio [WCZ⁺23a], text and music [HJL⁺22], [MBQF22a], image and audio [WSKB22], [GRHD22]). This requires projecting them into a shared latent space, which requires modality-specific encoders. Furthermore, as batches are better represented as from the same

modality, we reformulate the contrastive loss for its general usage in MMCL, with modalities $m_1$ and $m_2$ corresponding to encoders $E_1, E_2$ which map samples from modality 1 $x^{m_1}$ and 2 $x^{m_2}$ to a shared latent space: $E_1 \mapsto z_{m_1}, E_2 : x \mapsto z_{m_2} | z_{m_1}, z_{m_2} \in \mathbb{R}^d$, where $d$ is the dimensionality of the shared latent space.

$$\mathcal{L}_i^{m_1 \to m_2} = -\log \frac{\exp(sim(z_{m_1,i}, z_{m_2,i})/\tau)}{\sum\limits_{z \, \in \{z_{m_2}\}} \exp(sim(z_{m_1,i}, z)/\tau)}$$

and

$$\mathcal{L}_i^{m_1 \leftrightarrow m_2} = \mathcal{L}_i^{m_1 \to m_2} + \mathcal{L}_i^{m_2 \to m_1}$$

The first implementation of MMCL is CLIP [RKH⁺21b], which has been instrumental in the development of text-image generation models as well as captioning and retrieval approaches. Analogous applications have resulted from CLAP [WCZ⁺23a] and derivative work, which applies MMCL between audio and text. The original CLAP [WCZ⁺23a] uses a simple PANN [KCI⁺20] as the audio encoder and BERT as the text encoder [KT19]. The following studies scale up both the encoders to an HTSAT architecture [EDAIW23], [WCZ⁺23b], [CDZ⁺22] for audio and Respectively RoBERTA [WCZ⁺23b] and GPT-2 [EDAIW23], [RNS⁺18] for text, as well as the pre-training data using large-scale datasets [WCZ⁺23b]. In addition to keyword-to-caption augmentation, [WCZ⁺23b] introduces feature fusion using attention pooling in the audio branch of the model.

In music specifically, key models include MusCALL [MBQF22a] and MuLan [HJL⁺22], which have been leveraged in generation and have been shown to hold meaningful representations for general music understanding tasks. MusCALL implements a ResNet50 encoder and a simple BERT text encoder and includes intramodality supervision and semantic weighing of negatives by textual similarity as training design devices [MBQF22a]. MuLan, in contrast, uses a straightforward multimodal contrastive training setup, with an Audio Spectrogram Transformer [GCG21] audio encoder and a BERT text encoder. However, MuLan drastically scales up training data to close to 40 million text-music pairs. While MusCALL [MBQF22a] uses average + attention pooling of audio as the contrastive embedding, MuLan uses the `[CLS]` token from the AST and the BERT encoder for the contrastive task. In the field of symbolic music, recent work on contrastive learning between textual descriptions and textual representations of music has bridged the gap between symbolic music and text [WYTS23a]. The encoders used for this work are a fine-tuned DistilRoBERTa and a MAE-style BERT encoder pre-trained on a large-scale ABC notation dataset (see sections II-B, V-A1) [WYTS23a]). Other modalities have been leveraged for MMCL for music such as metadata [FDVS20] and playlists [AJFF⁺23], but never to the scale of text-music contrastive approaches.

Finally, efforts have been made to bridge the representations learned by multimodal contrastive models by either distilling CLIP knowledge into audio or creating any-to-any

joint embedding spaces. Key approaches include Wav2CLIP [WSKB22] and AudioCLIP [GRHD22]

*c) Instance Contrastive Learning:* Instance contrastive learning has had high success rates in the field of speech representation learning, where numerous foundation models apply contrastive learning loss granularly between individual tokens from extracted sequential representations. The similarity of such tokens can lead to pseudo-labels for pre-training, making the contrastive learning loss as classification of mask modelling discussed in IV-A3.

The first implementation of such approaches was Contrastive Predictive Coding, which uses a contrastive loss between an aggregated context and future samples as positives and past samples as negatives [vdOLV18], which was applied for speech and image representation learning. The following approaches for speech and general audio include Wav2Vec [SBCA19], which uses larger front-end and context encoders as well as a convolutional context model rather than a GRU. VQ-Wav2Vec follows by introducing a quantisation module after the convolutional frontend (Gumbel-softmax or K-means clustering) for downstream BERT training. Wav2Vec2.0 [BZMA20b] leverages a transformer architecture as a context model. Embeddings from the frontend encoder are masked before the transformer module and a contrastive loss is applied to identify the GS-quantised frontend embedding corresponding to each masked transformer output SSAST [GLCG22b] also applied an auxiliary contrastive loss to masked modelling between masked transformer tokens and reconstructions of other tokens using an AST decoder-only setup. W2V-BERT [CZH+21a] extracts quantised convolutional embeddings in the same fashion as [BZMA20b]. the unquantified embeddings are masked and passed to a conformer stack, from which an embedding sequence is extracted after $N$ out of $N + M$ blocks. A contrastive loss is applied between tokens of this intermediary sequence and corresponding target embeddings.

Finally, instrumental to MERT, a key recent music foundation model, is the contrastive learning approach in HuBERT [HBT+21b]. HuBERT and musicHuBERT[MYL+23] apply masked modelling on discrete tokens to learn representations (see section IV-A3). To obtain these discrete tokens, the authors use an offline convolutional encoder to obtain continuous embedding sequences. the distribution of embeddings from these sequences is parameterised with a contrastive loss and token indices are obtained through K-Means clustering. MERT [LYZ+24] applies the same contrastive approach to discover its token vocabulary for one of its acoustic teachers (see Section IV-A3).

*2) Generative Pre-training:*

*a) Autoencoders:* Autoencoders (AEs) [RHW86] are designed to learn lower-dimensional feature representations of data in an unsupervised manner. The basic architecture consists of an encoder compressing the input data into a lower-dimensional latent space and a decoder attempting to reconstruct the original input from it. Formally, given an input $x$, an encoder function $E$, and a decoder function $D$, an autoencoder aims to minimise the reconstruction loss

$$\mathcal{L}(x, D(E(x))),$$

where $\mathcal{L}$ is typically a mean squared error loss for continuous data or cross-entropy loss for discrete data. Given the sole focus on reconstruction, conventional AEs don't inherently model the data distribution, nor do they guarantee that their latent space is continuous or well-structured. As such, sampling arbitrary points from their latent space and decoding them is not likely to generate sensible outputs in many cases.

Variational Autoencoders (VAEs) [KW14] address these limitations by introducing a probabilistic framework scalable to large datasets. Instead of encoding inputs to fixed points in latent space, VAEs encode them as probability distributions, typically Gaussian. The encoder outputs parameters of this distribution (mean $\mu$ and variance $\sigma^2$), and the decoder reconstructs from samples drawn from this distribution. VAEs are trained to minimise both reconstruction loss and the Kullback-Leibler divergence between the encoded distribution and a prior (usually a standard normal distribution):

$$\mathcal{L} = \mathbb{E}[\log p(x|z)] - \mathbb{D}_{\text{KL}}[q(z|x)||p(z)],$$

where $x$ is the input, $z$ is a latent variable sampled from the encoded distribution, $\mathbb{E}[\cdot]$ denotes the expected value, and $D_{\text{KL}}(\cdot||\cdot)$ is the Kullback-Leibler divergence term acting as a regulariser, encouraging the encoded distributions to be close to the prior distribution. This formulation allows VAEs to generate novel samples by sampling from the prior distribution and decoding, making them more promising than conventional AEs for generation. Still, because they optimise for the average reconstruction error over the latent distribution, they often tend to produce blurry or averaged-out reconstructions. They also suffer from posterior collapse, because the KL divergence term can be minimised by making the encoder output match the prior regardless of the the input. VAEs are typically used in combination with diffusion models, providing better results [RBL+22], [FM22b].

Vector Quantised VAEs (VQ-VAEs) [OVK17] instead use a discrete latent space, in which the encoder output is mapped to the nearest vector in a learned codebook, preventing collapse by forcing the encoder to output semantically meaningful, discrete codes. The discrete latent space can better capture structured representations (which is particularly useful for modalities with inherently discrete elements like symbolic music) and often produces more detailed reconstructions. The VQ-VAE training process involves optimising a loss function with three main components:

- Reconstruction loss: $\mathcal{L}_{recon} = \log p(x|z_q(x))$, where $x$ is the input and $z_q$ is the quantised latent vector.
- Codebook loss: $\mathcal{L}_{codebook} = ||sg[z_e(x)] - e||_2^2$, where $z_e$ is the encoder output, $e$ is the selected codebook vector, and $sg[\cdot]$ is the stop-gradient operator.
- Commitment loss: $\mathcal{L}_{commit} = \beta||z_e(x) - sg[e]||_2^2$, where $\beta$ is a hyperparameter.

The total loss is thus $\mathcal{L} = \mathcal{L}_{recon} + \mathcal{L}_{codebook} + \mathcal{L}_{commit}$. During training, the encoder learns to map inputs to continuous latent vectors that are close to codebook entries. The codebook itself is trained by moving its vectors towards the encoder outputs that select them. This is achieved through an exponential moving average update or by directly optimising the

Table II: Music Foundation Model with Generative Pre-training including VAE, GPT, multimodal GPT, and Audio Diffusion.

| Model | Modality | Application | Training Paradigm |
|---|---|---|---|
| Jukebox, JukeMIR | Audio (Music) | Both | VAE, GPT |
| MusER | Symbolic (MIDI) | Generation | VAE |
| Singsong | Audio (Music) | Generation | GPT |
| AudioLM | Audio (Sound), Text | Generation | GPT |
| MusicGen | Audio (Music), Text | Generation | GPT |
| MusicLM | Audio (Music), Text | Generation | GPT |
| Music Transformer | Symbolic (MIDI) | Generation | GPT |
| pop music Transformer | Symbolic (MIDI) | Generation | GPT |
| Jazz Transformer | Symbolic (MIDI) | Generation | GPT |
| MelodyGLM | Symbolic (MIDI) | Generation | GPT |
| MUPT | Symbolic (ABC) | Generation | GPT |
| SpeechGPT | Audio (Sound), Text | Both | GPT |
| LauraGPT | Audio (Sound), Text | Both | GPT |
| Audio-PaLM | Audio (Sound), Text | Both | GPT |
| MuseCoCo | Symbolic (MIDI), Text | Generation | GPT |
| ChatMusician | Symbolic (ABC), Text | Both | GPT |
| AudioLDM | Audio (Sound), Text | Generation | Diffusion |
| AudioLDM2 | Audio (Sound), Text | Generation | Diffusion |
| Make-An-Audio 1 | Audio (Sound), Text | Generation | Diffusion |
| Make-An-Audio 2 | Audio (Sound), Text | Generation | Diffusion |
| Stable Audio Open | Audio (Sound), Text | Generation | Diffusion |
| CRASH | Audio (Music), Score | Generation | Diffusion |
| Noise2Music | Audio (Music), Text | Generation | Diffusion |
| Mousai | Audio (Music), Text | Generation | Diffusion |
| MusicLDM | Audio (Music), Text | Generation | Diffusion |
| TANGO | Audio (Music), Text | Generation | Diffusion |
| JEN-1 | Audio (Music), Text | Generation | Diffusion |
| Diff-A-Riff | Audio (Music), Text | Generation | Diffusion |
| GETMusic | Symbolic (MIDI) | Generation | Diffusion |
| whole-song-gen | Symbolic (MIDI) | Generation | Diffusion |

codebook loss. The commitment loss encourages the encoder to "commit" to codebook vectors, preventing its output from fluctuating excessively.

Residual Vector Quantised VAEs (RVQ-VAEs) [JG82] address limitations of VQ-VAEs, namely the fixed codebook size and potentially limited codebook utilisation. RVQ-VAEs extend the VQ-VAE concept by using multiple VQ layers in a hierarchical fashion. Each layer quantises the residual error from the previous layer, allowing for more fine-grained representations. This approach aims to increase representational capacity by using multiple codebooks, encourage better use of all codebook entries across different levels of abstraction, and allow for better preservation of fine details through the residual nature of the quantisation.

Some works in the music domain have explored the use of conventional AEs and particularly VAEs for primarily the tasks of music generation [RER+18b], [YWW+19], [WZZ+20b], [JXC+20], [CE21], [TH20] and style transfer [BKWW18], [WY23b], although they don't fit most definitions of being a 'foundation model' due to their relatively small size and limited amount of training data used. These approaches are primarily in the symbolic music domain [RER+18b], [BKWW18], [YWW+19], [WZZ+20b], [JXC+20], [CE21], [WY23b], leveraging architectures such as RNNs, GRUs, LSTMs, and even transformers. Some work exists on audio as well [LAH19], [TH20] using a Gaussian Mixture distribution rather than a Gaussian one in order to further aid the disentanglement of musical attributes in the VAE's latent space.

More relevant for music foundation models has been the work on neural audio codecs (NACs), which are a machine learning-based alternative to audio compression. NACs aim to learn a codebook of discrete audio tokens that can efficiently represent audio signals with the highest fidelity possible for a given codebook size. Representations from NACs are often used for quantisation and tokenisation, as described further in Sections IV-A1 and IV-A3. One of them, EnCodec [DCSA22], uses an RVQ-VAE with convolutional layers and a time- and frequency-domain reconstruction loss along with multiple codebooks to support variable bandwidth training. It further uses a multi-scale spectrogram discriminator as a perceptual loss, and, optionally, entropy coding with a small transformer to further compress the RVQ representation. SoundStream [ZLO+21] uses a similar architecture but proposes an end-to-end training procedure that incorporates the reconstruction and adversarial losses. The Descript Audio Codec [KSL+23] employs a range of further optimisations to increase compression, such as a periodic activation function, a modified quantiser dropout, a multi-scale STFT discriminator, and a multi-scale mel reconstruction loss with varying mel bin sizes.

*b) Diffusion models:* Diffusion models are a class of probabilistic generative models that have recently become popular due to the generation quality, particularly in the domain of image generation, for their capacity to generate high-quality synthetic data [ERO21], [YXK+22], [YLK+22], [HJA20]. The core idea behind diffusion models is to model the data distribution by simulating a diffusion process that gradually transforms data into noise, and then learn to reverse this process. This approach contrasts with other generative models like Generative Adversarial Networks (GANs) [GPAM+20] and Variational Autoencoders (VAEs), which directly map random noise to data samples. Diffusion models have made notable progress in the fields of audio and music generation.
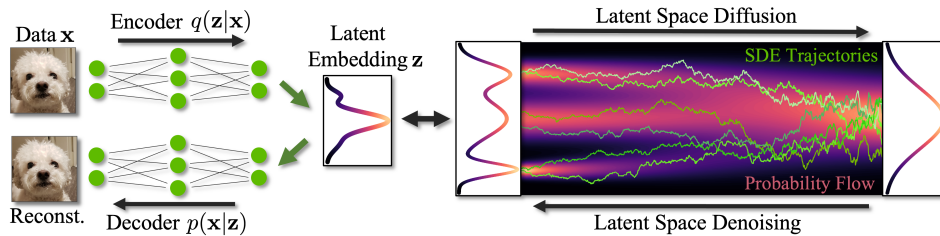
Figure 6: Paradigm of latent diffusion models. [KK23]

Section III-B explores these applications thoroughly. In this paragraph, we focus on the fundamentals of diffusion models and delve into design choices for prominent foundation models utilising diffusion for audio and music.

**Fundamentals of diffusion models**: diffusion models generate samples by learning to iteratively reverse a diffusion process: The diffusion process can be understood as a sequence of $T$ steps where a data sample $\mathbf{x}_0$ is progressively transformed into a noisy sample $\mathbf{x}_T$ through a series of Gaussian noise additions. The forward process is defined by a Markov chain:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t, \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1-\alpha_t)\mathbf{I})$$

which can be reparameterised as:

$$\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_0 + \sqrt{1-\alpha_t}\epsilon_t, \epsilon \sim \mathcal{N}(0, \mathbf{I})$$

where $\alpha_t$ are determined by a noise schedule controlling the amount of noise added at each step. The reverse process, which is the core of the diffusion model, aims to denoise $\mathbf{x}_t$ to $\mathbf{x}_0$. It is parameterised as:

$$p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t))$$

The parameters $\theta$ are learned to approximate the reverse diffusion process using variational inference. The objective is to minimise the variational bound on the negative log-likelihood of the data, which translates to:

$$\mathcal{L} = \mathbb{E}_q \left[ \sum_{t=1}^{T} D_{KL}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1} \mid \mathbf{x}_t)) \right]$$

Denoising diffusion probabilistic models [HJA20] introduce reparameterisation devices that largely simplify the formulation of the objective and the sampling strategy by reframing the model as a noise-predictive model $\epsilon_\theta$ :

$$\mathcal{L}_{\text{simple}} = \mathbb{E}_{t,\mathbf{x}_0,\epsilon} \left[ \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\|_2^2 \right]$$

Among key considerations for diffusion models, the noise schedule (i.e. $s : t \in [1 \cdots T] \mapsto \alpha_t \in \mathbb{R}$) , the network architecture, the conditioning mechanism (covered in a following paragraph), and the guidance mechanism.

**Accelerated sampling strategies**: As diffusion models typically require numerous diffusion steps to reach generated samples of a desired quality, a prevalent focus of research on diffusion models has been requiring fewer sampling steps to reach a given generation quality. Approaches such as truncation [ZHCZ22], [LXY+22] and knowledge distillation [SH22], [HZZ+24] have been generally successful. However, perhaps the most widely used is the Denoising Diffusion Implicit Model (DDIM) sampling procedure [SME21], which is also largely used in audio diffusion models [LCY+23c], [LTY+23a], [CWL+23a], [SJS23b], [GMMP23a]. Briefly, DDIM approximates the stochastic sampling process of DDPM by a simple probability flow ODE:

$$d\bar{x}(t) - \epsilon_\theta^{(t)} \left( \frac{\bar{x}(t)}{\sqrt{\sigma^2+1}} \right) d\sigma$$

Where $\sigma = \sqrt{(1-\alpha)/\alpha}$, $\bar{x} = x/\sqrt{\alpha}$. The deterministic nature of solving this equivalent ODE greatly speeds up the sampling process with minimal quality loss.

**Latent diffusion models**, which have been leveraged to great effect in image generation [RBL+22] and further in audio and music generation [LCY+23c], [LTY+23a], [SJS23b], [GMMP23b], [GMMP23a], [HHY+23], apply the same principles as diffusion models with a key difference that has been foundational in recent studies for music and audio diffusion. While traditional DDPMs model diffusion directly in the data space (i.e. pixel or sample space), latent diffusion models do so in the latent space of a pre-trained autoencoder. This provides the distinct computational advantage of not modelling expensive diffusion steps in a high-dimensional space and has sped up audio and music generation systems considerably, making diffusion systems tractable from an application standpoint.

**Conditional diffusion models**: Diffusion models, like many generative models, can be conditioned on various signals to guide the diffusion process. Most notable, perhaps, is the use of text conditioning for image, audio, and music. Generally, we consider a conditioning signal $C$, which is used to influence the generation process. Commonly, conditioning is applied through cross-attention in the diffusion architecture (U-Net [HJA20] or Transformer [PX23]).

Finally, diffusion models have largely adopted classifier-free guidance as a conditioning strategy for training and generation, after some works used "vanilla" guidance such as addition or self-attention [AFL23], [HLJK23] or classifier guidance [DN21]. which leverages a pre-trained classifier's gradient to change the denoising direction. Classifier-free guidance simplifies the guidance process by alternating between unconditional and conditional denoising during training with a "conditioning dropout probability" [HS22a]:

$$\nabla_x \log p(x|c) = w\nabla_x \log p(x|c) + (1-w)\nabla_x \log p(x)$$

This, at inference time, is simple to implement by denoising a weighted average of conditional and unconditional weights. This has been shown to produce superior results to other guidance mechanisms. In audio and music specifically, it is the predominant mechanism of guidance for diffusion models.

**Diffusion design choices in audio and music generation**: We now refer back to Section III-B for a thorough overview of key diffusion models for audio and music. This paragraph covers specific design choices of interest for these models.

Controllable Raw Audio Synthesis with High-Resolution (CRASH) [RH21] is a score-based generative model tailored for the unconditional synthesis of raw audio. Leveraging the early approach of diffusion processes modelled by stochastic differential equations, CRASH is engineered to produce high-fidelity drum sounds at 44.1 kHz.

Noise2Music [HPW$^+$23a] utilises an efficient 1D UNet conditioned on text through cross-attention and classifier-free guidance for diffusion on raw audio. It is the first diffusion model to generate 3.2kHz audio conditioned on text. It includes an upsampler diffusion model conditioned on the LoFi audio and text that upsamples the audio to 16kHz, and a super-resolution model produces the final 24kHz audio output.

Mousai [SJS23b] leverages a similar setup with latent diffusion: it uses cascaded diffusion U-Nets conditioned with cross-attention and classifier-free guidance on text encodings. Authors pre-train a diffusion autoencoder on magnitude spectrograms and then condition An efficient 1D U-Net on the obtained latent. At inference, the latent and the audio are generated through diffusion with a DDIM Sampler. This cascaded architecture allows for long-form coherence and fast generation.

In the field of audio generation, AudioLDM and [LCY$^+$23c] subsequently AudioLDM2 [LTY$^+$23a] have had influence on music generation through derivative work MusicLDM [CWL$^+$23a]. AudioLDM uses a mel-spectrogram VAE and a HiFiGan [KKB20b] vocoder as frozen components and trains a classifier-free guided diffusion UNet to generate audio. One specificity of AudioLDM which is often echoed in other work is the use of CLAP global audio and text conditioning during training instead of token sequence conditioning through cross-attention. This is due to the notable scale difference between text-image and text-audio datasets, which allows text-guided image generation to exploit text for training, while such augmentation strategies as proposed in AudioLDM are beneficial to alleviating the lack of captioned data. AudioLDM2 uses another approach to circumvent the lack of captioned data. By finetuning a GPT2 model to predict the sequential output of an AudioMAE encoder conditioned on text, audio or phoneme data, AudioLDM2 creates a shared sequential conditioning space dubbed "Language of Audio". They then employ cross-attention CFG conditioning to train the LDM. The authors of MusicLDM largely take inspiration from AudioLDM, but retrain CLAP on music data and propose a beat-synchronous mixup strategy to enhance diversity and novelty in output generations.

TANGO [GMMP23a] trains a VAE on mel-spectrograms and employs cross-attention to condition the same latent diffusion UNet as AudioLDM1 [LCY$^+$23c] on T5 text embeddings, with classifier-free-guidance. In contrast to AudioLDM, TANGO does not leverage CLAP encodings for text conditioning, removing the limitations induced by singe-element embeddings and trusting T5 to hold necessary information for the latent diffusion model to learn the appropriate intermodal mapping.

Diff-A-Riff [NPA$^+$24] is a latent diffusion model specifically designed to generate high-quality instrumental accompaniments suitable for different musical contexts. With integrated controls for audio referencing and text cueing, Diff-A-Riff enhances compositional flexibility, generates 48kHz pseudo-stereo audio, and significantly reduces inference time and memory usage.

Other key work with diffusion models includes JEN-1 [LCY$^+$23b], which employs a standard diffusion UNet conditioned on T5 embeddings with the added specificity that authors include an autoregressive training objective by including causal padding in self-attention UNet blocks, which allows for multiple generation use-cases. Make-an-audio 1 [HHY$^+$23] and Make-an-audio 2 [HRH$^+$23] focus on captioned data scarcity and temporal alignment respectively. Make-an-audio 1 leverages a standard cross-attention Diffusion UNet with classifier-free guidance and trains on CLAP-score selected captions generated from pre-trained automatic captioning models. Make-an-audio leverages temporal encoding of audio events through LLM augmentation to enforce temporal coherence of occurrences between prompts and generated audio. Most recently, Stable audio [ECT$^+$24] uses a proprietary-trained CLAP model and timing embedding information to condition a latent diffusion UNet for timing-conditioned audio generation. specifically, authors make generation faster by using a memory-efficient attention implementation. Specific details include the use of the next-to-last layer of the CLAP text encoder for conditioning and FiLM conditioning of the diffusion timestep. Finally, Stable audio 2 [EPC$^+$24] upscales stable audio 1 by using a Diffusion Transformer Architecture with 1.1 Billion parameters. Further, authors employ a DAC-like autoencoder [KSL$^+$23]. This architecture upscaling allows Stable Audio 2 to generate long-form (upwards of 3 minutes) music with better objective metrics on music generation compared to previous models.

Finally, Diffusion has also emerged as a promising paradigm in symbolic music generation. Authors of [Ata23] use a standard diffusion UNet pipeline with DDPM sampling to generate binary piano rolls representing MIDI data. [WMX24] leverages a similar approach with hierarchical generation to first generate form, lead sheets structure, lead sheets, and accompaniment with conditioning from previous stages. All four stages are diffusion UNets with classifier-free guidance from previous stages. Finally, latent diffusion is used to generate MusicVAE [RER$^+$18a] latent embeddings in [MEHS21]. Authors leverage a transformer as diffusion backbone for this approach with musical scale conditioning through a FiLM layer.

*c) Autoregressive Predictive Coding:* Autoregressive Predictive Coding (APC) has had immense success in the field of natural language processing, largely facilitated by the widespread usage of the transformer architecture [VSP+17b]. Key applications include the GPT language model family [RNS+18], [BMR+20], which uses APC among other techniques as a pre-training paradigm. Conceptually, APC trains a model to predict the next element in a sequence of discrete tokens using an autoregressive architecture. In doing so, the model learns contextual representations over the whole sequence that contain high-level semantic information. Overwhelmingly, the architecture used for APC pre-training is the transformer, although recent advances in state space models (SSMs) have begun encroaching on the territory of the transformer model [GD23]. For more information, please refer to subsection IV-D

**Fundamentals of APC:** Formally, consider a sequence of discrete tokens (codes) $X = (x_1, x_2, \ldots, x_n)$. The goal of Autoregressive Predictive Coding (APC) is to model the joint probability distribution of the next token given the previous tokens in the sequence using an autoregressive model $A_\theta$:

$$P(X) = \prod_{t=1}^{n} P(x_t \mid x_{<t}; A_\theta)$$

where $x_{<t} = (x_1, x_2, \ldots, x_{t-1})$ represents the tokens preceding $x_t$. The model $A_\theta$ is applied to the sequence $x_{<t}$ to predict $x_t$ $(x_{\leq t} = A_\theta(x_{<t}))$. Optionally, the model can be conditioned on an additional sequence of context tokens $C = (c_1, c_2, \ldots, c_m)$ pertaining to any conditioning signal or modality:

$$P(X \mid C) = \prod_{t=1}^{n} P(x_t \mid x_{<t}, C; A_\theta)$$

During training, the model is optimised to maximise the conditional log-likelihood, formulated as the cross-entropy loss:

$$\mathcal{L} = -\mathbb{E}_{X,C} \left[ \sum_{t=1}^{n} \log P(x_t \mid x_{<t}, C; A_\theta) \right]$$

**Model design in APC for audio, symbolic music, and acoustic music:** For audio and music specifically, a key design consideration beyond the usual context length, attention mechanism, model size, choice of positional encodings, and conditioning mechanism, is the choice of the discrete audio tokeniser. APC requires discrete tokens for prediction, so to perform APC pre-training on audio sequences, these sequences must be converted into discrete tokens (codes). Varieties of tokens and tokenisers are covered in Subsection IV-C, but in this section, we briefly delve into the design choices of foundation APC models in audio, speech, and music for which tokeniser to leverage for APC pre-training as well as specific training details. The range of applications for APC-based models in audio and music generation and understanding as well as multimodal approaches are well covered in Section III. Here, we cover specific architectural details as well as training devices used for these models.

In generic audio, AudioLM [BMV+23] utilises a Soundstream encoder [ZLO+21] to extract acoustic tokens and k-means-clustered w2v-BERT [CZH+21a] representations to extract semantic tokens. AudioLM then cascades three autoregressive transformer models to model semantic tokens, then coarse acoustic tokens conditioned on semantic tokens, and finally fine acoustic tokens conditioned on coarse acoustic tokens. The sample rate of acoustic vs semantic tokens for AudioLM is 50Hz for acoustic tokens and 25Hz for semantic tokens, both with a vocabulary size of 1024.

In acoustic music, a key foundation model that employs APC for pre-training is Jukebox [DJP+20b]. The authors train three VQ-VAE (vocabulary size 2048) models to reconstruct three temporal resolutions of music with losses accounting for reconstruction accuracy, codebook usage, and embedding stability. Three autoregressive transformers are used to model these sequences, each conditioned on upsampled tokens from the one-level-coarser transformer model. Jukebox is notable as the largest foundation model for music, sitting at 5B billion parameters and trained on the largest scale amongst other approaches. It has notably served as a foundation model for many downstream applications and has been shown to hold competitively informative representations for downstream probing, for which it is notably state of the art for many tasks to this day [CDL21]. Other APC-trained acoustic music models include MusicGen [CKG+23b], which is trained to model sequences of Encodec [DCSA22] tokens, with novel token interleaving patterns to alleviate the computational costs of generating multiple codebook streams. MusicLM [ADB+23] implements a similar setup to audioLM with soundstream acoustic tokens, w2vBERT semantic tokens, and possible text conditioning using MuLan [HJL+22] embeddings and a similar hierarchical transformer cascade to AudioLM [BMV+23] to predict semantic tokens then fine-grained acoustic tokens.

In symbolic music, APC has also been used for understanding and generation, where tokenisers are key to generation performance as they encode information from MIDI signals into tokens. Key approaches include Music Transformer [HVU+18a], pop music transformer [HY20b], Jazz transformer [WY20], And MuseCoCo [LXK+23a], which pretrains a BERT extractor for text-attribute pretext tokens and then trains an autoregressive transformer on text-conditioned symbolic music generation to great effect.

**Model design in APC for multimodal audio and music understanding:** APC-trained audio and music language models have also been leveraged for LLM-based audio and music understanding. leveraging representation sequences from these models has proven to be effective for multimodal music understanding and captioning models, also trained or fine-tuned using APC. Most notable examples in the audio domain include listen, think, understand [GLL+23], which makes use of LoRa adapters applied to a pre-trained LLaMa [TLI+23b] LLM with acoustic tokens obtained from a pre-trained AST model to generate responses to user text queries including descriptive questions and chain-of-thought reasoning. Audio-PaLM [RAN+23] use a pre-trained PaLM [ADF+23] model adapted to generate text as well as audio tokens, which is done by adding audio-token rows to the token embedding matrix and

fine-tuning the whole architecture. In doing so, AudioPaLM is able to generate both audio and/or textual responses to audio and/or textual user inputs, augmenting its versatility. Qwen-audio [CXZ+23] integrates generative pre-training with an audio encoder and frozen LLM as a natural language decoder. This architecture strategically leverages multi-task training to handle a variety of audio types and tasks, such as speech transcription and translation. By adjusting the decoder of hierarchical label sequences, Qwen-Audio can effectively cope with the challenges caused by label changes in different datasets. Besides, it can also handle basic music classification and description tasks on singer identification, emotion, genre, instrument, etc. Qwn2-audio [CXY+24] has a similar pre-training process with expanded data volume and better natural language prefix for different data and tasks, providing better results in audio.

APC has also been used for LLM-based music understanding with great success. Music-Understanding LLaMA [LHSS23b] adopts a pre-trained MERT encoder and a LLaMA LLM. The system is trained with APC with only the music understanding adaptor from the MERT embedding space to the LLaMA embedding space being unfrozen, and all other elements of the system frozen. LLaRK [GDSB23b] uses Jukebox [DJP+20b] as an audio encoder and a fine-tuned LLaMA2 model to generate responses to textual queries. One specificity of LLaRK is that it augments training data with extracted musical features such as chord sequences, Key, tempo, and tags. GPT-4 is then used to create question-answer pairs using the extracted features to bolster the size of the training data for LLaRK. M$^2$UGEN [HLSS23a] generalises music multimodal understanding with LLMs beyond audio and text to video and images. M$^2$UGEN also employs a frozen LLaMA model to generate textual responses. For this approach, live modules include adaptors from frozen MERT, ViT, and ViViT encoders for various modalities, as well as a text-audio token adaptor trained from scratch to generate audio tokens for either MusicGen or AudioLDM2. Specifics of M$^2$UGEN include special understanding/generation tokens appended as prefixes to generate responses to distinguish between audio-generative and text-generative tasks. Finally, recently, in the field of Symbolic Music, ChatMusician [YLW+24a] makes use of LoRA adaptors on a pre-trained LLaMA model to generate answers including ABC notation as a second language for the model in the training set. It is clear through these studies that a key consideration for multimodal LLMs that use pre-trained textual LLMs for music understanding is the adaptation of encoded audio tokens to the LLM vocabulary, which might include LoRa adaptors, from-scratch external adaptors, or fine-tuning the LLM itself.

*3) Masked Modeling:* The main idea behind masked (language) modelling (MM/MLM) pre-training is to mask a portion of the input data and train the model to predict the original content of the masked parts using the remaining context. This is usually done with Transformer models, which naturally deal with sequential data and can effectively capture long-range dependencies. Practically, input sequences are represented as a series of tokens, with different tokenisation strategies existing, ranging from learnable embeddings to a token codebook learned from a different model (see Subsection IV-C). Depending on whether the tokens are continuous or discretised, the masked modeling objective is typically formulated as a regression loss between the predicted and true tokens or a classification loss between the predicted probabilities and the true tokens, respectively. As will be described in the following paragraphs, this loss is often used in conjunction with others, such as a contrastive loss.

*a) Continuous Masked Modelling:* Continuous Masked Modelling approaches generally operate on frequency representations of audio, learning continuous (i.e. not discretised) "tokens" of the input. At their core, they are inspired by the Vision Transformer (ViT) [DBK+20], which divides an input image in fixed-size 16x16 patches, flattens them and linearly projects them to 768 dimensions to obtain patch embeddings, adds positional embeddings to maintain spatial information, and processes them through a standard Transformer encoder used for supervised pre-training. ViT was later adapted to the Audio Spectrogram Transformer (AST) [GCG21] for audio classification, which demonstrated the necessity of image pre-training to achieve good performance. For that reason, SSAST [GLCG22b] proposed a self-supervised adaptation of AST by introducing two two-layer MLPs at the Transformer output. The first, a reconstruction head, tries to predict the original masked patch and, as is now typical in continuous masked modelling, uses the Mean Squared Error (MSE) loss, defined as $\mathcal{L} = (r - x)^2$, where $x$ is the true patch and $r$ is the reconstructed patch. The other, a classification head, tries to match the correct spectrogram patch for each masked position, which is a token-wise contrastive loss (see Subsection IV-A1). The losses are summed, with the reconstruction loss being weighted ten times the classification loss. SSAST also proposes a clustering method for choosing which patches to mask with controllable granularity, rather than randomly masking them, which aims to guide modelling at both a local and a global level.

Many audio-based approaches followed replacing the reconstruction head with a transformer decoder and using high masking ratios (between 70% and 85%). Among other design choices, these approaches defer in the exact patch masking strategy, which patches are fed to the encoder, and what loss is used. MaskSpec [CWZZ23] and MSM-MAE [NTO+22a] use a simple encoder-decoder architecture, where both masked and unmasked embedded patches are given to the encoder. MAE-AST [BPH22] uses a similar regression and contrastive loss to SSAST, but only unmasked tokens are given to the encoder. The encoder outputs are then given to the decoder together with the masked embedded patches. This approach results in a significant training speedup and memory usage reduction. Audio-MAE [HXL+22a] only uses unmasked patches in the encoder and the MSE loss, but also shows improvements by using a local attention mechanism that groups patches into local windows for decoding. More recently, Dasheng [DYW+24] scales up the MAE paradigm to 1.2 billion parameters and 272,356 hours of audio, using densely extracted mel spectrogram frames as input and learnable position embeddings, achieving several state-of-the-art audio classification results.

Some other approaches further tweak the patch masking

Table III: Music Foundation Model Trained with Masked Language Modeling.

| Model | Modality | Application | Training Paradigm | Tokenizer | Architecture |
|---|---|---|---|---|---|
| MAE-AST | Audio (Speech&Sound) | Understanding | MLM | Spectrum | Transformer Encoder Decoder |
| Audio-MAE | Audio (Speech&Sound) | Understanding | MLM | Spectrum | Transformer Encoder |
| SSAST | Audio (Speech&Sound) | Understanding | MLM | Spectrum | Transformer Encoder |
| Beats | Audio (Sound) | Understanding | MLM | Spectrum | Transformer Encoder |
| DiscreteBERT | Audio (Speech) | Understanding | MLM | vqwav2vec | Transformer Encoder |
| WavLM | Audio (Speech) | Understanding | MLM | 1-D CNN | Transformer Encoder |
| w2v-BERT | Audio (Speech, Audio, Music) | Understanding | MLM, Contrastive Learning | Spectrum | Transformer Encoder |
| ampNet | Audio (Music) | Generation | MLM | Discrete Tokens (DAC) | Transformer Encoder Decoder |
| MidiBERT-Piano | Symbolic (REMI) | Understanding | MLM | REMI, compound word | Transformer Encoder |
| MusicBERT | Symbolic (MIDI) | Generation | MLM | MIDI (OctupleMIDI) | Transformer Encoder Decoder |
| MRBERT | Symbolic (MusicXML) | Generation | MLM | MusicXML Note Event, Compound Word | Transformer Encoder Decoder |
| EAT | Audio (Sound) | Understanding | MLM (Online Distillation) | Spectrum | Transformer Encoder |
| A-JEPA | Audio (Speech&Sound) | Understanding | MLM (Online Distillation) | Spectrum | Transformer Encoder |
| data2vec | Audio (Speech) | Understanding | MLM (Online Distillation) | 1-D CNN | Transformer Encoder |
| MT4SSL | Audio (Speech) | Understanding | MLM, MLM(Online Distillation) | 1-D CNN | Transformer Encoder |
| data2vec 2.0 | Audio (Speech) | Understanding | MLM (Online Distillation) | 1-D CNN | Transformer Encoder |
| M2-Duo | Audio (Speech, Audio, Music) | Understanding | MLM (Online Distillation) | Spectrum | Transformer Encoder |
| music2vec | Audio (Music) | Understanding | MLM (Online Distillation) | 1-D CNN | Transformer Encoder |
| MuLaP | Audio (Music), Text | Understanding | MLM | 1-D CNN | Transformer Encoder |
| JMLA | Audio (Sound), Text | Understanding | MLM (Online Distillation) | Spectrum | Transformer Encoder Decoder |
| MusIAC | Symbolic (REMI), Text | Generation | MLM | REMI | Transformer Encoder Decoder |
| AV-HuBERT | Audio (Speech), Image | Understanding | MLM | 1-D CNN | Transformer Encoder |

strategy and architecture used. Audio-GMML [AAAK23] uses a Group Masked Model Learning (GMML) strategy [AANK23] to mask patches that are deemed significant in the input, and combines intermediate representations from the encoder with its outputs as an input to a shallow decoder. MW-MAE [YTHT24] uses a multi-window, multi-head attention module in the decoder aiming to capture local and global context levels simultaneously. Finally, ASiT [AAW+24] employs GMML and self-distillation with a teacher-student architecture, using a reconstruction loss alongside a local and global similarity contrastive loss.

More recently, interest has been shown in predicting latent representations rather than reconstructing the original input, particularly to encourage modelling of more abstract, higher-level features. data2vec [BHX+22] proposed this concept with a self-distillation setup where a teacher transformer encoder model processes the full, unmasked input, while a student model with an identical architecture learns to predict the teacher's latent representations from a masked view using MSE. After each student model update, the teacher's weights are updated using an Exponential Moving Average (EMA) of the student's weights. data2vec demonstrated the potential of this approach for self-supervised learning in any data modality. MAP-Music2Vec [LYZ+22] adapted this framework specifically to music, using a waveform input with a 1D convolutional embedder for both the teacher and student transformer models. They pre-train various versions of the model with different input length, mask span, and masking probability with 130k hours of music audio, and probe various transformer layers on clip-level music tasks, achieving comparable performance to other approaches.

A similar latent prediction architecture was proposed in JEPA [ADM+23] for images, which was later adapted to audio [FFH24]. JEPA contains a so-called context encoder and a target encoder, similar to a student and teacher setup, respectively. A large part of the image or spectrogram is given

to the context encoder, while the target encoder sees the whole input. A predictor network, which is a lighter-weight ViT, then attempts to reconstruct the target encoder output via the context encoder output and positional conditioning. Again, the target encoder weights are updated through an EMA of the context encoder weights. The audio adaptation of JEPA also introduces a curriculum masking strategy, where masking in pre-training is initially done randomly, but gradually time-frequency aware masking takes over, which should be a harder task for the model. Finally, M2-Duo [NTO+24] uses a similar setup to JEPA in the audio domain, but masking is also applied to the target encoder input, opposite to that applied to the context encoder input, which they claim encourages the teacher to also model the input well and produce better encodings.

Masked modelling has also been adapted to a multimodal setting. AV-MAE [GFI+23] proposes a simple audio-video adaptation using just the reconstruction loss. They experiment with different encoder and decoder architectures, including how cross-modal fusion happens, and investigate finetuning applications, particularly unimodal ones. Specifically, they identify and experiment with *early fusion*, where audio and vision embeddings are concatenated before being passed to the transformer, *separate* processing, where the two streams are encoded separately and are fuse before being decoded (*late fusion*), *shared fusion*, in which weights are shared between the two encoders, and *mid-fusion*, in which the output of two separate encoders is fused and fed to a joint encoder. pre-training with random initialisation, they find minor differences between fusion strategies for audio, although larger ones are identified for video and audiovisual applications. Those were favored by parameter sharing in the decoder and mid- or late-fusion in the encoder, possibly because they allowed stronger coupling between the modalities.

CAV-MAE [GRL+23] proposes an audiovisual approach using mid-fusion, with separate audio and video encoders

proceeded by a joint encoder. Three partially masked streams are passed through the joint encoder: an audio encoding, a visual encoding, and a concatenated audiovisual encoding, all stemming from patch representations of the inputs. The inter-modal contrastive loss is then computed between mean average pooled embeddings of each single-modality stream. The encoder output of the joint encoding is decoded with a joint decoder, with the reconstruction loss computed between its outputs and the original input patches. MaViL [HSX$^+$24] follows a similar mid-fusion encoding setup. In addition to using the inter-model contrastive loss, it also uses an intra-model contrastive loss between different views of the input. pre-training is first done by training MaViL with the input reconstruction objective for each modality, and then the model is used as a teacher to a student that learns to predict the aligned, contextualised audiovisual representations produced by the teacher.

Finally, MuLaP [MBQF22b] proposes a music audio multimodal approach using weakly-aligned captions. Specifically, it follows a similar setup to ViLBERT with one branch for audio and one for language, which are encoded separately using a CNN encoder and a BERT-style tokeniser and encoder respectively. The outputs are then fed to two co-attentional layers, in which the key and value vectors are exchanged between modalities, allowing each modality to attend to relevant information in the other modality. In addition to a reconstruction and a classification loss for the audio and text masked modelling respectively, MuLaP also employs an audio-text matching objective by processing both true and fake audio-text pairs.

*b) Discrete Masked Modelling:* Many music modelling approaches work with discretised representations. This is the case with symbolic music but also with several notable music foundation models that operate on audio by transforming it into a sequence of discrete tokens. Several approaches to discretised masked modelling have been proposed, including for audio, differing in their discretisation method, as well as whether the input or only the prediction targets are discretised. Unlike continuous masked modelling, which typically uses a regression loss, discrete masked modelling often employs a classification loss like the categorical cross-entropy:

$$\mathcal{L} = -\sum_{c=1}^{C} y_c \log(p_c),$$

where $C$ is the number of classes (discrete tokens), $y_c$ is the true label (1 if the sample belongs to class $c$, 0 otherwise), and $p_c$ is the predicted probability of class $c$ for the sample.

The first approach to BERT-like discretised masked modelling for audio is HuBERT (Hidden-Unit BERT) [HBT$^+$21c], applied to speech processing. The input to HuBERT's transformer encoder is still continuous, being a waveform encoded through CNN layers. The encoder outputs are then used to predict a discretised target. During the first iteration, k-means is run on MFCC features, and the assigned cluster IDs are used as the targets. The loss is then computed between predictions and assigned cluster IDs only for masked regions. During subsequent iterations, k-means is instead run on features

from an intermediate layer of the transformer. The authors also demonstrate performance improvements with a multitask learning setup, in which clustering ensembles with multiple k-means models are used to create multiple codebooks corresponding to different granularities. WavLM [CWC$^+$22b], with the same architecture, effectively improves performance by increasing the data scale and mixing noise in the input during pre-training. Fast-HuBERT [YMZ$^+$23] accelerates the training of HuBERT using Fbank features with a lower sampling rate and simplifies the HuBERT loss to cross-entropy loss. HuBERT-AP [RLW$^+$22] and CTCBERT [FWGL23] change how the Transformer output is aligned with the unsupervised pre-training targets, making the pre-training process more like the practice of ASR task. PBERT [WWW$^+$22], MonoBERT [MZY$^+$23], and Poly-BERT [MZY$^+$23] introduce phoneme-like information to the SSL pre-training to improve performance.

w2v-BERT [CZH$^+$21b] uses a similar setup to HuBERT, but proposes and end-to-end training process. It utilises a learnable wav2vec 2.0 quantisation module using the features encoded by the CNN as input and trained on the contrastive loss to the context vectors (outputs of the context encoder, which in this case is a conformer). The context vectors are passed through another conformer that predicts the token ID, using the target IDs learned by the quantisation module in the same iteration. DiscreteBERT [BAM20], on the other hand, shows that discretisation of the input, not just the targets, can be more effective. Specifically, they use a pre-trained vq-wav2vec model to quantise the input, which is subsequently masked, as well as to create the targets, and perform standard BERT pre-training. They also try using k-means cluster assignments from MFCCs and filterbank features, but vq-wav2vec significantly outperforms them.

More recently, EncodecMAE [PRF24] combines the Hu-BERT approach with the EnCodec quantiser and a masked encoder-decoder setup. Specifically, it uses EnCodec's encoder and RVQ layer with 8 out of its 32 codebooks to produce discrete targets. The MAE, operating on spectrogram frames rather than patches to increase its time resolution, is then trained to predict those targets. EnCodecMAE is shown to perform well across audio modalities, including music. BEATs [CWW$^+$22] proposes a different setup, where a random codebook is initialised, and on the first iteration projected patches are matched to the nearest vector in the codebook. Their assigned ID is then used as a target for masked modelling. In subsequent iterations, the model acts as a teacher, producing prediction targets, with the tokeniser trained to predict them using cosine similarity.

A few approaches have been proposed specifically in the music domain. VampNet [GSKP23] operates on discretised inputs and targets using the Descript Audio Codec (DAC). DAC compresses the 44.1kHz music audio input to a bitrate of 8kpbs using 14 codebooks, 4 of which are for coarse information and 10 for detailed. VampNet then uses two transformers, one for each token type, to learn to predict the masked set of tokens conditioned on the other token type. MERT [LYZ$^+$24] instead proposes a multitask approach, trying to predict discrete RVQ EnCodec tokens from 8 codebooks and

trying to reconstruct continuous CQT frames with masked modelling. It uses a continuous input of 24kHz waveforms encoded through a CNN. As an alternative to EnCodec tokens, MERT experiments with k-means clustering of log mel spectrograms and chroma features, however those generally perform slightly worse. It also shows that mixing short audio excerpts from the same batch to an input as an augmentation can improve the robustness of the model to common audio perturbations. Finally, another approach [WHL24b] investigates various discrete token prediction setups using BEST-RQ [CQZ+22a], a method where the projection matrix and codebook are randomly initialised and not updated during training. BEST-RQ provides a simpler alternative to discretisation that doesn't require a separate training phase, while performing comparably when enough training data is used.

While symbolic music can express continuous values (e.g. pitch bending between notes, or gradual dynamics changes such as crescendos), it's effectively usually treated as a discretised representation with a given level of rhythm and pitch quantisation. As such, discrete masked modelling has prevailed for symbolic music. MusicBERT [ZTW+21], aimed at symbolic music generation, uses the BERT setup with cross-entropy loss with a custom symbolic tokenisation called OctupleMIDI. The tokenisation creates separate tokens for time signature, tempo, bar number, position within the bar, instrument, pitch, duration, and velocity, which they show performs better than other tokenisations at both phrase- and song-level tasks. Instead of random masking, they mask tokens of the same type that belong to a bar, which they claim encourages the model to learn better representations than it would with adjacent token prediction. Another approach [CCC+24] uses the same setup as MusicBERT, but uses the more established REMI and CP symbolic representations rather than OctupleMIDI, and tries to establish a symbolic piano music classification benchmark on public data, rather than the private dataset used to train MusicBERT.

MRBERT [LS23b] follows a similar setup, but treats the melody and rhythm separately by using a "semi-cross attention" mechanism to exchange information between the two streams. This is done by allowing the query from one stream to interact with the other through a softmax operation, and then proceed to attend to the key-value pairs of the original stream. MRBERT is then finetuned for autoregressive, conditional, and seq2seq generation. Finally, MusIAC [GSK+22], focused on music infilling, uses an adapted version of the REMI tokenisation with multiple levels of control tokens such as track-level and bar-level controls. It uses BERT-like pre-training, but with an encoder-decoder transformer setup.

Finally, a few approaches have applied discrete masked audio modelling to a multimodal setting. AV-HuBERT [SHLM22] proposes an audio-visual adaptation with speech data using cluster assignment targets. The two modalities are initially encoded separately, with the encoded features then concatenated and passed through a joint transformer encoder. During the first iteration, MFCC features from the audio are clustered, and the assignment IDs are used as prediction targets. During subsequent iterations, features from an intermediate layer of the preceding model are instead used.

The authors note the model's over-reliance on speech audio during pre-training, and counter it with two main methods: firstly, while visual data is encoded with a ResNet, audio data is projected with a fully connected network, forcing the audio encoder to learn simple features, and, secondly, during pre-training, an entire modality is randomly dropped out, with a likelihood of 0.5. u-HuBERT [HS22b], TESSP [YRC+22] and VATLM [ZZZ+23] further extend this approach by introducing a text stream as well as fusion techniques to allow pre-training with examples that don't have all modalities available.

### B. Music Domain Adaptation for Foundation Models

At the core of foundation models' versatility lie their few-shot, zero-shot, and in-context learning capabilities. Supervised finetuning paradigms enable foundation models to perform tasks without needing extensive task-specific training or instruction data, enabling dealing with use cases in the music domain that traditionally have limited data available. In this subsection, we will discuss finetuning methods for music FMs such as prefix tuning and adaptor tuning, followed by discussion on the zero-shot and in-context learning capabilities of FMs. We note that limited work in the music domain investigates the impact of projection layers like MLP, cross-attention, and Q-formers in prefix tuning, as well as prompting for zero-shot learning, such as chain-of-thought.
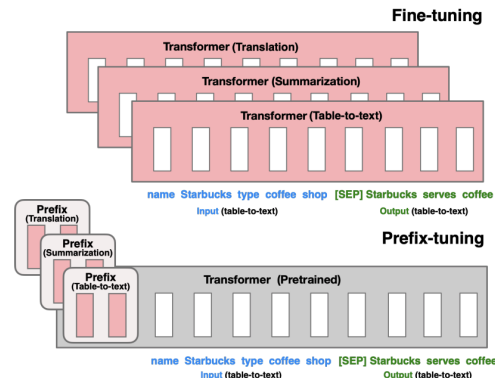


Figure 7: The top subfigure demonstrates the finetuning that updates all parameters shown in red. The button one shows prefix-tuning, which only optimises the parameters of input prompts or prefix blocks[LL21].

*1) Prefix Tuning & Prompt Tuning:* Various approaches, including prompt tuning [LARC21] and prefix tuning [LL21], have been developed for large FMs to enable tuning on a relatively small set of downstream tasks or for modality alignment. Prefix tuning involves adding a series of trainable vectors (prefixes) to the input, thus effectively preparing the model's context for task-specific processing without changing its underlying parameters. This approach allows for a focused tuning of the model's behaviour while retaining its extensive pre-training knowledge. It can also be used to connect an image or music encoder to the LLM to enable it to understand information in other modalities. On the other hand, prompt tuning refines model performance by crafting and tuning task-specific prompts that guide model response generation. This

Table IV: Music Foundation Model Trained with Multimodal Adapters.

| Model | Modality | Application | Training Paradigm | Tokenizer | Architecture |
|-------|----------|-------------|-------------------|-----------|--------------|
| Qwen-Audio | Audio (Speech, Sound & Music), Text | Understanding | prefix tuning, GPT | 1-D CNN | Transformer Encoder Decoder |
| LLaRK | Audio (Music), Text | Understanding | prefix tuning, GPT | Pre-trained model (CLAP, Jukebox) | Transformer Decoder |
| Musilingo | Audio (Music), Text | Understanding | prefix tuning, GPT | Pre-trained model (MERT) | Transformer Decoder |
| MU-LLaMA | Audio (Music), Text | Understanding | adapter tuning, GPT | Pre-trained model (MERT) | Transformer Decoder |
| M2UGen | Audio (Music), Image, Text | Both | adapter tuning, GPT | Pre-trained model (MERT) | Transformer Decoder |
| SALMONN | Audio (Sound & Speech), Text | Understanding | adapter tuning, GPT | Pre-trained model (Whisper, BERT) | Transformer Decoder |
| LTU | Audio (Sound), Text | Understanding | adapter tuning, GPT | Pre-trained model (Whisper) | |

approach exploits the inherent ability of LLM to generate context-appropriate responses based on the cues provided in the prompts and, therefore, requires minimal data for effective learning.
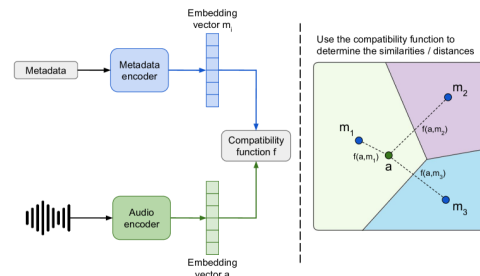
MusiLingo [DML+24] and Llark[GDSB23b] seamlessly integrate pre-trained music encoders into LLMs using prefix tuning. MusiLingo uses a single projection layer to connect a pre-trained music model, MERT, to a Llama, while Llark uses CLAP and Jukebox to Llama2, enabling efficient task-specific adaptation through prefix adjustment. These allow such models to excel at generating music subtitles and processing music instruction, demonstrating a scalable approach to augmenting LLM in music multimodal applications. While these approaches use an MLP projection that aligns different modalities, some other visual-language models such as BLIP-2[LLSH23] and Flamingo[ADL+22] also use Q-former and cross-attention for the modality alignment, respectively.

*2) Adaptors and Full Tuning:* Model adaptation [FVRA21], [Csu17] involves adapting a pre-trained model to perform better in specific, often narrower, domains and tasks. Adaptors[HGJ+19], [ZHL+23a] have been widely adopted in language modelling and other modalities like computer vision. For example, MINERVA [LAD+22] and Galactica [TKC+22] use continued pre-training to adapt existing LLMs to the domains of science and math. Unlike prefix tuning or prompt tuning that freezes or only finetunes the initial transformer layers of LLMs, adaptors generally modify the deeper transformer layers.

In the realm of music foundation models, an example of domain adaptation is when a model trained on a wide variety of music is adapted to specialise in particular types of music with, for example, unique rhythms or instrumentation. On the other hand, task adaptation usually takes place when a model is adapted to solve new tasks, often narrower than the original. One such example demonstrates the use of a music tagging model, which hopefully captures rich semantic information about various musical aspects, to other music classification and regression tasks [CFSC17]. In terms of finetuning PLMs in the music domain, MertTech [LMW+24] uses multi-task finetuning to improve instrument playing technique detection in world music, of which limited examples existed in the pre-training dataset. For music-related multimodality, Mu-llama [LHSS24] utilises llama-adapter [ZHL+23a] to align music recordings and their text captions. SALMONN [TYS+24b] utilises a window-level Q-Former to convert variable-length sequences from speech and audio encoders output into variable numbers of augmented audio tokens for input to the Vicuna

LLM and to finetune the LLM's parameters using Low-Rank Adaptation (LoRA)[HWAZ+21], allowing for a range of audio-speech and music-understanding tasks to be performed. Its music comprehension capability even exceeds that of a model designed for music [WMB+24].

*3) Zero-Shot Learning & Instruction Tuning:* Zero-shot learning [RPT15], [XSA17] takes the concept of minimal examples even further, allowing models to perform tasks without having been explicitly trained on them. In music and audio, supervised learning has been well studied to tackle different classification tasks such as scene classification, environmental sound classification, etc. However, the current supervised learning techniques require large amounts of annotated data from target sound classes. It can become expensive for humans to collect sufficient labelled data. Therefore, zero-shot learning has been favoured in many realistic scenarios without human annotation. As shown in Fig 8, traditional zero-shot classification in audio leverages semantic embeddings to enable sound recognition without direct training examples. This method maps pre-trained audio features and textual descriptions into a unified semantic space, utilising acoustic embeddings from VGGish and textual embeddings from pre-trained language models like Word2Vec, GloVe, and BERT [CLPN19], [XRV21]. By evaluating compatibility between these embeddings, the approach classifies sounds effectively, even for untrained classes, significantly improving classification performance.



Figure 8: Zero-shot learning with PLM[TTG+24].

Instruction tuning is typically used to fine-tune language models on a wide variety of tasks, enabling the language model to generalise to new tasks with only instructions and no examples. FLAN (Finetuned Language Net) [WBZ+21] utilises the concept of instruction tuning to enhance the zero-sample learning capabilities of language models. By pre-training LLMs with finetuned 137B parameters on more than 60 NLP tasks, FLAN demonstrates significant improvements

in handling unseen tasks. The approach organises NLP tasks into clusters and selectively trains the model on some tasks while retaining others for evaluation. This approach allows FLAN to outperform larger models such as 175B GPT-3 on a variety of datasets, demonstrating its efficacy in generalising instruction tasks without prior direct contact. Supernatural Instruction [WMA$^+$22] introduces a comprehensive benchmark of 1,616 NLP tasks of 76 different types, providing a powerful platform for assessing the generalisation capabilities of NLP models in an instructional setting. The tasks are described through expert-written natural language instructions across multiple languages and task formats, enabling detailed evaluation of models such as Tk-Instruct. This smaller-scale Transformer model outperforms traditional large-scale models by being trained specifically for these different instruction sets, revealing the potential of targeted instruction-based training to achieve superior cross-task generalisation. Instruction Tuning in music such as ChatMusician [YLW$^+$24b] represents a significant advancement in the application of FMs to music. ChatMusician, through continuous pre-training and supervised finetuning of text and ABC notation, has gained the ability to perform complex music-generation tasks such as generating music based on given chord progressions, keys, motifs, and music structures, along with understanding music theory at the zero-sample level, effectively surpassing the musical capabilities of LLMs such as GPT-4. Prompt tuning with multi-tasking has also been shown to help with zero-shot learning[SWR$^+$21], while prompt writing techniques such as chain-of-thought [WWS$^+$22] have remained underexplored.

*4) Few-shot Learning & In-context Learning:* Few-shot learning refers to the ability of a model to learn and adapt to new tasks by observing only a small number of examples. As shown in Fig 9, traditional few-shot algorithms evaluate the pre-trained features of all examples and merge them into an anchor embedding for classification. Few-shot learning has been widely adopted in fields such as natural language processing and computer vision. In the context of audio and music, few-shot learning enables foundation models to understand different properties of audio and music to perform certain tasks on sound event detection, transcription and source separation [WSBB20], [WSC$^+$20], [WSBB22].
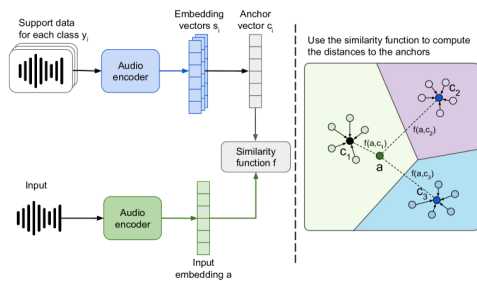


Figure 9: Few-shot learning with PLM [TTG$^+$24].

Recently, in-context learning (ICL) was proposed, in which a model uses exemplars in the input to make predictions or generate outputs for unseen tasks [BMR$^+$20], [DLD$^+$22]. This approach has not yet been fully explored in music. ICL typically requires a generative model capable of handling interleaved textual and musical data, which is a way of combining instructions and musical examples in a format in the input sequence. However, the development of generative pre-trained models that effectively handle this interleaved musical text format remains scarce, with both AnyGPT and ChatMusician lacking the nuanced understanding required to interpret and generate complex musical content based on contextual cues.

### C. Audio Tokenisers

The length of a music waveform sequence is typically long. A music recording may have a sampling rate of up to 48 kHz and several minutes or hours long. Such long sequences pose challenges for training machine learning systems. To address this issue, researchers use tokenisers to compress audio signals into latent representations with shorter sequence lengths than raw audio signals. Audio tokenisers can be categorised from different perspectives as follows.

*1) Hand-crafted versus Learning-based Tokenisers:* Before the advent of learning-based tokenisers, researchers developed hand-crafted tokenisers for audio processing. The most representative tokeniser is the Short-Time Fourier Transform (STFT) [GL84]. An audio signal is denoted as $x \in \mathbb{R}^L$, where $L$ is the number of samples. The audio signal $x$ is split into short-time frames. Then, a Fourier transform is applied on each frame to obtain the *spectrogram* with a shape of $\mathbb{R}^{T \times F}$, where $T$ is the frames number and $F$ is the frequency bins number. Later on, the *mel spectrogram* was proposed to compress the STFT by applying the mel scale [SVN37]. The mel spectrogram has a shape of $T \times M$ where $M$ is the mel frequency bins and there is $M \ll F$. Mel spectrogram has been widely used in audio pattern recognition tasks such as automatic speech recognition. Based on mel spectrograms, researchers proposed MFCCs features (c.f. subsection II-B). The advantage of hand-crafted features is that they have explicit explanations - for example, MFCCs decouple speech contents and the timbre of speakers. However, hand-crafted features may lead to information loss. Recently, learning-based representations have been proposed to address audio understanding and generation problems.

*2) Continuous Audio Tokens :* to include more info Audio tokens can be continuous or discrete. Continuous tokens are usually used for audio understanding and diffusion model-based audio generation tasks. For example, continuous tokens can be wav2vec [SBCA19], wav2vec 2.0 [BZMA20b], HuBERT [HBT$^+$21a], and HIFI-codec [YLH$^+$23]. The continuous tokens usually have a shape of $\mathbb{R}^{T \times C}$, where $C$ is the channels of the representation. Later on, modified discrete cosine transform (MDCT) methods inspired from audio compression algorithms were proposed [DVE$^+$23]. FunCodec is a frequency-domain codec that can achieve higher compression ratio than time-domain codecs. Recently, semantic audio codecs such as [LXY$^+$24] have been proposed to introduce both semantic and acoustic information into audio representations by using AudioMAE and k-means clustering algorithms.

Many codecs are trained in an unsupervised way. We denote the input to a mask-based regression system as $X$

and the mask as $M$, where $M$ has the same shape as $X$. The mask-based regression system takes $(1 - M) \odot X$ as input to predict $(1 - M) \odot X$ $M \odot X$. The represented works include SSAST [GLCG22b], Audio-MAE [XLB$^+$22], MAE-AST [BPH22], MSM-MAE [NTO$^+$22b], Audio-MAE [XLB$^+$22], MAE-AST [BPH22], BEATs [CWW$^+$22], and EAT [CLM$^+$24]. The advantage of regression-based training is that they can reconstruct fine-grained audio signals. On the other hand, the contrastive training strategy uses semantically rich discrete label prediction to encourage the tokens to have high-level semantic information and discard redundant details of audio signals. For example, Wav2Vec [SBCA19] and Wav2Vec 2.0 [BZMA20b], HuBERT [HBT$^+$21a], w2v-BERT [CZH$^+$21a], MERT [LYZ$^+$24] apply contrastive pre-training for automatic speech recognition and music recognition. Recently, SemantiCodec [LXY$^+$24] has been proposed by using a k-means clustering algorithm to extract audio semantic information and use a diffusion model to reconstruct audio. BEST-RQ [CQZ$^+$22b] applies a random projection quantiser to train audio tokenisers in an unsupervised way. MT4SSL [MZT$^+$23] combines both k-means clustering offline targets from HuBERT and the teacher model's online targets from data2vec as the pre-training pretext task.

*3) Discrete Audio Tokens:* Discrete tokens are commonly used for language model-based audio generation tasks. For instance, in VQ-VAE [VDOV$^+$17b], the encoder outputs discrete tokens. These discrete tokens can be input to language models such as Transformers to predict the next discrete token for audio generation. In SoundStream, a residual vector quantizer (RVQ) [ZLO$^+$21] is introduced to represent audio signals in a hierarchical structure. Discrete tokens usually have a shape of $\{0, 1, ..., D - 1\}^T$, where $D$ is the vocabulary size of the tokeniser. Discrete audio tokens have the advantage of representing audio signals with a finite vocabulary. Therefore, discrete audio tokens can be seamlessly integrated with large language models in the natural language processing domain. DiscreteBERT [BAM19] applies a VQ-wav2vec vocabulary to extract audio tokens. RVQGAN has been proposed for music generation in [KSL$^+$24]. VampNet [GSKP23] applied Descript audio codec (DAC) to extract tokens for music generation. MusicGen [CKG$^+$24] proposes a simple and controllable music generation system by using Encodec.

*4) Symbolic Music Tokeniser:* Symbolic music tokenisers encode music information into a structured and symbolic format, such as MIDI [YCY17], [CCC$^+$21], MusicXML [Goo01], ABC notation [GMP20], Humdrum [Hur02], Lily-Pond [NN03], and Octopus [LLY$^+$20]. Symbolic music datasets have the advantage of containing compact information of music for editing than audio representations. Symbolic music representations high-level musical information that is useful for tasks such as music theory analysis and algorithmic music generation. Music Byte Pair Encoding (BPE) [FGCB23] compresses symbolic tokens into a latent space to further reduce the latent dimension. Those representations have been used in symbolic music generation systems such as ChatMusician [YLW$^+$24a].

## D. Model Architectures

After tokenising the audio signals, the representations or tokens are input into audio understanding or generation models. Many of these systems have an encoder-only, decoder-only, or encoder-decoder architecture. The audio encoders can be trained separately on audio tasks or jointly trained with the decoders. We describe the encoders and decoders as follows.

*1) Encoder:* An audio encoder is used to extract representations or tokens of audio signals. We describe the architecture details in this section. For the tokeniser of decoder-only models such as Encodec for MusicGen, please refer to the previous subsection.

*a) Convolutional Neural Networks (CNNs)-based Audio Encoder:* Convolutional neural networks have been widely used in audio pattern recognition [CFS16], [HCE$^+$17], [KCI$^+$20]. We denote the audio signal as $x \in \mathbb{R}^L$ and its time-frequency representation such as log mel spectrogram as $X \in \mathbb{R}^{T \times F}$. CNNs consist of convolutional layers, downsampling layers, and pooling layers to extract high-level features fromthe mel spectrogram. The CNNs have the advantage of extracting both time-domain patterns and time-frequency domain patterns. CNNs have been widely used in automatic speech recognition, music tagging, audio tagging, and source separation [CMJG17]. A variety of CNNs, including ResNet and EfficientNets [VBV22] have been used to improve the vanilla CNN architectures. CNN-based encoders also include Musicset Unsupervised Large Embedding (MULE) [MKO$^+$22] and universal audio representations [WLW$^+$22]. The advantage of CNNs trained on time-frequency domain representations is that they are relatively easier to train compared to time-domain systems and Transform architectures.

Researchers have also investigated time-domain-only audio encoders. The time-domain audio encoders are directly applied on the waveform $x \in \mathbb{R}^L$ without transforming it into a time-frequency representation such as a spectrogram. The time-domain encoders have the advantage of extracting non-harmonic time-domain features. For one-dimensional CNNs, time-domain convolutional layers are applied to the waveform to extract high-level features. To reduce sequence length, time-domain CNNs use dilated convolution layers to capture a larger receptive field, such as WaveNet [VDODZ$^+$16]. SampleRNN uses hierarchical recurrent neural networks to split long sequences into short sequences [MKG$^+$16]. The advantage of time-domain encoders is that they can be used to reconstruct audio signals without estimating the phases of spectrograms. Time-domain encoders have been widely used in source separation such as ConvTasNet [LM19] and 1D convolutional LSTMs [DUBB19]. Audio codec algorithms such as SoundStream [ZLO$^+$21] also apply convolutional layers to extract latent representations from waveforms.

*b) Transformer-based Audio Encoder:* The Transformer [VSP$^+$17b] is an architecture that utilises self-attention mechanisms to capture contextual relationships in input sequences. The core part of a transformer is self-attention: $Attention(Q, K, V) = Softmax(\frac{QK^T}{\sqrt{d_k}}V)$, where $Q$, $K$, and $V$ are query, key, and value with shapes of $N \times d_k$. The variable $N$ is the sequence length and $d_k$ is the dimension of

the tensor. Since the successful applications of Transformer in natural language processing and computer vision, Transformers have also been widely used in audio processing. In audio tagging, Transformer-based systems usually take the time-frequency representation of an audio file such as a spectrogram $X \in \mathbb{R}^{T \times F}$ as input. Audio spectrogram transformer (AST) [GCG21] proposes to split a spectrogram into $16 \times 16$ patches. The patches are linearly projected to a sequence of 1-D embeddings. The embeddings are inputted to a transformer for classification. HTS-AT [CDZ+22] is a hierarchical audio transformer with semantic tokens for sound event classification and detection. To reduce the sequence length of embeddings, efficient training of Transformers with patch-out [KSEZW21] was proposed. SpecTNT [LWW+21] is a frequency-dependent architecture. In each SpecTNT block, a Transformer is used to extract frequency-dependent features into the tokens. In source separation, band-split Transformer [LWKH24] was proposed to extract frequency-dependent features as input to Transformers.

*2) Decoder:*

*a) Linear Projection:* is typically applied to the encoder output to predict pre-training targets like tags or tokens in encoder-only models. Early works of decoders apply linear projection on the latent embeddings, such as PANNs [KCI+20], and HTS-AT [CDZ+22]. Multiple linear projection (MLP) has also been used as a decoder in audio tagging tasks. We denote the embedding extracted from the audio encoder as $h$ and the output of the decoder as $y = f_{\text{dec}}(h)$. The advantage of applying linear projections is that they are lightweight and require little data to train or finetune. The projection has been widely used in audio tasks with simple output formats, such as audio tagging.

*b) Transformer.:* A linear projection is not enough to map from audio representations to a desired task. To address this problem, Transformer-based decoders were proposed as decoders to map from the audio representations to desired tasks and support flexible output formats. Inspired by the T5 architecture [RSR+20a] in natural language processing, SpeechT5 [AWZ+21] applies a Transformer-based symmetric encoder and decoder to autoregressively decode audio representations to the outputs of a variety of tasks. Sparse attention [ZTS21] methods such as Deepspeed [ARA+22] have been proposed to reduce the computational cost in Transformers. Long attention models such as block-wise self-attention [QML+19] and local attention [DDHB19] have been proposed to model long sequences.

*c) State-Space Models:* have also been used to model audio signals. In [GGDR22], a SaShiMi system was proposed to introduce state into audio representations. SSamba [SDJM24] is an unsupervised and attention-free system that uses states. Audio Mamba [ESFC24] investigated using state space models to replace self-attention for audio classification. Those state-based systems have the advantage of removing the quadratic computation complexity in Transformer models.

Recently, decoder-only architectures such as LLaMA [TLI+23a] have been proposed to address the natural language processing and multimodality problems. The decoder-only systems interleave the audio representations with nat-

ural language. The representative works include Listen, Think, and Understand (LTU) [GLL+23], Pengi [AC87], Salmonn [TYS+23], LauraGPT [CCG+23] and Qwen-Audio [CXZ+23], where automatic recognition systems and audio tagging systems are used as an encoder, and a language model is used to predict the outputs according to the input questions. Decoder-only architectures have the advantage of allowing using arbitrary interleaved modalities as input and output.

*E. Interpretability & Controllability on Music Generation*

From a musical perspective, the *goal* of controllable music generation is *to make the music in **certain locations** follow **certain features** during the generation*. The location of music is multidimensional – it usually contains time-axis, stem-axis, and pitch-axis. For example, a specific location can be "notes within the first four bars in the piano track between C2 and C5", or "the last ten seconds of the drumset". The features of music include both information *of* music and information *about* music. Information about music is usually described by intrinsic music language, such as rhythm, pitch, chord etc., while information about music is usually described by natural language, such as "a pop style", "a happy emotion", "syncopated with the melody track", or even the synesthesia of a particular visual scene.

To achieve the goal, most machine-learning methods follow a conditional generation approach, modelling $P(target|control)$ where $target$ is the music within the specified location and $control$ is the information based on which based on the *methods* of control includes to *add*, *subtract*, and *edit* notes or features of certain locations.

Interpretability could lead to better controllability of foundation models. For an end-to-end model, the controllability can be added when the music control coincides with the mathematical properties. These mathematical constraints can be used to guide the gradient descent path[LGW18], to evaluate the sampling process [DJGD21], or constrain the learned distribution [HPN17]. These can be used to control the factors that we can mathematically define, but the performance differs in cases. For end-to-end models, another approach is to use prefix control which defines tokens, such as MuseCoCo[LXK+23b]. All these show that control can only be applied to those concepts that can be explicitly defined.

In the interpretability-oriented approach, one of the common approaches is to use representation learning models to learn a latent space of implicit music concepts, such as pitch contour, accompaniment texture [WZZ+20b] and timbre [LKJX21]. These concepts are usually hard to define by rules, but we encode them in the latent space. It has been shown that these learned representations can be recombined for style transfer, and we can produce variations or new music via sampling or interpolating the latent space. We have also shown these latent codes can be used in longer-term prediction, and infilling [ZX21], [WX22]. As another approach, interpretability can be achieved by defining an interpretable workflow. For example, whole-song [WMX24] defines a general hierarchical music language, so that the generation process is interpretable with respect to the defined workflow. Accomontage3 [ZXW23] uses

an interpretable architecture to generate two layers of latent codes in order to achieve a multi-track arrangement. Under current-generation fashion, the control of chords and rhythm is still external, and current methods cannot be learned well. The interpretability of such a concept is for future research.

On the other hand, we believe that aligning these external controls with implicit music knowledge in a non-interpretable foundation model helps enhance its interpretability. In recent years, we have seen significant progress in text-to-music generation models. These models generate symbolic or acoustic music based on given text descriptions. Many of these models leverage language models or diffusion models. For instance, MusicGen [CKG+23a] stands out as an exemplary audio text-to-music model combining a T5 encoder [RSR+20b] for text descriptions, a pre-trained Encodec [DCSA22] as a compressor for music audio signals, and an acoustic transformer decoder for generating Encodec tokens. Riffusion [FM22a] represents a notable diffusion-based audio text-to-music model, employing a UNet-based stable diffusion [RBL+21] framework. Additionally, recent prominent models like MuseCoCo [LXK+23c], MusicLM [ADB+23], and MusicLDM [CWL+23a] demonstrate promising results in symbolic and acoustic text-to-music generation.

All these models are trained in a supervised manner on extensive datasets comprising pairs of textual descriptions and corresponding music. This training enables them to generate symbolic or acoustic music based on textual inputs during inference. Despite their impressive capabilities, these advanced models remain black boxes, making it challenging to extract embedded musical knowledge from them or determine their understanding of musical concepts.

The opacity of these models presents significant challenges for interpretability. Understanding how these models implicitly translate textual descriptions into musical concepts and subsequently produce music remains an unsolved task. However, recent developments in content-based controllable music generation models offer a promising solution to enhance interpretability. This approach involves explicitly defining musical contents and fine-tuning the models to generate music aligned with these predefined content-based controls.

This methodology assumes that large models, when generating music, employ musical elements as internal intermediaries; in other words, interpretable prior musical knowledge exists hidden within the large models. By finetuning these models to generate music in a flexible and controlled manner based on diverse musical elements, researchers aim to leverage the interpretability of these models in music generation tasks.

A notable related work is the emergence of Parameter-Efficient Fine-Tuning (PEFT) methods, including Low-Rank Adaptor (LoRA) [HWAZ+21] and LLaMA adapter [ZHL+23b], which provides efficient ways to adjust pre-trained large models via tuning just a few parameters. These methods offer cost-effective ways to manipulate large pre-trained generation models, paving the way to transforming a standard text-to-music model into a content-based controllable generation model.

Based on PEFT techniques, recent works such as Coco-Mulla [LXZJ24] and AirGen [LXJZ23] focus on generating music from chords, drums, piano rolls, and text inputs. Similarly, MusicControlNet [WDWB24] considers rhythm, melody dynamics, and textual cues to generate music. These efforts demonstrate a promising way to improve the model's interpretability, facilitating more nuanced and controlled music generation based on explicit musical controls. Meanwhile, they also suggest the growing trend toward enhancing the controllability of music generation through interpretable model adjustments.

### F. Foundation Models as Music Agents

*a) What is an agent?:* Broadly speaking, an agent is a kind of entity that has desires, beliefs, intentions, and the ability to take actions [ZNAP95]. In the realm of AI, an agent is an artificial entity that can set objectives autonomously or semi-autonomously, perceive the environment, work out plans, reason and make decisions, and take actions to achieve these objectives, powered by the technologies of artificial intelligence [XCG+23]. In terms of music agents, they refer to AI agents dedicated to music and related domains, which can understand, generate, and convert music in symbolic, acoustic, and other modalities.

The development of AI agents has undergone a long process. Early attempts at building AI agents focus on improving the specific capabilities of an agent, such as sensing the environment, symbolic reasoning, taking actions with reinforcement learning, or handling specific tasks. However, the keys to AI agents are the general capabilities of reasoning, planning, perception, and action. While focusing on specific aspects of AI agents can result in steady progress, it might not be able to make fundamental enough improvements to push AI agents to practical usage.

Large language models (LLMs) that emerged in recent years have revolutionised the era of language and multi-modality and brought strong capabilities in reasoning, planning, perception, and generation, making them a good fit as the core module of AI agents. Some pioneer investigations such as AutoGPT [Teab], HuggingGPT [SST+23a], and Visual ChatGPT [WYQ+23] leverage LLMs to understand user requests, reason and plan according to objectives, decompose complicated tasks into subtasks, revoke external tools for task execution, and response generation. On the one hand, LLM-based agents extend the capabilities of LLMs to handle complicated and/or multimodality tasks. On the other hand, LLMs empowered AI agents to a new level with strong general capabilities that have never been possible using traditional AI agent technologies. Due to the great potential of LLM-powered AI agents, a lot of research efforts have been made in this area, such as improving the reasoning and planning capabilities [WWS+22], [Teaa], benchmarking AI agents [SST+23b], enhancing tool use [SDYD+23], [YSC+24], designing agents to handle challenging tasks [SST+23a], and applying AI agents to more domains, including audio [LZL+23], [LZL+24] and music [YSL+23]. In the next paragraph, we will review the typical AI agent in music and audio domains.

*b) Examples of music and audio agents:*

- MusicAgent [YSL+23], the first LLM-powered AI for music, integrates diverse models and an autonomous workflow to address various music tasks like generation, transcription, and conversion. It simplifies the complex process for professionals and amateurs by analyzing requests, decomposing them into subtasks, and invoking external tools to fulfil these tasks.

- AudioGPT [HLY+23] focuses on audio modality and leverages a large language model to process different audio modalities (speech, music, sound) and handle different audio understanding and generation tasks. For music tasks, it supports singing voice synthesis by calling external music models.

- SpeechAgents [ZLW+24] is a multimodal multi-agent system based on LLMs to simulate human communication. In SpeechAgents, each agent leverages a multimodal LLM as the decision centre and uses multimodal signals to exchange with other agents.

- Loop Copilot [ZMX+23] introduces an innovative system combining large language models with specialised AI music models to streamline the collaborative creation of music loops. This system utilises a conversational interface for dynamic, iterative music editing, and a global attribute table to ensure consistency throughout the creation process. Besides, it is not limited to creating music based on vague text inputs but allows for fine-grained musical edits, including adding or removing tracks and making localised adjustments to modes and tempos. This capability enhances the system's utility in detailed music production tasks.

- ComposerX [DYY+24] introduces a novel multi-agent framework for polyphonic notated music composition, utilising the reasoning power of large-scale language models as well as extensive knowledge of music history and theory. This approach produces high-quality, coherent compositions better than traditional single-agent systems, and requires no specialised training or services, making it a cost-effective alternative.

- ByteComposer [LLD24] pioneers a human-like melodic composition process using a four-step agentic framework: conceptual analysis, draft composition, self-assessment, and aesthetic selection. ByteComposer combines the interactive and knowledge-understanding capabilities of LLMs with symbolic music modelling to achieve performance comparable to that of a human composer and is extensively validated through professional feedback.

*c) Future Work of Music Agents:* Although there are some preliminary investigations on music agents, there are still a lot of space to improve. We introduce some future work of music agents bellow.

- Improve reasoning and planning ability. The key capability of the music agents is enabled by the foundation model behind it. When handling complex tasks, the reasoning and planning ability of the foundation model plays a critical role in understanding the user request, decomposing the whole task into subtasks, organizing complicated working flows, and calling suitable models and tools for task execution [HC23], [HGM+23]. Thus, improving the reasoning and planning capability has always been the top priority in building music agents as well as general artificial intelligence.

- From semi-autonomous to full-autonomous agents. Nowadays, most music agents are semi-autonomous, requiring a human in the loop, such as initiating user requests, setting system configurations, providing suitable model and tool sets, and strong interactions in the agent workflows. While, Autonomous agents based on LLM are popular research directionss[HGM+23], [WMF+24], [YLL+24], [YLY+24]. In order to play more important and critical roles in helping musicians and music consumers, music agents should evolve from semi-autonomous to full-autonomous, gradually reducing the degree of user interference when fulfilling complex user requests.

- Versatile music agents supporting multiple modalities and tasks. Music does not appear alone but usually engages with other modalities, such as text in lyrics and comments, tags/taxonomy in genres and styles, symbols/MIDI in music scores, and image/video in album covers and music videos. Modelling music with other modalities and tasks together could greatly extend the scenarios that music agents can support. This either requires the foundation model to be a powerful multimodal model or the external tools and models to cover multiple modalities and tasks.

## G. Scaling Laws for Music

Scaling laws [KMH+20], [HKK+20], [HBM+22], [AAA+23] characterise how the performance of large language models (LLMs) (usually measured with a cross-entropy loss on a held-out set) scales as a power law of model size, dataset size, and training compute. Scaling laws can predict the loss of LLMs with a given number of model parameters and data size (the number of training tokens, i.e., batch size × training steps) and can determine the optimal configuration of model size and number of training tokens given a computational budget. There are different formulations for scaling laws, such as OpenAI's scaling laws [KMH+20], Chinchilla scaling laws [HBM+22], and data-constraint scaling laws [MRB+24]. OpenAI's scaling laws are more precise, while Chinchilla's laws are simplified based on some assumptions. For simplicity[14], we briefly describe Chinchilla scaling laws as follows:

$$\mathcal{L}(N, D) = \frac{A}{N^\alpha} + \frac{B}{D^\beta} + E, \qquad (2)$$

which describes the loss $\mathcal{L}$ with regard to the model size $N$ and the number of training tokens $D$. It consists of three terms:

---

[14]Strictly speaking, Chinchilla scaling laws have obvious errors. For example, given two models with model size $N_1$ and $N_2$ where $N_1 > N_2$, when increasing the same amount of training tokens $\Delta D$, the loss reductions of the two models are the same according to Chinchilla laws as shown in Equation 2. However, it is evidently not the case in practice since the larger model $N_1$ will have more loss reduction than $N_2$. But for ease of explanation, we just use Chinchilla scaling laws as an example here.

leftmargin=2em

- Functional approximation error ($\frac{A}{N^\alpha}$), which describes a perfectly trained model with $N$ parameters underperforms the ideal generative process [HBM$^+$22]. As $N$ is increasingly large, this term will approach 0.
- Convergence error ($\frac{B}{D^\beta}$), which reveals that the model is not trained to convergence, as we just use a certain number of training steps and training tokens. As $D$ is increasingly large to cover the true data distribution, this term will approach 0.
- Minimal achievable loss ($E$), which corresponds to the entropy of the natural text and characterises the loss for an ideal generative process on the true data distribution.

In Equation 2, $\alpha$ and $\beta$ describe the effectiveness of the model size and data size in reducing loss. A larger $\alpha/\beta$ means the model/data is more effective in loss reduction. $A$ and $B$ are constants that are dependent on many factors such as data corpus, tokenisation, and vocabulary size, and do not have a fundamental meaning. Estimating the parameters $(A, B, E, \alpha, \beta)$ of the above scaling law usually requires collecting some data points $(N, D, \mathcal{L})$ and fitting the predicted loss and observed loss using some criterion and optimisation algorithms (e.g., Huber loss [COR30] and L-BFGS algorithms [Noc80]).

There are several typical uses of scaling laws: 1) predict the loss of LLMs with a given model size and training tokens; 2) allocate the optimal model/data size given a compute budget; 3) determine the relative scaling of model and data with an increased compute budget; 4) determine the ratio of data to model given a compute budget.

Traditional scaling laws are mostly studied in language language models for text [KMH$^+$20], [HKK$^+$20], [HBM$^+$22], [AAA$^+$23], [MRB$^+$24], which are not applicable for music. MuPT [QBM$^+$24] is the first work to study the scaling laws for music, with dedicated design on symbolic music data representations and scaling laws in the data-constraint and over-fitting scenarios. While the great benefits of MuPT in guiding the modelling scaling of symbolic music generation, there is still a lack of scaling laws for music understanding and generation in the audio domain. Later, some studies emerged investigating scaling laws for multimodality [HKK$^+$20], [AYC$^+$23] and especially for speech-language model [CM24]. Since multimodal data such as images, videos, and speech are high-dimensional and complex and not usually modelled by language models, it is not straightforward to study the scaling laws on multimodal data. Therefore, they usually tokenise multimodal data into tokens and model them as the next token prediction. In this way, they could study the scaling laws similarly to those in text-based large language models.

In [CM24], two different data representations for speech are studied: 1) semantic tokens extracted by HuBERT [HBT$^+$21a]; 2) semantic BPE tokens which are obtained by applying byte-pair encoding [SHB16] on the semantic tokens from HuBERT. The scaling laws show that given the same amount of tokens, semantic BPE tokens could consume more model capacity than semantic tokens. This is because semantic BPE tokens are semantically richer and have a larger coefficient $\beta$ as shown in Equation 2, which means by increasing the same amount of tokens, semantic BPE tokens could cause more loss reduction.

In this way, in order to better coordinate with the data and reduce the loss, the model size should increase. An intuitive explanation is that as the data is semantically richer, it should consume more model capacity to process them. This also aligns with our intuition: in the audio domain, the vocoder or codec model usually generates the waveform data with either auto-reconstruction or from mel-spectrogram, which processes less semantic information, so the vocoder or codec model size is usually small. However, the language models that predict semantic tokens are usually very large.

Since there is not much research on scaling laws for music in the audio domain (although the scaling laws from multimodal and especially speech domain could provide much insight into the audio music domain, they are not directly applicable), and the emergent ability of music FMs is under-explored, we call for actions to proactively study the scaling properties of different data representations and model paradigms for audio music understanding and generation in the future. Such research topics include, but are not limited to, which type of downstream capability of large models can be inferred by the small models, how the properties of audio tokenisers relate to the parameters of scaling law, and what the minimal number of model parameters or training data is for a specific type of downstream tasks etc.

### H. Additional Future Improvements

In this subsection, we will shortly discuss potential research topics related to foundation models for music. Other than the foundation model techniques discussed in previous subsections, the following paragraphs will include domain knowledge, long sequence modelling and causal modelling. For discussions on fairness, transparency, and potential issues from model bias, please refer to subsection VI-A2.

*1) Domain Knowledge:* The model architecture described in this section so far, apart from small modifications, can be construed as 'standard' in the sense that it has been explored for a plethora of other tasks both somewhat linked with music (e.g. speech) as well as entirely different from it (e.g. machine translation). Although such an approach is quite appealing [HBL12], the limited data circumstances and inadequacies of architecture/optimisation in the music domain have led to a substantial interest in incorporating domain knowledge. This reversal towards an appreciation of domain expertise has brought about a revival in domain-specific engineering of machine learning approaches in many areas, including music [Ser13]. There are quite a few disciplines from which domain knowledge for music foundation models can be sourced. Acoustics, for example, gave rise to spectrum-based representations, which continue to find their use in deep learning architectures [DJP$^+$20a]. The same applies to psychoacoustic chroma tokens [Mï5] which can assist deep learning should it fail to automatically infer the importance of timbre information from the data [LYZ$^+$24]. Music theory provides a further rich source of information which has inspired numerous approaches and will possibly continue to do so.

One key challenge here is how to integrate music theory and related concepts into existing standard architectures rather

than designing custom architectures to fit one or more music theories, as was done before [Wid98] and still remains popular in the area [HOB23], or uncover new ones by means of automatic theoretical analysis [YV16]. One popular direction for integrating music theory into the Transformer and related architectures is input modification. Although a simple approach, it provides opportunities for encoding a wide range of music-theoretical information. So far, the work in this direction has focused on encoding lower-level relative attributes, such as pitch interval, duration and onset [GKH23], and some higher-level attributes, such as bars [WERA16]. There remains to be an extension of this work to other higher-level relative attributes. The same applies to absolute attributes such as melody, harmony and rhythm. Methods that do not rely on input modifications to incorporate music theory are interesting but likely challenging due to changes required to the core Transformer implementation as was the case with the work of [HVU$^+$18b] on relative positional encoding.

Another key challenge lies in bridging the gap between approaches more naturally formulated at the symbolic token level and approaches suitable to be used with acoustic tokens. As in other areas of structured audio processing, any attempt to split audio into the underlying units (e.g. phones in speech or chords in music) is a highly non-trivial endeavour requiring expert annotators to act under the uncertainty of determining boundaries among those units. In the case of music, this is further compounded by the need to decide on the level of some of those units (e.g. degree of loudness). Thus, approaches capable of automatically inferring segmentations or marginalising over all possible segmentations would need to be developed.

*2) Long-Sequence Modelling:* The length of a typical music recording (several minutes) and its resolution (48 kHz and above) firmly suggest that music belongs to a class of long-sequence modelling problems. For example, a 3-minute recording with a resolution of 48 kHz packs approximately 9 million samples. The need to handle long sequences has inspired a large body of approaches, which can be grouped into those that modify existing architectures and those that propose novel architectures that are better suited to handling long sequences.

Approaches that modify the widely popular Transformer primarily focus on the two most important shortcomings of this architecture when handling long sequences: the complexity of attention and positional encoding. The quadratic complexity of the standard attention mechanism has been primarily tackled using sparsification in structured and unstructured forms leading to a drastic decrease in complexity with minimal or no loss in performance [CGRS19], [DYY$^+$19], [GQL$^+$19], [YGG$^+$19], [BPC20], [KKL20], [QML$^+$20], [ZGD$^+$20], [XOG$^+$22]. The standard (absolute) positional encoding seemingly appropriate for short sequences found in many speech and natural language processing tasks has been extended to various relative positional encoding schemes [SUV18], [DYY$^+$19], [RSR$^+$20b], [WLQ$^+$22] and hybrids of absolute and relative encodings [SAL$^+$24], [SDP$^+$23]. Although many of these approaches have been integrated into Transformer-based foundational music models [HVU$^+$18a], [YLW$^+$22], their ability to handle truly long music recordings remains

to be seen. In this context, the recent work on long-context speech recognition [FR23], where context length reached 1 hour of speech, suggests this to be a promising direction.

Recently there has also been substantial interest in developing novel approaches with linear or near-linear complexity with respect to sequence length. Among them, structured state-space models [GGR22] appear to yield promising results for sequences reaching 1 million tokens [GD23]. These models can be interpreted as a combination of recurrent and convolutional architectures which additionally integrate principled mechanisms for modelling long-range dependencies. Another popular line of research explores hierarchical architectures where sequences are segmented into shorter units [DYY$^+$19], [PZV$^+$19], [RPJ$^+$20]. This allows making use of separate mechanisms for modelling short-term (e.g. attention) dependencies within segments and long-term (e.g. implicit [BKKB24] or explicit [PZV$^+$19] recurrence) dependencies across segments. Some of these approaches showed promising results on synthetic tasks with sequences reaching 2 million tokes.

*3) Causal Modelling:* The discussion about domain knowledge earlier in this section stopped short of mentioning perhaps one of the strongest kinds of domain knowledge - causal relationships. Such relationships we cultivate all of our lives. In childhood, many of us learned the impact of key pair variations on our perception of sound by tinkering with the piano keys. As we grow older and start watching movies we begin to judge choices made by sound directors in particularly dramatic scenes. The knowledge of these relationships and our ability to enact them enables us to conduct interventions (e.g. fix a key, increase tempo) and create counterfactual pieces (e.g. imagine a known piece in a new key). Do machine learning engineers need to know the causal relationships that the pianist Horowitz formed and his ability to manipulate them in order to automatically generate Horowitz-quality piano pieces?

Even in simpler scenarios than music, determining causal relationships is a non-trivial endeavour. In many cases this is hardly a possible task due to the lack of access (e.g. astronomy) or understanding (e.g. neurology). In other cases, the full set of causal relationships is not known or available. Latent variables (e.g. hidden Markov models [Rab89], auto-encoders [Kra91]) emerged as a powerful but implicit surrogate aimed at capturing some of these relationships. However, the lack of appropriate parameter estimation methodologies capable of learning causal rather than merely correlational relationships meant that latent variable models have so far achieved limited success in this area. Causal learning emerged in response to these deficiencies of modern machine learning. Starting from the work done by Pearl and colleagues [Pea09] it has now emerged into a new field full of exciting developments [KLL$^+$22]. Despite much progress made in causal learning over the years, work done in music and related areas remains extremely limited. The key issue stems from the lack of data with labelled causal information and/or methods capable of automatically extracting causal relationships. Although some causal models have emerged for simple image generation tasks, where causal factors could be limited to shape, texture and background [SG21], it is not clear how to extend them to

other areas where causal relationships are much less trivial to specify, extract, and model.

An interesting alternative trend recently emerged where large non-causal models are manipulated into exhibiting causal relationships [LPP+20], [MBAB22], [OWJ+22], [WWS+22], [LPV+23]. One prominent example is the work of Wei and colleagues on chain-of-thought prompting for eliciting factually correct reasoning of large language models in question-answer type tasks [WWS+22]. Unlike standard prompting [BMR+20], where an example question is supplemented with a short factual answer (e.g. giving a numerical answer, such as 11, to the question about how many cans of food are left after a certain number of them go missing and some restocked), the chain-of-thought prompting supplies an answer that showcases how a human could arrive to that conclusion [LYDB17]. This approach has been shown to be effective in arithmetic, common sense and reasoning tasks. Whether similar prompt manipulations can be extended to music foundation models remains to be seen.

## V. DATASETS & EVALUATION

### A. Datasets

Data plays a central role in training the underlying model. The diversity, quantity and quality of the training data contribute greatly to the generalisability and robustness of the model[XPD+24], [SKB+24], [ZQL+24]. For machine learning and computer music researchers, training music-based models may require hundreds of billions of tokens or even larger datasets, and selecting high-quality datasets of different music recordings is a great challenge. In the realm of computer music, datasets are categorised into composition-level, typically sheets represented in MusicXML and ABC notation, performance-level symbolic music exemplified by MIDI formats, and datasets comprising raw audio waveforms. This section explores extensive open-source music datasets that are beneficial for training LLMs or LDMs tailored for music applications. Additionally, we will discuss Python libraries that are pivotal for music processing. Subsequently, the focus will shift to multimodal music datasets. Finally, the evaluation methodologies for foundation models developed on these datasets will be examined.

*1) Composition-level Symbolic Music Datasets:* In this subsection, we focus on introducing symbolic music datasets that contain more than 1,000 music excerpts. We will exclude smaller datasets though they are potentially useful for evaluation or supervised fine-tuning in specific domains. For more information on music datasets, readers can refer to the datasets website page of the ISMIR community [15]

*a) MusicXML:* Besides representing traditional Western Classical music scores, MusicXML is widely utilised for storing lead sheets of popular music, encompassing the melody, lyrics, harmony, and various markings of a song. The melody is notated in modern Western music notation, with lyrics presented as text beneath the score, and harmony indicated by chord symbols positioned above the score. Notably, lead

sheets provide little details regarding instrumentation or accompaniment.

The Wikifonia Lead Sheet Dataset [LRL17], originally from the now-defunct public repository Wikifonia.org, includes 5,533 MusicXML lead sheets of Western music across various genres. Before the service ceased in 2013, Lim et al. secured 2,252 lead sheets in major keys with typically one chord per bar. If a bar contained multiple chords, only the first was chosen. The dataset is converted and distributed to CSV format on the website[17]. Besides ABC notation, it also provides a version of MusicXML and MIDI files.

To complement the Yamaha Signature MIDI Collection, Jeong et al. developed the MuseScore Lead Sheet Dataset [JKK+19] by collecting MusicXML scores from MuseScore, a community-driven music score web platform, along with including their own transcriptions. These scores were transcribed voluntarily by community members. Although the dataset includes just 226 MusicXML pieces by 16 composers, it corresponds to 1,052 piano performances in MIDI format. Collectively, these pieces comprise a total of 666,918 notes in MusicXML and 3,547,683 notes in MIDI.

Hooktheory Lead Sheet Dataset (HLSD) [YHF+21] include 11,329 music lead sheet samples collected from the TheoryTab forum on Hooktheory's website, a resource specializing in music education software. These samples, whose melodies are transcribed by users with corresponding chord progressions in high quality, are denoted by both chord symbols and functional labels. The dataset included an excessive 704 chord class. Due to copyright restrictions, only song snippets are shared on TheoryTab, with annotations such as structural and genre labels.

*b) ABC Notation:* The ABC format is a plain text method for documenting music. Besides representing single-line melodies, it has the capability to fully represent polyphonic classical scores.

The Irish Massive ABC Notation (IrishMAN) dataset [WLYS23] is a comprehensive collection of 216,284 Irish tunes in ABC notation. It is primarily sourced from an online traditional Scottish and Irish collection [18] and the official website of ABC notations[19], To achieve uniform formatting, all tunes were converted from ABC to XML and back to ABC using automated scripts, with non-musical fields such as titles and lyrics removed to maintain focus on musical content. Additionally, each tune in the dataset is annotated with control codes for musical forms and structures.

Further, an open-source website has Western folk music in ABC format such as the Nottingham Music Dataset (NMD), which is comprised of 1,200 British and American folk songs. A recently refined version is available online [20]. Another example is the Henrik Norbeck collection that includes more than 2,800 ABC scores with lyrics, for Ireland and Sweden [21]. More information can be found in the dataset section of the music generation survey [JLY20].

[15]https://www.ismir.net/resources/datasets/.

[17]http://marg.snu.ac.kr/chord_generation/
[18]https://thesession.org/
[19]abcnotation.com
[20]https://ifdo.ca/seymour/nottingham/nottingham.html
[21]http://www.norbeck.nu/abc/

Table V: Open-source music dataset for pre-training.

| Dataset | Modality | n files | Description |
|---|---|---|---|
| Wikifonia | MusicXML | 2,252 CSV samples | CSV of MusicXML from Wikifonia.org |
| MuseScore Lead Sheet Dataset | MusicXML, MIDI | 226 piece with 336k notes | Derived from MuseScore website |
| Hooktheory Lead Sheet Dataset | MusicXML | 11,329 lead sheet samples | Derive from TheoryTab music theory forum [16] |
| IrishMAN | ABC, MIDI, MusicXML | 216,284 | Scottish & Irish folk songs |
| Nottingham Music Dataset | ABC notations | 1,200 | Online corups of British & American folk songs. |
| ABC tune book of Henrik Notebook | ABC notations | 2,800 | Irish & Swedish folk songs |
| Lakh MIDI Dataset | MIDI | 176,581 files | Mainly pop and rock music |
| Yamaha Signature MIDI Collection | MusicXML, MIDI | 1.4k | Piano performance, mainly Romantic pieces |
| DoReMi | Image, MusicXML, MEI, MIDI | 6k | Steinberg's Dorico |
| ADL piano dataset | MIDI | 11,086 | Pop, classical and jazz piano pieces |
| Symphonies | MIDI | 46,359 files, 650 hours | Classical symphony, multi-instruments |
| NES-MDB | MIDI | 5,278 | NES games BGM |
| MAESTRO | MIDI, audio | 1.2k files | Classical Piano |
| GiantMIDI-Piano | MIDI, audio | 10,855 pieces, 1237 hours | Machine transcribed classical piano |
| Meta-MIDI | MIDI, audio | 436,631 MIDI files | 10M match to Spotify music tracks |
| Free Music Archive (FMA) | audio | 106,574 tracks, 8.2k hours | Collected from FMA website |
| MTG-Jamendo | audio | 55,701 tracks, 3.8k hours | Collected from Jamendo website |
| Music4ALL | audio | 109,269 tracks, 911 hours | Collected from YouTube |
| Million Song Dataset (MSD) | audio feature | 1,000,000 | 1M pop song provided by Echo Nest |
| AudioSet | URL of audio | 1,011,305 music clips | 2,084,320 clips including general audio |
| AcousticBrainz | audio feature | 2,524,739 | 2M audio features with MusicBrainz metadata |
| Disco-10M | feature & URL of audio | 15,296,232 | 10M features with diverse genres and artistis |

*2) Performance-level Symbolic Music Datasets:*

*a) MIDI:* There are many MIDI datasets and we only include the largest and the widely used portion.

The Lakh MIDI Dataset (LMD)[Raf16] is a large MIDI polyphonic music corpus with inconsistent expressive characteristics, consisting of various genres, instruments, and periods of time. The "LMD full" version includes the whole 176,581 MIDI files after removing duplication, and the "LMD matched" subset refers to the 45,129 files that are aligned with items in the Million Song Dataset (MSD).

The Yamaha Signature MIDI Collection[22] comprises 1.4k high-quality solo piano MIDI performances recorded during an international competition for junior pianists. Predominantly featuring late Romantic pieces by Chopin and Liszt, along with Mozart sonatas, this dataset is essential for research in performance generation. A corresponding MusicXML version is available through the MuseScore Lead Sheet dataset, or the (note-)Aligned Scores And Performances dataset (nASAP) [PCCF+23].

DoReMi dataset [SF21] has around 6k pages of the score images, with the corresponding MusicXML, MEI, and MIDI scores. Its primary objective was optical music recognition (OMR) but can easily be used in other tasks. The MIDI is exported using Steinberg's Dorico [23] software.

The Augmented Design Lab (ADL) Piano MIDI dataset [FLW20] is derived from the LMD by retaining only one of the multiple versions of each song and only tracks featuring piano-family instruments identified by MIDI program index 1-8, resulting in 9,021 unique rock or classical piano MIDI files. To enhance genre diversity with additional styles like jazz, they incorporated 2,065 files sourced from Internet sources,

culminating in a total of 11,086 pieces in the final dataset.

Symphonies dataset [LDC+22] is a large corpus of symphonic music from various online sources. The collection consists of 46,359 MIDI files, primarily symphonic works, with a total duration of about 650 hours.

The NES-MDB dataset [DMM18] encompasses 5,278 songs from the soundtracks of 397 NES games, involving compositions from 296 composers and containing over two million notes. This dataset is extracted from the assembly code of NES games, capturing precise timings and parameters for authentic chiptune renditions. The NES synthesiser includes five instrument voices—two pulse-wave, one triangle-wave, and one noise generator—with the complexity of its audio synthesis chip abstracted for researchers. NES-MDB is available in multiple formats to facilitate research in NES music, with additional details provided for those interested in deeper technical exploration.

Additionally, there are many online corpora for MIDI files such as BitMidi[24] Classical Archives[25], and FreeMidi[26].

*b) MIDI-audio dataset:* Apart from what has been mentioned above, there are multiple music transcription or generation datasets that not only include symbolic music but audio as well.

MIDI and Audio Edited for Synchronous TRacks and Organization (MAESTRO) dataset [HSR+18] consists of about 200 hours of professional piano performances including both MIDI annotations and audio recordings.

The GiantMIDI-Piano dataset [KLCW20] is the largest classical piano dataset, encompassing 10,854 MIDI files by 2,784 composers, totalling 1,237 hours in duration. Such

---

[22]http://www.yamahaden.com/midi-files
[23]https://github.com/steinbergmedia/DoReMi

[24]https://bitmidi.com/
[25]https://www.classicalarchives.com/
[26]https://freemidi.org/

collection is transcribed using a high-quality, open-sourced piano transcription system [KLS+21b]. ATEPP [ZTR+22] is a similar transcribed classical piano dataset but with an emphasis on classical performance variety. It contains 11,674 performance tracks by 49 virtuoso pianists.

MetaMIDI Dataset (MMD)[EP21] is a large dataset of 436,631 MIDI files with metadata. It also includes a subset with 168,032 files that are matched to the audio clips in the MusicBrainz database.

*3) Acoustic Music Datasets:* The Free Music Archive (FMA) [DBVB16] is an expansive dataset that comprises 343 days of audio. It includes 106,574 tracks from 16,341 artists and 14,854 albums, organised into 161 genres. The dataset offers high-quality audio, audio features, and comprehensive metadata at the track and user levels, along with free-form text like artist biographies. Notably, a significant portion of the FMA consists of experimental music, which markedly differs from typical music used in downstream tasks. Therefore, it is advisable to exclude experimental tracks during pre-training to optimise relevance and applicability.

The MTG-Jamendo [BWT+19] Dataset is compiled from music available on Jamendo²⁷, and features tags applied by the uploaders. This dataset contains 55,701 full audio tracks, each at least 30 seconds long and encoded at 320kbps MP3, totalling 509 GB of audio. The average track length is 224 seconds, summing up to 3,777 hours of audio. The tracks are annotated with 692 tags, derived from genres, instruments, and mood/theme categories. To streamline the dataset, variant spellings and semantically identical tags were consolidated, resulting in 99 remapped tags.

Music4ALL[SPD+20] database encompasses 109,269 songs from 15,602 anonymous users along with their listening histories, each represented by 30-second audio clips, lyrics, and multiple metadata attributes like genres and tags.

Aside from the aforementioned, there are many large-scale datasets which only provide the feature or URL of the audio, allowing users to collect the tracks from the internet. Some such examples include the Million Song Dataset (MSD) [BMEWL11], AudioSet [GEF+17], AcousticBrainz [PBK+15], and Disco-10M [LGFW24] etc.

*4) Software for Music & Audio Processing:* Developing foundation models for music necessitates leveraging open-source software for tasks ranging from preprocessing and feature extraction to training the models and evaluation in both symbolic and audio domains. In symbolic music processing, libraries such as `mido` and `pretty_midi` offer extensive capabilities for MIDI manipulation, whereas `note_seq`, supported by Google's Magenta team, provides advanced functionalities for processing and training symbolic music. For audio processing, tools like `librosa`, `Essentia`, and `madmom` are crucial for their comprehensive audio analysis and feature extraction capabilities.

Importantly, for audio I/O operations, the use of `torchaudio` with the `sox_io` backend is advised due to its superior speed and performance compared to alternatives like the `soundfile` backend.

²⁷https://www.jamendo.com/

*5) Multimodal Datasets: Present & Challenges:*

*a) Text-Audio:* LP-MusicCaps-MSD[DCLN23b] originates from the ECALS subset[DWCN23] of the Million Song Dataset[BMEWL11] and includes 520k 30-second music clips. It features a diverse vocabulary of 1,054 labels spanning various music-related categories, such as genre, instrument, vocal attributes, mood, theme, and culture. On average, each clip is tagged with 10.2 labels. The 3 captions are generated by GPT-3.5 for each audio clip, involving one caption, one summary, and one rephrased version.

The Song Describer Dataset (SDD) [MWD+23a] introduces a curated dataset of 1.1k natural language descriptions paired with 706 music recordings, all licensed under a Creative Commons license, designed to provide a powerful tool for music and language models. This dataset helps evaluate three key tasks: music captioning, text-to-music generation, and music language retrieval, highlighting the critical role of cross-dataset evaluation. Each description in SDD contains rich musical characteristics such as genre, mood, and orchestration to provide semantically and syntactically comprehensive captions. SDD is unique in that it provides longer audio track segments and multiple descriptions of each track, thereby enhancing the robustness of the assessment. This dataset, combined with metadata from MTG-Jamendo, not only enriches the field of music research but also sets new standards for realistic and sustainable dataset practices in the field of music language.

MusicQA [LHSS24] is a dataset of music question-answer pairs, generated by MPT-7B model [Tea23]. The dataset is created by using MusicCaps and MagnaTagATune that contain music captions or tags related to descriptive and inferential questions about the music. Each audio has five open-ended question-answer pairs focusing on aspects like music emotion, tempo, and genre. It include 12,542 music clips in the training set, 76.15 hours in total, with 112,878 question-answer pairs.

MusicInstruct [DML+24], also derived from the MusicCaps dataset's music-caption pairs, contains Q&A pairs generated with few-shot learning techniques by GPT-4. It includes 5,521 YouTube ID of audio clips, a short version with 27,540 Q&A pairs focused on detailed aspects such as emotion, instruments, vocals, tempo, and genre, typically eliciting concise one or two-sentence responses, as well as a long Version that contains 32,953 Q&A pairs with broader questions about musical pieces, often resulting in more elaborate, paraphrased responses.

MusicBench [MGG+23] is a text-to-music dataset, also based on the MusicCaps dataset, that consists of 5,479 samples. MusicBench is structured into various training and testing sets: the initial sets (TrainA and TestA) are further enhanced by splicing four control sentences that describe music features in chord, key, beat and tempo information into the original prompts to create TrainB and TestB. Subsequently, TrainB is rephrased using ChatGPT to generate TrainC. Additionally, an audio augmentation strategy that adjusts pitch, speed, and volume is implemented to improve data diversity. A comprehensive training set after augmentation consists of 52,768 text-audio pairs.

The Multimodal Album Reviews Dataset (MARD) [OEAL+16] includes metadata and customer reviewers'

Table VI: Open-source multimodal music dataset.

| Dataset | Modality | n files | Tasks |
|---|---|---|---|
| LP-MusicCaps-MSD | audio URL, text | 520k audio, 1.5M text | music captioning. |
| Song Describer Dataset (SDD) | audio, text | 706 audio, 1.1k | music captioning, text-to-music, retrieval |
| MusicQA | audio, text | 12,542 clips, 112,878 Q&A | acoustic music instruction following |
| MusicInstruct | audio URL, text | 5.5k clips, 60,493 Q&A | acoustic music instruction following |
| MusicBench | audio, text | 52,768 text-audio pairs | text to music |
| MARD | audio URL, text | 65,566 albums, 263,525 reviews | album music description |
| MUEdit | audio pairs, text | 10815 text, 60.22 hours | music editing with text prompt |
| WikiMusicText (WikiMT) | ABC, text | 1,010 | text to music, music captioning |
| IMAC | audio URL, image URL | 85k images, 3,812 songs | affective music-image correspondences |
| URMP | MIDI, Audio, Video | 44 pieces | audiovisual symphony separation |
| URSing | audio, video | 65 pieces, 4 hours | audiovisual singing voice separation |
| RAVDESS | audio, video | 7356 piece | speech & songs in different emotion and intense. |
| EmoMV | audio, video | 5986 pairs | affective music-Video correspondencs |
| SymMV | MIDI, audio, video | 1140 pairs, 76.5 hours | video background music generation |
| MUImage | audio, image | 9966 text, 27.72 hours | image to music |
| MUVideo | audio, image | 13203 text, 36.72 hours | video to music |
| AnyInstruct | text, audio, images | 108k instruction-following entries | instruction following w/ interleaved format. |
| V2M | audio, video | 190k pairs, 6403 hours | video to music |
| MMtrail | text, audio, video | 20m pairs, 27.1k hours | text to music, video to music |

comments elements from Amazon, enriched by metadata from MusicBrainz and descriptors from AcousticBrainz [PBK+15]. This dataset includes 65,566 albums, accompanied by 263,525 customer reviews, supporting multimodal frameworks for analysing human preferences and behaviours to music across different genres.

The MUEdit dataset [HLSS23a] is curated to facilitate the development and training of models capable of understanding and editing music based on specific prompts. The MUEdit dataset includes 10-second music pairs, each selected based on rhythmic characteristics such as tempo and pitch to ensure consistency and relevancy, totalling 55.69 hours. The creation process involves generating descriptions for these music pairs using the MU-LLaMA model, followed by the generation of editing instructions via the MPT-7B model. These instructions are designed to mimic a conversational interface where the model not only receives detailed descriptions but also generates editing commands based on these descriptions, simulating a realistic user-model interaction for music editing tasks.

WikiMusicText (WikiMT) [WYTS23a] contains 1010 lead sheets in ABC notation format from Wikifonia.org, each with metadata including title and artist extracted directly from the scores. Additionally, it derives genre and description from Internet information. The genre labels are assigned by analysing related Wikipedia content and categorising it into 8 GTZAN taxonomy classes: Jazz, Country, Folk, R&B, Pop, Rock, Dance, and Latin. Descriptions are synthesised from Wikipedia using the BART-large model. To emphasise musical content, extraneous text within the ABC notation has been removed.

*b) Visual-Audio:* Image-Music Affective Correspondence (IMAC) database [VDG19] is developed to advance research in crossmodal emotion analysis. This extensive dataset comprises the URL of over 85,000 images and 3,812 songs, totalling approximately 270 hours of audio. Each item in the database is categorised under one of three emotional labels: positive, neutral, or negative, facilitating detailed studies of emotional responses across different media.

The University of Rochester Music Performance (URMP) [LLD+18] collection presents a comprehensive set of multimodal data for 44 distinct musical compositions. This dataset encompasses a wide range of resources, including synchronised audio and visual elements along with detailed annotations. For each musical piece, the dataset provides a dedicated folder containing various components: MIDI scores, visually enhanced PDF sheet music and audio recordings. These audio files are available in both isolated instrument tracks and mixed versions, stored in high-fidelity WAV format with 48 kHz sampling rate and 24-bit depth. Moreover, the dataset includes video recordings of performances, captured in 1080P resolution at 29.97 frames per second, with careful alignment to the musical score. To facilitate in-depth analysis, the URMP dataset also incorporates frame-by-frame pitch information and note-level transcriptions, both presented in ASCII text format.

Creating the University of Rochester Music Performance (URMP) dataset posed significant synchronisation challenges due to the need for both high-quality recordings and efficient dataset creation. To tackle this, various methods were tested to balance expressiveness and synchronisation without extensive joint rehearsals. Initial attempts included using a metronomic beat for synchronisation, which proved too rigid, and a pre-recorded piano guide, which failed to provide adequate synchronisation cues. Further approaches involved using rehearsal recordings with a conductor to enhance expressiveness and synchronisation, leading to better results but at the expense of scalability. Subsequent strategies explored visual cues by having players follow video recordings of the first performer, which improved timing after long rests but still fell short of desired synchronisation levels. The most effective method involved a conductor directing the performance along with a pianist, recorded on video, which players followed during their individual recordings. This approach minimised the need for joint rehearsals while maintaining high synchronisation and expressiveness, proving to be efficient and scalable for creating

the dataset. Each method's efficacy was carefully evaluated, considering both quantitative synchronisation measures and players' subjective feedback, ensuring a balanced and practical approach to dataset creation.

The University of Rochester MultiModal Singing Performance Dataset (URSing) [LWD21] is designed to evaluate realistic singing performance models. It comprises 65 songs totalling four hours of audiovisual recordings from the University of Rochester students. The dataset creation involved singer recruitment with auditions for tuning accuracy, recordings using an AT2020 microphone and iPhone 11 within a semi-anechoic booth, and professional post-processing for audio quality. Annotations of mouth regions were added using the Dlib library, and the dataset includes solo and accompaniment audio tracks in WAV format, MP4 video recordings, and a benchmark set of 54 instrumental and 26 vocal-backed excerpts for detailed evaluation. URSing provides a rich resource for exploring audiovisual music information retrieval techniques.

The Ryerson Audio-Visual Database of Emotional Speech and Song (RAVDESS) dataset [LR18] provides a validated collection of 7356 multimodal emotional voice and song recordings performed by 24 professional actors. This gender-balanced dataset includes multiple emotional expressions of varying intensities, provided in a combination of facial video and sound or in separate formats, focusing on multimodal sentiment recognition and analysis in speech and music.

The Lyra dataset [PVG$^+$22] is a collection of traditional and folk music in Greek, containing 1,570 YouTube timestamp links of music videos and approximately 80 hours in total. Derived from a Greek documentary film series, the dataset provides a wealth of metadata on musical genres, instruments, and geographic origins. The dataset supports a wide range of computational musicological tasks, making it an invaluable resource for ethnomusicological research.

The EmoMV [TRH23] project is focused on assembling a trio of datasets aimed at facilitating the study of emotional correspondence between music and video modalities. The datasets, named EmoMV-A, EmoMV-B, and EmoMV-C, incorporate music video segments with EmoMV-A and EmoMV-B deriving content from existing datasets, while EmoMV-C features self-collected music videos from YouTube. Each pair within these datasets is annotated as either emotionally matched or mismatched.

Human experts conduct the annotation process for EmoMV-A, while a pre-trained deep neural network generates annotations for EmoMV-B and EmoMV-C. The reliability of these emotional labels is evaluated through a user study that examines their accuracy. In conjunction with the dataset development, researchers have created a benchmark deep learning model. This model is designed for binary classification of affective correspondence between music and video and has been optimised for affective music-video retrieval applications. This advancement provides a valuable resource for investigating how music and visual elements interact to express emotional content.

SymMV dataset [ZWW$^+$23b] is a larger video-music dataset that pairs music videos (MV) with piano MIDI and audio. It includes 1140 pop piano music pieces, along with their corresponding official MV, totalling 76.5 hours. Each MV pair is annotated for chord progression, tonality, and rhythm, enhancing its utility for research in MV emotional correspondence. The dataset functions as a valuable resource for developing algorithms that generate background music for videos. It enables the exploration of inherent relationships between visual and musical elements, leveraging the emotional and stylistic correspondence between video content and its soundtrack. The creation of this dataset involved a meticulous process: sourcing professional-grade piano renditions from select YouTube channels, aligning these with the original music videos, and utilising sophisticated transcription algorithms to convert the audio into MIDI format. To maintain high standards, professional musicians conducted a thorough review of the compiled dataset.

MUImage and MUVideo[HLSS23a] are text/image/video-to-music datasets that facilitate the development of models for generating customised music to complement accompanying visual content. These datasets were constructed by carefully generating detailed descriptions of images, videos and music. The MU-LLaMA model is used to add captions to music files, while the BLIP and VideoMAE models provide captions for images and videos, respectively. These descriptions are then used as input to the MPT-7B model, which produces music-generated cues. This structured approach supports the creation of integrated arbitrary music datasets

The AnyInstruct[ZDY$^+$24] dataset comprises 108k multimodal instruction-following entries, uniquely integrating text, speech, images, and music. Developed using GPT-4 for generating textual dialogues, and enhanced with images from DALL-E 3, music from MusicGen, and voices synthesised via the Azure Text-to-Speech API, this dataset features diverse modalities in an interleaved format. The dataset is organised into two parts, containing comprehensive multimodal dialogues along with approximately 205k images, 503k voice recordings, and 113k music tracks.

The V2M [TLY$^+$24] dataset comprises video and acoustic music pairs on a large scale. It includes three subsets: V2M-190K, which consists of approximately 190K video-music pairs totalling around 6,403 hours; V2M-20K, containing about 20K pairs with approximately 597 hours of content; and V2M-bench, which comprises 300 pairs totalling 9 hours in duration.

MMTrail [CWC$^+$24] is a comprehensive multimodality dataset featuring over 20 million trailer clips accompanied by visual captions, as well as 2 million high-quality clips annotated with multimodal music captions.

*6) Future Work:* Advances in music modelling require more expansive and diverse datasets. Current resources, in spite of their extensive size, predominantly feature Western popular music, limiting model versatility. To develop robust models that grasp and generate music authentically across various cultures, it's crucial to integrate a broad spectrum of global musical styles and genres into these datasets.

The current Multimodal Music Instruct dataset needs substantial enhancement to support foundation models effectively. It should expand not only in volume but also in the com-

plexity of its content, particularly in areas of sequential MIR downstream tasks like beat recognition and chord progression. Furthermore, introducing a chain-of-thought component is essential for enabling models to perform complex reasoning and explain the musical theories underpinning their responses, thus fostering deeper analytical capabilities in music understanding.

Additionally, datasets that simulate real interactions between users and music generation systems are necessary. These should include multi-round conversations with interleaved text-music data, including user queries about music recommendations or variations and the corresponding model responses. Such datasets would improve model responsiveness and adapt algorithms based on realistic user engagement, enhancing the overall user experience in music applications.

In addition to the data collection mentioned above, pre-training data cleaning strategies, mixing strategies of pre-training data from different sources, and curriculum learning with a specific priority of different types of datasets are still open issues to be explored.

Last but not least, there is a need to find reasonable ratios of balancing downstream tasks and pre-training data to develop effective tuning strategies for injecting new knowledge or expertise while avoiding catastrophic forgetting, e.g., allowing the GPT to learn more about music, adapting the model to new music genres, or make it more suitable to user preferences.

### B. Evaluation

Music understanding and generation are of increasing interest to both academic and industry researchers and have the potential to broadly impact the economics and culture of music. However, evaluation remains an open and difficult research problem, facing challenges like inconsistent downstream evaluation conditions, data leakage, model bias, etc. In the following subsection, we will review commonly used evaluation protocols and metrics for both music understanding and music generation.

*1) Evaluation for Music Understanding:*

*a) Overview:* In this section, we elucidate the predominant evaluation protocols followed in assessing foundation models for music understanding. As seen in Section III-A, these types of models can either serve as common backbones for solving traditional MIR tasks or can engage multiple data modalities, most often visual or textual modalities, to perform a variety of novel tasks that require multimodal understanding. Evaluation, therefore, varies depending on model design, tasks supported, and input-output modalities. Here, we distinguish between the evaluation of two main types of tasks: traditional MIR tasks and multimodal tasks, providing an overview in Table VII.

*b) Traditional MIR tasks:* We start by considering traditional **MIR tasks**. In this scenario, the goal is to measure music understanding by determining a model's ability to correctly recognise specific musical properties at different levels of granularity, from frame to track level, by adopting the canonical MIR evaluation approach for the corresponding task. For models operating in the audio domain, tasks typically included are a subset of key and mode detection [GDSB23a],

[LYZ+24], tempo detection [GDSB23a], [LYZ+24], genre classification [GDSB23a], [LYZ+24], [DESW23], instrument recognition [GDSB23a], [LYZ+24], [DESW23], music tagging [LYZ+24], emotion recognition [LYZ+24], pitch detection [LYZ+24], [DESW23], beat tracking [LYZ+24], and, in the less common performance domain, vocal technique detection [LYZ+24], singer identification [LYZ+24], and performance assessment [ZLD24], [ZJJH21], among others. MIR-based evaluation can be executed in two distinct ways: via probing [LYZ+24], [DESW23], [CDL21] or via text-based instruction [GDSB23a]. Probing is used on models that produce audio representations: the audio backbone, which is either kept frozen or fine-tuned, is used as a feature extractor, and a lightweight probing network, typically a shallow MLP, is trained on top to solve the downstream task of interest. Instead, in models equipped with a text interface, MIR-based evaluation can be realised by prompting the model with an adequate text input (e.g. *"What is the key of this song?"*), and casting the output text (e.g *"This song is in the key of F minor."*) to the corresponding label in the dataset.

Recently, there have been several efforts to unify this kind of evaluation methods across the full spectrum of music understanding tasks, by creating **benchmarks** that either include the music domain as a subset of interest [TSK+22] or fully specialise in it [YML+24], [PAJB23]. The most extensive of these is MARBLE [YML+24], which proposes a unified protocol for evaluating music audio representations across tasks at different levels of hierarchy, including high-level description (key detection, music tagging, genre classification, and emotion recognition), score-level (pitch tracking, beat tracking, melody extraction, chord estimation, lyrics transcription), performance-level (ornament and technique detection), and acoustic-level tasks (singer identification, instrument classification, source separation).

*c) Multimodal tasks:* In the context of **multimodal** foundation models and beyond traditional MIR tasks, there is a similarly wide variety of tasks that are useful in evaluating music understanding. Among these, one that stands out as a common approach to evaluate fundamental capabilities in this domain is **cross-modal retrieval**, defined as the task of retrieving data samples of one modality (e.g. a music clip), based on a query sample of another modality (e.g. a textual description or a video). Models that are amenable to this type evaluation are most typically those trained to learn joint embedding spaces between two or more music-related data modalities, such as music-text contrastive models [HJL+22], [MBQF22a], [WYTS23a], or those trained to learn correspondences between music and images or videos [MSSR23], [TRH23], [SAFN24], [SVRS22], [VDG19], [YZT+24]. In the case of music-text retrieval, this evaluation approach is not only useful for directly measuring performance on the task itself, but it also offers a framework for evaluating models on music classification tasks by casting them as text-based retrieval. In this case, cross-modal retrieval results in a further variation of the MIR-based evaluation approach presented earlier [DWCN23].

A further type of multimodal evaluation comprises tasks that are centred around language generation. As such, this is

Table VII: Evaluation of music understanding foundation models. * indicates human evaluation.

| Model | Traditional MIR Tasks | | | | | | | | | Multimodal Tasks | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Beat | Chord | Pitch | Key | Tempo | Instrument | Genre | Emotion | Tags | Retrieval | Captioning | MQA |
| M2-Duo [NTO+24] | | | ✓ | | | ✓ | ✓ | | | | | |
| Music2Vec [LYZ+22] | | | | ✓ | | | ✓ | ✓ | ✓ | | | |
| MusicHuBERT [MYL+23] | | | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ | | | |
| MusicFM [WHL24a] | ✓ | ✓ | | ✓ | | | | | ✓ | | | |
| MERT [LYZ+24] | | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | | |
| COLA [SGZ21] | | | | | | ✓ | | | | | | |
| CLAR [ATM21] | | | ✓ | | | ✓ | | | | | | |
| MULE [MKO+22] | | | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ | | | |
| JukeMIR [CDL21] | | | | ✓ | | | ✓ | ✓ | ✓ | | | |
| MuLaP | | | | | | | | | | | | |
| MusCALL [MBQF22a] | | | | | | | ✓ | | ✓ | ✓ | | |
| MuLan [HJL+22] | | | | | | | | | ✓ | ✓ | | |
| CLAMP [WYTS23a] | | | | | | | ✓ | ✓ | | ✓ | | |
| M$^2$UGen [HLSS23b] | | | | | | | | | | | | ✓ |
| SALMONN [TYS+24a] | | | | | | | | | | | ✓ | |
| MuLLaMa [LHSS23a] | | | | | | | | | | | ✓ | ✓ |
| MusiLingo [DML+24] | | | | | | | | | | | ✓ | ✓ |
| Pengi [DESW23] | | | ✓ | | | ✓ | ✓ | | | | | |
| LLark [GDSB23a] | | | | ✓ | ✓ | ✓ | ✓ | | | | ✓* | ✓* |

only suitable to music understanding models that take audio or (audio, text) pairs as inputs and generate text outputs. Under this setting, music understanding is measured by assessing language outputs that typically encompass several musical concepts. In practice, this is articulated in two subtasks: **music captioning** [TYS+24a], [LHSS23a], [DML+24], [GDSB23a] and **music question answering** (MQA), which can be either open-ended [TYS+24a], [HLSS23b], [LHSS23a], [DML+24] (sometimes also referred to as *music reasoning* [GDSB23a]), or multiple-choice [WMB+24]. In music captioning and open-ended MQA, the model is instructed to produce a language output that either describes the audio input (e.g. *"Describe the contents of the provided audio in detail."*) or answers a question about it (e.g. *"What are some possible uses for this music in a film or TV show?"*). Performance is then measured in one of three ways: via automatic metrics, via human evaluation or via another LLM, using a strategy called LLM-as-a-judge [ZCS+23]. In the first case, automatic evaluation typically consists of computing text generation metrics such as BLEU [PRWZ02], METEOR [BL05], ROUGE [Lin04], and BERT-score [ZKW+20], which measure the linguistic or semantic similarity between the generated and reference text. We note that MQA evaluation bears similarities with the flavour of MIR-based evaluation that makes use of (question, answer) pairs. A key difference, however, is that the former is characterised by open-ended questions, while the latter focuses on a single musical aspect and assumes a one-to-one mapping between model output and ground truth.

While evaluation on traditional MIR tasks leverages standard metrics and benchmarks from the MIR literature and beyond, language-based evaluation follows less established protocols and is supported by only a handful of datasets. For the music captioning task, most works [TYS+24a], [GDSB23a], [DML+24] make use of the MusicCaps dataset [ADB+23], while others [LHSS23a] employ ad-hoc evaluation datasets created with the help of LLMs. For the MQA task, no standard dataset exists for the open-ended

variant, and all evaluations are carried out on ad-hoc datasets [LHSS23a], [HLSS23b], [DML+24], [GDSB23a]. The MuChoMusic benchmark [WMB+24] supports instead music understanding evaluation via multiple-choice MQA. In this setup, models are provided with an audio clip accompanied by a question and a set of answer options (A, B, C, or D) to choose from. The closed format of this task allows for more standardised evaluation, overcoming some of the limitations of match-based metrics.

*d) Measuring music-domain knowledge in LLMs and Visual-Language Models:* Finally, we look at assessing music understanding abilities in non-music foundation models, through the lens of recent benchmarks that allow to evaluate music-domain knowledge encoded in LLMs [YLW+24a], [LYT+24] and visual-language models [YNZ+23].

For LLM evaluation, we identify two key benchmarks, MusicTheoryBench [YLW+24a] and ZIQI-Eval [LYT+24], both employing ABC notation (see Section II-B). MusicTheoryBench [YLW+24a] focuses on evaluating music knowledge and reasoning via a set of 372 multiple-choice questions designed to align with college-level music education standards. The questions span a variety of topics including notes, rhythm, chords, orchestration, and music-related cultural and historical knowledge and are aimed at assessing both the analytical and inferential skills of LLMs. ZIQI-Eval [LYT+24] is instead designed to test LLMs' musical abilities in both comprehension (focusing on harmony, melody, and rhythm) and music generation or continuation. Extensive testing of 17 LLMs through this benchmark reveals that current models generally possess insufficient knowledge in the music domain and exhibit poor compositional understanding and generation abilities, underlining the need for targeted improvements in musical training for LLMs.

Considering instead visual-language models, the Massive Multi-discipline Multimodal Understanding and Reasoning Benchmark (MMMU) [YNZ+23], among its broad evaluation of multimodal models across various disciplines, includes a

dedicated subset focused on music sheet understanding. This features 369 music sheet images accompanied by targeted questions which probe knowledge of musical theory, including harmonic intervals, instrument tuning, and chords within specific bars.

*e) Challenges and open questions:* Despite the flourishing of research in this field, evaluating music understanding models in a consistent and reliable fashion remains a challenge due to limited benchmarks dedicated to the music domain. As a consequence, existing works follow inconsistent evaluation practices which make use of different datasets, metrics, experimental settings, and even task formulations. This becomes particularly prominent in language-based evaluation, where, as highlighted earlier in this section, there is little in the way of standardisation, and automatic metrics become unsuitable due to the open-ended nature of the tasks [GDSB23a]. Text generation metrics are also known to deviate from human judgement when there isn't a sufficient number of reference sentences or, similarly, when there are several potentially correct outputs, as is often the case in music understanding. This has prompted some to conduct human studies to obtain more reliable evaluations on open-ended tasks [GDSB23a].

Finally, comprehensive evaluation should go beyond assessing task-specific performance and include other key factors like robustness, safety, efficiency, and alignment with human preferences. While this is becoming commonplace in other domains, it remains unexplored in the context of music understanding.

*2) Objective Evaluation for Music Generation:*

*a) Overview:* Similar to the field of music generation itself, research on objective evaluation for music generation models is in a nascent stage–The field has yet to agree on a set of standard metrics and/or benchmarks, and new metrics have been constantly proposed in an ad-hoc manner with new methods. Objective evaluation of machine-generated music is challenging for several main reasons: (i) The aesthetics of music is multifaceted, e.g., harmony, melody, instrumentation, and how they are assembled together, which are both hard to quantify individually and highly interrelated with one another. (ii) Music generation is, in fact, a family of tasks that can take a variety of user inputs, ranging from a short text description [ADB+23], [CKG+24] to fine-grained musical concepts like melody [CKG+24], chords [vRBKH23], beats [WDWB24], and even music composition with rendering (performance) objective [PCCKW23], [ZCCC+24], hence necessitating task-specific metrics.

We provide an overview of the landscape of objective evaluation metrics for music generation according to three axes: (i) **functional type**, (ii) **purpose**, and (iii) **input domain**.

A **functional type** (or **type** in short) defines the high-level workflow to map a set of generated (and/or reference) samples to a metric value. Let one sample of *reference* (i.e., human-composed) music be denoted by $\boldsymbol{x}$, and one sample of *generated* music by $\tilde{\boldsymbol{x}}$. Also, let $X = \{\boldsymbol{x}^{(i)}\}_{i=1}^{N}$ denote the set of $N$ reference samples, and $\tilde{X} = \{\tilde{\boldsymbol{x}}^{(i)}\}_{i=1}^{\tilde{N}}$ the set of $\tilde{N}$ generated samples. We further define $f(\boldsymbol{x})$ or $f(\tilde{\boldsymbol{x}})$ to be any function that extracts an *attribute* or some *features* from music, either global ones (e.g., a text description for the music) or

local ones (e.g., a chord sequence), which may sometimes also be independently given by users as conditions, and $\mathrm{score}(\cdot, \cdot)$, or $\mathrm{score}(\cdot)$ to be a deterministic function that computes the reported metric value. Commonly adopted metrics may be categorised as follows:

- **Pair-wise aggregate** (*P-A*) – $\mathrm{score}(\{f(\boldsymbol{x}^{(i)})\}_{i\in[N]}, \{f(\tilde{\boldsymbol{x}}^{(i)})\}_{i\in[\tilde{N}]})$. *P-A* metrics seek to compare the distribution of outputs induced by music generation models to a sample of the distribution of human-created music. To compute *P-A* metrics, two sets of samples (generated and reference) with or without direct correspondence are needed. Attributes/features are extracted from each sample, and the $\mathrm{score}(\cdot, \cdot)$ function would then calculate some aggregate statistics of the attributes/features in each set, between which some notion of similarity or distance is computed. Frechét audio distance (FAD) [KZRS19] is a representative example of a *P-A* metric.

- **Pair-wise individual** (*P-I*) – $\mathrm{mean}_{i\in[\tilde{N}]}(\mathrm{score}(f(\boldsymbol{x}^{(i)}), f(\tilde{\boldsymbol{x}}^{(i)}))$, or $\mathrm{mean}_{i\in[\tilde{N}]}(\mathrm{score}(f_1(\boldsymbol{x}^{(i)}), f_2(\tilde{\boldsymbol{x}}^{(i)}))$, where $f_1$, $f_2$ are different attribute, or feature, extraction functions. *P-I* metrics are used when there is a clear *target* (e.g., a piece of human-composed music, some musical features, or semantic attributes) that each *individual* machine-generated piece is expected to resemble or reflect. *P-I* metrics require sample-wise direct correspondence between reference and generation, and hence assert $N = \tilde{N}$. CLAP score [WCZ+23b] is a widely adopted *P-I* metric.

- **Single individual** (*S-I*) – $\mathrm{mean}_{i\in[\tilde{N}]}(f(\tilde{\boldsymbol{x}}^{(i)}))$. *S-I* metrics typically serve as proxies for musical aspects that are hard to quantify, e.g., how "good" the harmony, rhythm, or repetitive structure is. *S-I* metrics are computed without a reference. Therefore, unlike *P-A* or *P-I* metrics, where we may specify higher/lower ($\uparrow$/$\downarrow$) is better, *S-I* metrics are often specified as the closer to reference samples, i.e., $\mathrm{mean}_{i\in[N]}(f(\boldsymbol{x}^{(i)}))$, the better. For example, pitch-class histogram entropy [WY20] is an *S-I* metric. We note that, while it is more desirable to compare the distribution of $f(\tilde{\boldsymbol{x}}^{(i)})$ induced by machine-generated music to that of human-composed music, i.e., $f(\boldsymbol{x}^{(i)})$, effectively turning *S-I* metrics into *Single aggregate (S-A)* metrics following our taxonomy, researchers have gravitated to comparing simple $\mathrm{means}$ of features $f(\tilde{\boldsymbol{x}}^{(i)})$ and $f(\boldsymbol{x}^{(i)})$.

The **purpose** of a metric is associated with either of the two following major focuses:

- **Quality**-focused: Such metrics try to characterise how close generated samples are to human-composed music, and can focus on a variety of aspects, like audio quality, musical diversity, rhythm, melody, harmony, repetitive structures, and so on.

- **Control**-focused: Such metrics examine how well a generative model adheres to the conditions given by users. Frequently used conditions are text descriptions, musical genres, emotions, instrumentation, tempos, chord

sequences, melodies, etc.

The **input domain** of an evaluation metric largely depends on the domain the evaluated model/system targets(cf. II-B). Audio-domain metrics accept audible waveforms as inputs, while symbolic-domain metrics apply to note/event sequences. Audio-domain metrics may be used on symbolic-domain samples with the help of audio synthesisers, but not vice-versa in general as music transcription [BDDE19] remains a challenging task. We note that, compared to audio-domain models, evaluation of symbolic-domain models is relatively less standardised – While some works follow the mechanisms proposed in evaluation-focused papers [WY20], [YL20], the majority define their own metrics for evaluation.

*b) Audio-domain, quality-focused metrics:* **Frechét audio distance** (**FAD**; type: *P-A*, ↓) [KZRS19] is the most widely adopted metric to assess the quality of generated audios, which has been shown to correlate well with human perception. FAD employs an audio encoder pre-trained on, e.g., multilabel audio classification [GEF+17], as $f(\boldsymbol{x})$ to map an input audio to a feature vector. Then, the $\text{score}(\cdot, \cdot)$ function first estimates a Gaussian covariance matrix from all feature vectors obtained from generated audios, and likewise for reference audios. Finally, the Frechét distance between the two estimated Gaussian distributions [DL82] is computed. As the procedure suggests, FAD conflates audio quality, musical quality, cross-sample diversity, and potentially many more aspects together, and measures a notion of *overall realisticness* of audios.

As for metrics that focus on one specific aspect, the **Structureness indicator** (type: *S-I*) [WY20] is a representative example, which examines the presence of *repetitive structures* across all time granularities. The Structureness indicator leverages Fitness Scape-plots [MJ12] as $f(\boldsymbol{x})$, and then applies a $\max$ operation on the scape-plot as the $\text{score}(\cdot)$ function, which is meant to describe how much the most-repeated excerpt within a specified time granularity (e.g., 10∼20 seconds) is repeated throughout the entire audio.

*c) Audio-domain, control-focused metrics:* With the recent rise of text-to-audio music generation models [LCY+23d], [ADB+23], [CKG+24], the most popular control-focused metrics currently are concerned with *adherence to text* (or genre/emotion/instrument tag-based) inputs:

- **CLAP score** (**CLAP**, type: *P-I*, ↑) [WCZ+23b] depends on a text encoder and an audio encoder as feature extractors $f_1, f_2$. The two encoders should have been jointly pre-trained via contrastive learning [vdOLV18] to form an aligned latent feature space. Then, CLAP score is simply the cosine similarity between the encoded features of the input text and those of the generated audio. MuLan score [HJL+22] is an almost identical metric, only with a different bi-encoder backbone.
- **KL divergence** (**KL**, type: *P-I*, ↓), as opposed to CLAP score, measures text input adherence indirectly. It requires a pre-trained audio classifier [ZTdCQT21], [KSEZW21] as $f(\boldsymbol{x})$. The KL divergence is computed between the output class distribution of the reference audio (for the text input) and that of the generated audio.

To evaluate adherence to more fine-grained and musically-centred types of controls, the metrics are often proposed in individual works and hence not yet standardised. Commonly considered musical controls and their corresponding metrics are as follows:

- **Melody**: **Chroma cosine similarity** (type: *P-I*, ↑) [CKG+24] leverages the chromagram (or pitch-class profile) [Fuj99] as $f(\boldsymbol{x})$ to extract the relative strength of 12 semitones from reference and generated audios, from which the frame-wise (or chunk-wise) cosine similarity is computed. If the chromagram is frame-wise binarised to be one-hot, the metric is equivalent to frame-wise **chroma accuracy** [WDWB24].
- **Chords**: **Chord recognition accuracy** (type: *P-I*, ↑) [HDS+23b], [MGG+23] use a trained chord recognition model (e.g., [CCK+22]) to predict the chord sequence from reference and generated audios. Then, an alignment (with various levels of tolerance [MGG+23]) between the two sequences is performed to compute the score.
- **Tempo**: **Tempo bin accuracy** (type: *P-I*, ↑) [MGG+23] relies on a beat detection model [BKS+16], [HCD21], whose predictions can be converted to beats per minute (BPM). The metric then checks whether the BPMs of the reference and generated audios fall in the same bin, e.g., Adagio (60∼70 BPM).
- **Rhythm**: **Beat F1 score** (type: *P-I*, ↑) [WDWB24] also employs a beat detector and follows the evaluation for beat detection models [DDP09]. It aligns the timestamps of beats predicted from reference and generated audios with a 70ms tolerance. **Beat count** (type: *S-I*) [MGG+23] is a simpler version that only counts the number of detected beats.
- **Intensity**: **Dynamics correlation** (type: *P-I* or *P-A*, ↑) [WDWB24] computes the smoothed frame-wise loudness in decibels, and then calculates the Pearson's correlation between the frame-wise values from reference and generation.

*d) Symbolic-domain, quality-focused metrics:* This family of metrics often fall under the *S-I* type, and use the corresponding $f(\boldsymbol{x})$'s to aggregate/summarise some feature (concerning aspects like pitch, rhythm, and tonality) from the token sequence representing $\boldsymbol{x}$, and then compute the score on the feature.

- **Pitch-class histogram entropy** (type: *S-I*) [WY20] collect the counts of 12 pitch classes (i.e., C, C#, . . . , B) over a certain period (e.g., 4 bars), and computes the entropy on the resulting histogram as a measure of the diversity of pitch classes used.
- **Groove consistency** (type: *S-I*) [WY20] divides each bar into segments (e.g., in 32th note increments), and construct a binary *groove vector* for each bar by checking whether each segment has note onsets. Then, some similarity (e.g., Hamming similarity) between groove vectors of neighbouring (or all) bar pairs is computed.
- **Scale consistency** (type: *S-I*) [Mog16] computes the fraction of pitch classes that belong to the 7 pitch classes of each of the 24 major/minor scales, and takes the maximum over the 24 scales as the scale consistency.
- **Similarity error** (type: *P-A*, ↓): was proposed by

[YLW+22] and encompasses both pitch and rhythm to examine whether generations are structurally similar to human-composed music. It first constructs a *note set* for each bar, whose elements describe the (pitch, duration, onset time relative to downbeat) of each note, and then computes the mean intersection-over-union (IoU) similarity of all bar pairs that are $t$ bars apart (where $t$ is preset). Finally, the difference between the mean IoUs between human-composed and generated pieces is taken.

*e) Symbolic-domain, control-focused metrics:* Unlike audio-domain models, where text descriptions are the prevailing control form, control inputs for symbolic models are much less standardised, and hence so are the evaluation metrics. Therefore, we group the metrics roughly according to controls that targeted aspect, and refer readers to relevant papers for further details.

- **Pitch/melody:** melody matchness (*P-I*, ↑) [HLYY21b], pitch distribution L2 distance (*P-I*, ↓) [ZRZY23].
- **Chords:** chord matchness [HLYY21b], chord F1 (between controls and generated music; *P-I*, ↑) [vRBKH23], and chord progression L2 distance (*P-I*, ↓) [ZRZY23].
- **Time signature:** time-signature accuracy (*P-I*, ↑) [vRBKH23].
- **Instrumentation:** instrument F1 (*P-I*, ↑) [vRBKH23].
- **Intensity:** Pearson's correlation of rhythmic intensity or polyphony (*P-A*, ↑) [WY23c], and note density L2 distance (*P-I*, ↓) [ZRZY23].

*3) Subjective Evaluation for Music Generation:* Due to the highly artistic and subjective nature of music, subjective studies, i.e., listening tests with human participants, still remain the gold standard of evaluation music generation models [YL20], [JYL23]. Regardless of whether the model operates in the audio or symbolic domain, the samples presented to participants are typically (synthesised) audio waveforms. Below we explain two common protocols for subjective studies.

*a) Pair-wise Comparisons:* In each set of the study, a participant is provided with two samples $\tilde{x}_1, \tilde{x}_2$ (without being able to infer the underlying generative models), and asked to choose from the three options: (i) $\tilde{x}_1$ is preferred over $\tilde{x}_2$, (ii) $\tilde{x}_2$ is preferred over $\tilde{x}_1$, or (iii) no preference among $\tilde{x}_1, \tilde{x}_2$. Associating the responses with the underlying models known only to the study setters, the *win-loss-tie* counts for every pair of the involved models can be obtained, on which the Wilcoxon signed-rank test [Wil92] is performed to check whether one model is superior to another with statistical significance. Works that adopt this protocol include [HVU+18b], [vRBKH23], [ADB+23], [THDL23].

*b) Absolute-scale Ratings:* Under this protocol, participants are asked to independently rate each sample $\tilde{x}$, possibly with an extra given condition $c$, e.g., a melody, or a short prompt. The ratings are typically in 5-, 10-, or 100-point scales. The most frequently evaluated aspects are:

- **Overall quality (OVL)**
- **Relevance to input text (REL)**,

which are seen in [LCY+23d], [GMMP23b], [CKG+24], [LTY+23b], [HDS+23b], [MGG+23] and focus on quality and control respectively.

On the other hand, compared to pair-wise comparisons, more fine-grained aspects are often considered. Although various terms and questions have been used, the fine-grained aspects may roughly be grouped into several categories:

- **Audio quality**: Whether the audio is clear, without noises [SJS23a], [MGG+23].
- **Melody quality**: Whether the melody is prominent, pleasant-sounding, and leaves a strong impression [ZZQ+22], [SJS23a].
- **Rhythm quality**: Whether the rhythm is clear, stable, and reasonable [ZZQ+22], [MGG+23].
- **Harmony**: Whether the melody, accompaniment, and chord sound pleasant together [SJS23a], [MGG+23], [HDS+23b], [ZZQ+22].
- **Orchestration**: For multi-instrument models, whether the instruments are used and arranged properly [DCD+23], [LDC+22].
- **Coherence**: Whether the transitions between musical phrases are natural and well-connected [HLYY21b], [ZZQ+22], [LDC+22], [WY23a], [DCD+23].
- **Structure**: Whether the music exhibits clear repetitive structures and reasonable development [WY20], [HLYY21b], [ZZQ+22], [LDC+22], [YLW+22], [WY23a].
- **Richness**: Whether the music is rich, creative, and interesting [HLYY21b], [LDC+22], [DCD+23], [WY23a], [MGG+23].

To examine statistical significance with absolute ratings, a Student's $t$-test [Gos08] is usually performed.

## VI. ETHICS & SOCIAL IMPACT

The development and deployment of FMs for music applications raise many ethical and cultural issues. This section delves into four main aspects identified through related work aimed at uncovering potential issues within the broader space of Music Information Retrieval (MIR). While little of this work specifically addresses FMs, most of these aspects are critical to understanding the ethical landscape surrounding FMs in music applications. The first section focuses on ethical and social impact issues. The first subsection highlights how the lack of diversity within the space where FMs are developed can negatively impact various "ologies": ontology, epistemology, and axiology. The second subsection addresses specific issues related to fairness, bias, transparency, and explainability. The third subsection examines the potential adverse consequences of FMs on humans. Finally, the fourth subsection provides suggestions on taking responsibility for the work done in this field. Then, the second section will discuss music copyright issues.

### A. Ethical issues

*1) The Sectarian -ologies:*

Might a more diverse population of MIR practitioners favour awareness of and sensitivity to a wider spectrum of the world's musics? [Bor20]

A major concern is the threat to diversity embedded in the sociotechnical systems that form the backdrop of FMs. ML

models often embed specific and narrow values and cultural backgrounds, which fail to capture and support diverse musical provenance and heritages. There is a fundamental need to account for cultural pluralism in MIR and to include contributions from non-Western traditions [Cla21], [HSH21], [Mor21], [Ber24]. While FMs should be able to handle all kinds of music, they indeed often fail to capture or support cultural heritage and diversity. The consequences are tangible: biases toward Western popular music, fostering of commonalities rather than distinctiveness, and threats to regional creative work. A significant concern is the prevalence of Western popular music in training FMs, which exacerbates the historical trend of Westernisation and displacement of local content. This issue mainly arises from the narrow socio-cultural-economic background of FM creators: MIR practitioners and researchers come predominantly from WEIRD (Western, Educated, Industrialised, Rich, Democratic) demographics [HSC18b], which leads to cultural biases reproduced in their work. Against this backdrop of cultural homogeneity, [HHSK23] suggests that an ethical turn in MIR should involve rethinking the four ologies—ontology, epistemology, methodology, and axiology. We borrow this framing from them to structure our discussions around what we consider the sectarian -ologies: the narrow ontologies, epistemologies,[28] and axiologies that currently affect the FM space.

**Ontology:** When reporting on the issues around recommender systems, [GJFA21] highlights that research practices are shaped by those who create them. Thus, the ontological question "what music is" whose engagement is fundamental for the development of FMs, should account for various music forms and ideas (ibid). The theory of "what music is" that MIR embodies, which [Bor10] called analytical ontology, includes ontological assumptions about what counts in music. The prioritised ontologies are not only those musical aspects that are important for Western tonal music (harmony, melody, rhythm), but also "those sounds that, through recording, have been disembedded from originating bodies, socialities and locales" (ibid). [Bor20] invites MIR researchers to question, "Whose music and which music, among the vast ocean of sounds in the world, gets to be the focus of MIR's influential scientific practices?". While, in recent years, one may notice a rise in the use and focus on non-Western music, especially from India, there is a need to diversify our understanding of music further. Music is not just a sequence of notes; it carries significant cultural value, including the use of specific instruments and scales. The lyrics may convey the feelings or opinions of a specific group of people.

Notably, ontological endeavours are not limited to unpacking "what music is," but also need to account for "what *good* music is." [NML23] discusses how the Eurocentric Westernisation of aesthetics influences training data and is misaligned with much of the world. This aspect is particularly critical in FMs, as exemplified by the reliance on LAION-Aesthetics for image FMs like Midjourney and Stable Diffusion, which were rated visually appealing by WEIRD

individuals and AI-art developers [Mac24][29]. There is a need to dissuade the idea of "universal" foundation models in favour of culture-specific and culture-preserving models. It remains to be seen whether supervised fine-tuning techniques can build rich culture-specific models or merely project non-Western musical languages into a space of Western musical concepts. This comment is particularly relevant when considered against the growing body of literature that criticised FMs for the lack of universalism. For instance, as discussed by [BGMMS21], GPT-2 was trained on sources that mostly represent males in their twenties. Thus, music FM developers should avoid perpetuating hegemonic ontologies or at least acknowledge that their system is not representative of all music. The risk is that, in the case of music generation, representation will be further skewed as music libraries are inundated with synthetic music influenced by inequitably designed FMs.

**Epistemology** The second "ology" affected by the lack of diversity among FM researchers is epistemology, which deals with what constitutes knowledge and how to reach it. Cultural homogeneity in the field results in similar methodologies and methods being employed across the field, which result in a limited palette of epistemological tools that are considered "pertinent." [Bor20] asks "How could MIR equip itself with tools suited to analysing and modelling a greater diversity of musics?" [Mor21] proposes that MIR research "is not merely a matter of mathematical models and computational optimisation," and that the field should embrace an "epistemological turn" to broaden the skill sets of MIR researchers, in order to allow them to take ownership around discussions on effects of AI on music making and consumption. Similarly, [CKM+19] proposes that MIR researchers should be educated on an array of non-technical matters (i.e. ethical, cultural, and financial issues) related to use of music data. [HSC18b] discusses the adaptation of MIR developments, emphasising that they must be possible without requiring extensive engineering expertise. [Bor20] highlights the numerous issues with universalising non-universal techniques, noting that "the techniques and parameters employed in MIR tend to derive from, and reflect, commercially dominant areas of global popular music. Yet those techniques and parameters come to be applied in powerful technologies as though they were universal, with inevitable 'de-pluralizing' effects."

Notably, calls for more plurality and diversity in AI and MIR research [Ser11], [HSC18b], [Bor20], [Mor21], [HHSK23] align with the core tenets of feminist epistemology [SV24], particularly with the concepts of *situated knowledge* [Har88] and *strong objectivity* [Har95], [Har15]. These concepts are intended to surface how systemically marginalised groups may possess greater epistemic insight due to their different experiences and expertises compared to dominant groups. Although incorporating all perspectives is essential for a comprehensive understanding, placing a focus on marginalised standpoints helps to ensure that the influence of structures of oppression on scientific research is accounted for. Such attention towards diversity should not only be reflected and fostered in the composition of teams

---

[28]Different from Huang, this essay combines epistemology and methodology, viewing methodology as a way to reach knowledge.

[29]https://knowingmachines.org/models-all-the-way

of researchers, both in terms of demographics and fields of expertise, but also fundamentally embedded in study design [TEE+19], [SKP+11].

**Axiology:** Axiology, the study of value, is proposed by [HHSK23] as a critical consideration in MIR: the demographic and cultural narrowness of MIR researchers can impact what is deemed good or valuable research. This section focuses on the values and agendas embedded in FM research and who benefits from it [MSW23], [MB23]. [GJFA21] indicates that ontological assumptions extend beyond music itself to include specific theories of human subjectivity embedded in MIR. For instance, recommender systems frame users as individuals overwhelmed by choice, and seek to maximise the time they spend on the platform rather than, for instance, on delivering diverse music. The human subjectivity that underlines this representation opposes the view that users are sovereign individuals who are fully aware of the specific music they want to listen to. The ideology of recommender systems is thus based on serving something palatable to users while denying them the ability to make those choices. [GJFA21] notes naturally competing interests behind AI technologies, as corporate private interests do not necessarily match public interests. This comment applies to FM-related work. [MSW23] found that creators of datasets used in MIR tend not to prioritise musicians' rights and demands, often ignoring the extent to which ML models are fair for musicians. Similarly, [GJFA21] invites us to question "what masters/mistresses does MIR serve?"

Analysing the values embedded in FMs is particularly important as the field is now inundated with venture capital [Mor21]. The music industry has a large interest in AI and FMs, but understanding the industry's effect on culture, society, musicians, listeners, and researchers remains challenging. New technologies created opportunities for new industry classes, such as streaming services and generative AI. However, [GJFA21] suggests that when MIR research is tied to the drive for economic growth, issues of sustainability [Dev19] and human values (e.g., transparency, accountability, privacy, and security) are at stake. This research should seek to identify "other goals and values to guide future science and engineering" and consider how it would look under a different set of assumptions and incentives linked to human flourishing. [GJFA21] advocates for a logic of diversity (rather than similarity) and collectivism. An important aspect for FM developers to consider is whether non-Western music should be represented and made globally available through digital data in the first place. [Bor20] suggests that non-Some forms of Western music have not been incorporated into global digital music archives. These types of music have diverse ontologies that are often deeply embedded in social, religious, or cosmological traditions, and they differ greatly from the universal music ontologies assumed by MIR and related disciplines. [MBB+23] indicates that ML applications should resist misappropriating original intents, and [HHSK23] explains how Western-centric theories might be insensitive to local practices, traditions, knowledge, and values that form distinct music ecosystems. [HSH21] adds several relevant topics to this discussion. By turning to Chinese and Indigenous

philosophies, the authors urge MIR developers to broaden their scope and consider both human and non-human rights when discussing AI ethics[30]. For instance, AI systems that produce billions of songs "just because they can" clash with Mohist's condemnations of "wasteful productions and performances of music." This more-than-human ecomusicological perspective is particularly relevant to FMs, given the exponential environmental costs of energy-consuming neural networks [BGMMS21], [Cra24]. Furthermore, following [HSH21] invites inquiry into the effect of excessive productions and appoints music AI developers the responsibility to consider how their products impact our soundscape's health amidst environmental crises.

*2) Fairness, Transparency, and Bias:* Since there can be no accountability without transparency, and without accountability fairness in unachievable, here we approach an understanding of how the themes of fairness and transparency are operationalised and proceduralised in AI and MIR. We discuss how achieving fairness in AI requires rejecting the myth of value-free science and value-neutrality [ZAPX24], [Mil21], which can perpetuate inequalities by ignoring the underlying social structures that AI systems can replicate. By acknowledging the inherent values and impacts of technological artefacts, designers and developers must accept that their creations are not morally neutral and instead should consider integrating explicit ethical and moral considerations early in the design process to tackle concerns upfront rather than shifting blame to users later [Mil21].

Critics have long argued that technological development must consider cultural-historical dimensions to avoid unintentional consequences and ensure ethical development. As a matter of fact, AI has a long history of contentious interactions with its critics. As early as three decades ago, [Agr98] suggested that these conflicts might be due to the AI field's limited engagement with cultural-historical explanations that scholars from social studies instead take for granted:

> Without the idea that ideologies and social structures can be reproduced through a myriad of unconscious mechanisms such as linguistic forms and bodily habits, all critical analysis may seem like accusations of conscious malfeasance. Even sociological descriptions that seem perfectly neutral to their authors [i.e., the critics] can seem like personal insults to their subjects [i.e., the AI practitioners] if they presuppose forms of social order that exist below the level of conscious strategy and choice.

This misunderstanding and lack of involvement with social studies mostly stems from the longstanding myth that science is value-free. [Har15, p. 7-11] argues that this myth, which became the hegemonic understanding of "objectivity" since the end of World War II, has contributed to establishing a hierarchy of knowledge, with STEM at the top and humanities at the bottom, and brought about a sense of autonomy of science from politics[31]. This divide, at present, is

---

[30]Notably, this move towards the non-human or more-than-human has become increasingly popular even within Western schools of thought like posthumanism [Bra16] and even in scientific discourse [Bar07]

[31]Although certainly not from corporate interests [Har15, p. 7-11].

causing a lot of confusion within the AI community about what terms like "fairness", "diversity", "transparency", and "bias" really mean [ZAPX24], [HCVKL23], [BRGL+23], [MKKW19]. In addition, several scholars outside [Har15] and within [HSC18b], [HHSK23], [Mor21] MIR argue that value-free/neutral assumptions hinder a truly ethical development of science and technology, as these assumptions can at times be unintentionally [HSC18b] instrumental in preserving inequalities in the distribution of (algorithmic) power and resources. In other instances, it can even be deliberately directed towards "attempts to isolate science from sensitive questions" [Pro91, p. 267].

**Fairness:** Understanding the concepts of "fairness" and "diversity" [PCGV23] requires first restricting their scope within a precise subset of tasks that are relevant in MIR research: fairness in automated decision-making (ADM) and diversity in dataset creation and curation. According to [Alg19] ADM refers to the delegation of decision execution to systems that use algorithmic models to carry out actions based on data. Additionally, the term ADM encompasses not only complex machine learning applications but also simpler yet impactful rule-based systems like risk assessment scores. To date, the most predominant MIR application pertaining to ADM are music recommendation systems [HCVKL23], [DB22].

Researchers should reject value-free assumptions whilst implementing operationalisations (e.g., measurements) of "fairness" in ADM and dataset creation, and instead align with well-established theories of social justice[32] all the while adhering to principles of empirical verifiability, as for example in [KA21], [MBM+20], [HDSSL20]. A concrete example of value-free assumption would be adopting the notion of "merit" in ADM fairness [KA21].[33] Whereas, in MIR the most prevalent manifestation of the myth of value freedom is currently the lack of ethical training, resulting for example in critically overlooked aspects during dataset creation [MSW23], [SHD21]. In fact, Zhao et al. [ZAPX24] specifically identify value neutrality as one of the main issues in dataset curation in machine learning research at large, pointing out how practitioners often use value-laden terms such as "diversity" and "bias" without first giving a clear definition. In their research, Zhao and colleagues directly address this gap by applying principles from measurement theory to identify key considerations and provide recommendations based on social science insights. To the best of our knowledge our field has yet to produce precise guidelines for operationalisations of diversity for music datasets. In this regard, [ZAPX24] and [MLL+22] offer clear and potentially transferable insights for future MIR research.

Finally, it is important to highlight that (un)fairness in dataset creation and curation begins at data collection and annotation. In one respect, data is being collected to build (music) datasets without obtaining consent from content creators [MSW23], [MBB+23] raising clear ethical concerns related

to intellectual property and labour exploitation. Moreover, in terms of employment conditions, [MP22] argue that a combination of precarity and financial dependency leads to a situation where millions of *data workers* from the Global South [KLS21a], [WMG22] are alienated and made compliant [MP22]. This issue is exacerbated by the way historical *power relations* are embedded in artefacts like interfaces and performance metrics that restrict workers' autonomy [Whi23] and promote the standardisation of specific interpretations of data [MP22], which as a side-effect, also introduces bias. We thus believe that researchers, developers, and companies should avoid supporting platforms for large-scale data annotation that are well-known places of exploitation, such as *Amazon Mechanical Turk* [PS16], [AC11], where half of the workers in 2018 earned around or less than $2/h$, and only one in 25 workers earned more than $7.25/h$ [HAM+18].

**Transparency:** In their systematic review of transparency in music Gen AI, [BRGL+23] identified a lack of precise terminology as one of the main challenges in compiling their collection, so they adopted the policy-oriented definition of algorithmic transparency proposed by the IEEE [IEE17, p. 82]:

> Transparency is a characteristic which describes a process whereby information is requested and then disclosed completely within the limits of public law, without distortion, and with respect to the computational and cognitive capacities of the information recipient in order to enable those recipients to interpret the information so that they are able to make rational, informed decisions.

The authors propose that, based on the current literature in AI and MIR, transparency can itself be proceduralised in terms of explainability, interpretability and documentation.

According to [ADRDS+20], explainability methods are categorised into transparent models, which are inherently interpretable, and post-hoc explainability methods that use external techniques (model-agnostic or specific). Explanation methods and interpretable models can enhance audio applications throughout most MIR tasks, from music generation, to genre classification and recommendation systems. For a systematic review of explainability and interpretability methods in music (Gen) AI see [BRGL+23, p.13].[34]

While in terms of documentation, [BRGL+23, p.15], different methodologies are emerging in both research and industry. Proposed approaches involve both documenting datasets and models while others focus on AI services, reporting on purpose, performance, safety, and security. These documentation strategies improve data practices, facilitate optimal dataset selection, and enhance algorithm quality and trust in AI services. In European law, the AI Act addresses transparency by requiring all providers of general-purpose AI models to adhere to documentation obligations. These include:

> Information on the data used for training, testing and validation, where applicable, including the type and provenance of data and curation methodologies

---

[32]Such as diversity and inclusion [MBM+20], intersectional approaches to political economy [KA21], and critical race theory [HDSSL20]. All of these examples address intersectionality [Har15, p. 57], which, in addition to being itself a theory of justice, also serves as an analytical framework [CB20, p. 2].

[33]See [TRHH22] for how meritocracy developed into a value-free construct.

[34]Another valuable overview of XAI in MIR and related fields is offered by the first ICASSP workshop on Explainable Machine Learning for Speech and Audio https://xai-sa-workshop.github.io.

(e.g. cleaning, filtering etc), the number of data points, their scope and main characteristics; how the data was obtained and selected as well as all other measures to detect the unsuitability of data sources and methods to detect identifiable biases.[35]

Outside academia and policy-making, several initiatives are promoting transparency with regard to AI applications in the music industry. *Fairly Trained*, a non-profit organisation, aims to ensure fair treatment for human creators by distinguishing AI providers who obtain consent from data owners from those who do not, using a certification system based on data handling practices. *AI:OK* provides similar certifications for music products and services within the generative AI and emerging technologies landscape. Additionally, *Water&Music*, a collaborative platform, has developed the *Music AI Ethics Tracker*[36], a living document that examines and analyzes ethical statements on AI from the music industry.

**Bias:** As previously discussed also the term "bias" is surrounded by confusion about its meaning. [BHA+21] explain that, in order to understand the relationship between inequity and foundation models, we must recognise that FMs serve as intermediary assets adapted for user-impacting applications. They propose a coarse dichotomy distinguishing between intrinsic biases (inherent properties of foundation models that could cause harm) and extrinsic harms (harms arising from specific applications using these models). Whereas, in a seminal work from the early era of *value-sensitive design*, [FN96] provide a more nuanced taxonomy of bias. They use the term bias to refer to "computer systems that systematically and unfairly discriminate against certain individuals or groups of individuals in favor of others," and they categorise this into three types: technical, emergent, and preexisting bias.

*Technical bias* is due to limitations in technology, such as hardware constraints, flawed algorithms, issues with pseudo-random number generation, and issues with the formalisation of human constructs. *Emergent bias* in computer systems arises after the system's design is complete and is used in real-world contexts. This type of bias occurs because the system's original design assumptions no longer match the evolving conditions and characteristics of its user base, both at an individual and societal level, such as issues with user expertise[37], or cultural values. Building on Friedman and Nissenbaum's taxonomy, [DDGK18] interpret technical bias as an epistemological problem, as the methods AI practitioners use to analyze and formalise problems are themselves sources of bias[38]; and emergent bias as a dynamical feedback phenomenon, whereby it occurs within evolving societal changes and user interactions. Feedback loops in decision-making systems can reinforce and amplify biases over time,

necessitating a dynamic perspective to understand and mitigate these effects. For example, as also discussed in the following section, recommender systems are notoriously affected by a feedback loop on the emergent "popularity bias" [MAP+20].

Finally, *preexisting bias* arises from social attitudes and institutions, either through explicit intentions or implicit norms. Notably, it can cascade into every aspect of AI research and production, from the lack of diversity in developers' teams to sampling biases in datasets and model or system outputs. The previously-discussed Western-centric tradition in MIR risks to marginalise or incorrectly appropriate non-Western music, all the while a lack of engagement with historical ethnic and gender-based power relations in education and the music industry at large [THED+22], [WH19], [BAK23] can result in them being mirrored and amplified through the use of FMs in ADM, as for example in recommendation systems. Particularly to music recommender systems, [HCVKL23] found that there is minimal collaboration between the fields of academic computer science and critical social sciences and humanities. This, they argue, leads to a situation where specific preexisting biases, such as popularity and demographic biases (e.g., gender), are commonly identified in computational studies [DB22], whilst other biases like ethnicity and socioeconomic status are less studied. For this reason, research in recommender systems—of which music systems are a subset—is increasingly considering multidisciplinary and intersectional approaches [DJB+23]. In this context, increasing the presence of underrepresented demographics in STEM[39] is a crucial first step. However, as also discussed in the previous section (under "Epistemology"), true diversity goes beyond mere physical presence in research teams, as "what is desired is the kind of diversity that fully respects the values and interests of all citisens while protecting those of the most economically and politically vulnerable groups" [Har15, p. xi]. This, when developing FMs for music, first and foremost entails addressing the needs of (long tail) musicians and listeners.

*3) Adverse Effects on Humans:* FMs raise several questions on how they can affect the relationship between music and society. Music creation is not only a technical process but also, and mostly, a creative practice that binds itself to social and cultural contexts [BB18]. Thus, the adversarial impact of FMs on humans and society cannot be ignored. Current literature in MIR and related fields recognises the harmful impact of not centring design on humans (and, by extension, living beings) [FH19], [GJFA21], [HSH21], [HSC18b]. In this subsection, we call for a deeper exploration into the potential threats posed by FMs, extending a call made by [GJFA21] for music recommendation systems. We also explore how FMs that do not consider ethical issues at their design stage can amplify known effects such as isolation, replacement, reduced agency and taste manipulation.

> As platforms like Spotify become the primary means of distribution and circulation of music and other content, they begin to shape the ways in which music is produced and readied for them, either through

---

[35]Annex XI, Point (a), as referred to in Article 53(1), as approved by the European Council on May 14, 2024.

[36]Available for consultation at: https://www.waterandmusic.com/data/ai-ethics-tracker

[37]Friedman and Nissenbaum give the example of an ATM designed with extensive written instructions being used in a community with low literacy rates.

[38]Thus, bias can be even introduced in operationalisations of fairness, especially in the presence of value-free assumptions, such as "merit" in [KA21].

[39]Which includes addressing early on widespread gender bias in (music) technology education [Arm08]

explicit policies, rules, and guidelines they impose, or through more hidden acts of infrastructural and algorithmic politics [GJFA21]

As from the above quote about music platforms re-shaping the future of how we consume music, questions arise on more fundamental tools: the data-driven machine learning models that support such platforms. As applications of FMs are developed for such platforms, they too can shape the ways in which music is created, curated, and distributed. This shaping process can happen through the known assumptions, rules, and limitations surrounding the models, or through, what [GJFA21] describes as "infrastructural and algorithmic politics." Indeed, the opacity of such algorithmic decisions is evident not only at the deployment stage but also in understanding their long-term impact on individuals, the music industry, and society as a whole. Although [HSC18b] stated that MIR algorithms are not fatal to individuals, we argue that ill-defined FMs can bring about the erasure of cultures, sub-cultures and identities. This may occur in the hands of premature deployment of FMs without careful reflection on the scope of the models. [HSC18b] presents a grounding example where a community has to adapt its music style because existing software can only cater to a dominant culture's style, thereby threatening local music traditions and manipulating the future of the local music. Design decisions and constraints play an important role in evaluating the impact of FMs. It is necessary to not only express the scope of the model, but also evaluate how it might affect communities beyond it. For this, [GJFA21] asks of recommendation algorithms: "How will detaching people's taste-forming musical habits from their wider ecologies disrupt ordinary forms of social and cultural relationality, and the linking of music to other forms of experience?" Short-term goals such as performance gain are only as valuable as their term in the long run, it is indeed necessary to push for more human-centered long-term evaluation that can help understand the unknown effects that may persist and amplify without proper intervention. It should be noted that although not all unknown effects are adversarial, such as what is deemed as "emergent properties" in FMs, their exploration is necessary to not only avoid ill effects on humans, but also to promote the transparency of the algorithms.

One such reason for these adverse effects might be the disjoint relationship between MIR researchers and the end-users [HSC18b]. Lack of communication can drive a wall of isolation between the communities researchers take from versus give to. [HSC18b], [NML23] push for reducing this gap between MIR researchers and musician communities by engaging in conversation, implementing feedback, and connecting stakeholders to researchers. [HSC18b] provide another grounding example where a fictional company, Adaetal, scrapes data from an online community to generate music in the specific style of the community. Although the company is successful in the task, the contribution of Adaetal in preserving the traditional music remains in question. [HSC18b] asks "Is this research only benefiting Adaetal, to the detriment of the tradition they are using? How is the work of Adaetal

contributing to this tradition?" The implications of Adaetal's work remain adversarial to the preservation of the tradition; whilst Adaetal benefits from the recognition and commercial success of developing algorithms that mimic a traditional style, human artists from the community may face possible loss of livelihood.

Human replacement remains a major threat to human livelihood when it comes to the deployment of FMs. How do music FMs affect the role of musicians, producers, musicologists, or even the listener? [Mor21] suggests that the long-tail problem[40] might become more competitive for human musicians when the space is infiltrated by AI-musicians. This ideology is not limited to Western circles. As mentioned above, [HSH21] finds that the school of Mohism, an eastern philosophy, believes such excessive production of music is wasteful. It is evident that across the world, the idea of producing music for the sake of producing music could not only be meaningless, but also unfavourable to humans. Despite obvious ethical dilemmas, generative FMs are attracting increasing venture capital, given their potential to generate millions of tracks every day. In fact, [Mor21], [SIBT+19a] recognise that the reduced labour cost is a strong incentive for music platforms to adopt such generative models. It is necessary for researchers developing such FMs to reflect on these potential issues before they deploy their models.

While FM companies demand an act of faith that their models will create new roles and skills for musicians rather than replace them, a callback to the transparency of such models is important. As [GJFA21] points out, the lack of transparency regarding the fundamental decision-making done by algorithms might lead to a lack of agency for its users: the musicians its meant to support. This lack of control is an important issue that resonates in conversations with musicians [NML23]. Musicians recommend more open-source code for greater transparency of their controls as well as to enable changes as required for their task (ibid.). "Control" is not only limited to the capabilities of the FMs, but also speaks to the human intervention that can be supported. Research suggests that musicians do not find that creative pursuits with AI are meaningful without evidence of human labour in the production process (ibid.). "Intention and choice" are as valuable as the output as well, signifying the human touch and agency is necessary for musicians to connect with the music (ibid.). As one participant in a study on perception of AI music says, "For me, creativity also involves the decision-making in a big way. And then to determine where to end things. It doesn't seem that my experience with AI so far affords these possibilities" [NML23].

Lack of transparency and control are not only detrimental to musicians, but can also unknowingly impact the community of listeners. In its most obvious form, taste manipulation can adversely affect the perception and reception of music. As seen in another example developed by [HSC18b], a fictional Digital Audio Workstation (DAW) company adopts technology

---

[40]The long tail is a common problem, found especially in recommendation, where the distribution of musicians' popularity is largely skewed. Few musicians gain immense popularity, while the rest remain to have incomparable popularity (the long tail) [Mor21].

developed to find similar music. However, the limitation of the training data disallows any music from a time signature apart from 4/4 to work well with their model. This manipulates how music is produced, at the detriment of the local music in a region. Such design decisions can reshape what is acceptable and not acceptable. However, this raises questions on who has the right to determine what is acceptable music or not. Hence, it is not only necessary to carefully reflect on design decisions, but also to document them indiscriminately to allow understanding of any changes which may reflect certain biases and affect the taste of listeners [HSC18b]. The historical prevalence of western music training datasets can also promote culture homogeneity because FMs trained on such data can be deployed as "universal" models, algorithmically enforcing western styles to become the dominant, and the only, style of music available as AI-music. This might promote an individualistic perception of music listening, which stands in contrast to the real world where our music perception is not only a product of our individual characteristics but also of the culture and environment we are a part of [GJFA21]. At its worst, this hyper-focused attention on individuals and neutralisation of culture might increase the echo-chamber effect highlighted by [Mor21]. This radicalisation of "personal attention" has been discussed in other fields where urgent calls have been made to redefine and shift our motivations to develop technology that feeds on user engagement [dlTPVB24], [Her22].

*4) Take Responsibility:* In this subsection, we will discuss the responsibilities of researchers, music educators, policymakers and companies to address ethical concerns.

**Researchers in MIR** must recognise their ethical responsibilities regarding music foundation model development. Unfortunately, genuine investigations into these matters are unlikely from large corporations primarily driven by capital interests since corporate ethical inquiries often serve as strategic manoeuvres to buffer against criticism [Bie20]. Thus, many researchers suggest the research community should self-appoint as ethical regulators to face the potential social impact and ethical concern [Mor21], [Ter20].

In ML and MIR communities, scholars have begun addressing these ethical considerations. The data-centric attribute of ML can somehow perpetuate current social biases embedded in training data [BS16]. Researchers must critically evaluate and document the capability, bias and limitations of their algorithms, ensuring transparency for users and stakeholders. The IEEE community advocates processes guided by ethical principles to mitigate bias and increase transparency [Han18], [CH19]. Besides, there are many ethical discussions in MIR. Some have focused on recommendation bias and fairness [HSC18a], [GHMS19], gender representation in music streaming [EDCB20], and user data concerns [SRD14], [CKM+19]. Legal aspects of AI-generated music, such as the use of copyright-protected training sets and transparency regarding AI involvement in music creation, have also been examined [SIBT+19b].

Researchers can enhance the ethical standing of music AI by improving the interpretability of the foundation model. Interpretability [DVK17] also known as explainability [MRW19], refers to making AI behavior and outcomes understandable.

This improvement helps in correcting errors, resolving biases, and providing causal inferences [ZUR17]. Additionally, it addresses music copyright issues by clarifying how the system generates specific content, which aids in attributing credits appropriately to users, programmers, or musicians involved in the training set [SIBT+19b].

Another measure researchers can take in the context of foundation models for music is to prioritise fairness and transparency discussed in previous subsections.

**Music education.** Given the capabilities of foundation models in music creation and production, and considering the significant implications for career trajectories and opportunities in the music ecosystem, music educators are likely to incorporate AI-assisted composition and the use of foundation models and prompt engineering into their curricula. However, given the rapid pace of technological advancement, AI tools are expected to replace entry-level jobs in small media companies or recording studios. So such a measure is only a temporary measure. Developers and administrators must consider the long-term implications of deploying music AI and explore ways to mitigate adverse consequences. Additionally, the current legal framework might be necessary to be adapted to better support artistic innovation and creativity in this evolving environment [SIBT+19b].

**Policymakers** can take multiple measures to address the ethical concerns of foundation models to the music industry while boosting creativity, such as introducing rights and obligations on "AI-generated music" labels, training set record keeping for copyright supervision as well as enhancing personality rights to prevent DeepFake [41].

The first measure is to label AI-generated songs on streaming platforms. AI-driven tools, especially foundation models, have been increasingly used in music creation with better and better capabilities, from composing and mixing to streaming. These models shape users' music consumption behaviours and can foster unethical practices, such as promoting AI-generated "artists" to boost revenue at the expense of human artists [EFJ+19]. Transparency about AI's role in music creation and recommendation can empower both artists and listeners. For instance, flagging AI-generated music or offering options to avoid such music could be beneficial [AW18]. This measure is all compatible with consumer rights such as the UK Consumer Rights Act 2015, leaving enough information for users to make decisions on their own. However, defining AI involvement levels is challenging. Decomposing music creation into stages and accurately documenting AI's contribution is complex and varies individually. This process may stifle creativity and create barriers for artists, especially those lacking institutional support who rely on AI-assisted compositional tools. Moreover, it is unclear whether consumers understand or care about AI's extensive involvement in music creation. Despite these challenges, further research is needed to assess the potential harms of not informing consumers about AI involvement and to develop effective transparency measures. Researchers must balance transparency with practical considerations to ensure

---

[41]Artificial Intelligence and the Music Industry – Master or Servant? A report from British All-party Parliamentary Group (APPG) https://www.ukmusic.org/research-reports/appg-on-music-report-on-ai-and-music-2024/

ethical and fair practices in AI-driven music creation and consumption [SIBT+19b].

Besides, policymakers should require full documentation of all data used by developers. For commercial companies, explicit permission from musicians or music companies is required to use copyrighted datasets. Governments should consider market access regulations and include models developed in countries or regions where copyright protection is relatively lax especially when copyright compliance can hardly be demonstrated for model training data. Furthermore, policymakers shall introduce a specific human right to protect musicians from being affected by DeepFake including misappropriation and false endorsement of their voice, image, name, and likeness (VINL) [42].

**Streaming platforms**, though, can leverage AI-generated music to potentially reduce payouts to musicians, such practices could stifle creativity and long-term music diversity. In response, leading companies like Spotify and Deezer have implemented measures to reportedly better reward professional musicians [43] [44] In 2023, Deezer introduced an "Artist-Centric Model" that prioritises streams from established artists, actively searched-for songs, and combats fraud. This resulted in the removal of millions of low-quality tracks generated by AI or music amateurs and a projected 7-10% royalty shift towards real artists [45].

*B. Copyright*

Intellectual Property (IP) is a type of property that includes intangible creations. Intellectual creations such as names of products or brands, designs or looks of products, inventions, literary works, films, art, photographs or music compositions are examples of forms of intellectual property. Depending on the nature of the IP, different types of protection apply, such as designs and trademarks, patents or copyright. By preventing third parties from using or copying the IP without authorisation from the IP owner for a limited period of time, the intent of IP protection is to encourage creativity and innovation. Producing the initial creation requires an investment and the IP protection allows the creators to derive a benefit from their original creations, therefore providing an incentive to put in the effort in the first place. Analogous to a patent protecting an invention, copyright is an instrument that allows protection of creative works. The Statute of Anne (1709) [46] is often regarded as the first act of law to provide copyright regulation. Since then copyright protection has been introduced in most jurisdictions and copyright has been considered a fundamental human right since the Universal Declaration of Human Rights (1948). In the music domain, copyright typically applies to compositions and sound recordings.

Copyright is sometimes described as a negative right because it only limits the rights of others. Copyright law disenfranchises a community of derivative creativity. Longen argues that "term limit extensions, increased protectionist treatment of secondary works online, and the functional lack of access to proper licensing mechanisms have rendered users' rights impotent." [Lon23] Finding a balance between encouraging creative work through copyright and fostering innovation through a more open creative commons has been a long-running controversy. [Les01]

Composers, songwriters, performing musicians and record producers have benefited from the protection of their work and a way of earning income from it. Recent progress in generative AI saw this technology acquiring the ability to produce outputs of comparable quality to human creators. It, therefore, becomes a potential form of competition with human creative output and a threat to the related sources of income. This brings forth copyright concerns mostly in two areas: usage of copyrighted material for training AI systems and the eligibility of the output of Gen AI systems for copyright protection.

In the former case, the concern consists in determining whether or not training an AI system on copyrighted material without obtaining a license from the right holders is a breach of copyright. Artists and the wider creative community typically argue that it constitutes a violation of copyright, while some commercial organisations offering generative AI services argue it is not. Multiple lawsuits have already been filed for alleged copyright violations [Sam23], and more may follow as new companies and applications make AI-generated music widely available, and legal precedent has not been set yet. Interestingly, independently of legal obligations, some commercial providers of GenAI services have voluntarily adopted the position that they will only train on licensed content and offer compensation to the creators whose content is used as part of the training set. This is motivated by various factors: the recognition that training data is a major contributor to the value created by machine learning systems, the belief that artists whose creative works are used for training should be compensated for their contribution to the value creation, pressure from investors who cannot risk the legal exposure of copyright cases, and the high value of corporate reputation.

There is also the concern of determining whether content generated by Gen AI system can be protected by copyright. Given that copyright is traditionally aimed at protecting human creativity, there is usually a requirement that sufficient human creative input is part of the creation process for the work to qualify for copyright protection. It therefore appears unlikely that content fully machine-generated may be eligible for copyright protection. However, copyright institutions need to determine how much creative human input is enough to qualify. As an example, the US Copyright office in its AI policy guidance recalls its view that "copyright can protect only material that is the product of human creativity"[47]. As a result, "If a work's traditional elements of authorship were

---

[42][41]

[43]Modernizing Our Royalty System to Drive an Additional $1 Billion toward Emerging and Professional Artists  Website.+

[44]Universal Music Group and Deezer to launch the first comprehensive artist-centric music streaming model. Website.

[45]Music in the Air – Focus on monetisation, Emerging Markets and AI; updating global music industry forecasts Website.

[46]Statute of Anne, London (1710), Primary Sources on Copyright (1450-1900), eds L. Bently & M. Kretschmer, www.copyrighthistory.org

[47]Copyright Registration Guidance: Works Containing Material Generated by Artificial Intelligence https://copyright.gov/ai/ai_policy_guidance.pdf

produced by a machine, the work lacks human authorship and the Office will not register it." Based on this rationale, the US Copyright Office determined that in the case where a user only provides a prompt to a Generative AI system, the "traditional elements of authorship" are provided by the machine, not the user, and therefore that the resulting output cannot be registered for copyright protection. Although this determination may be easy to make in simple scenarios, it will be vastly more difficult in cases where the interplay between a user/creator and the GenAI system is more intricate. Auditing creations that potentially contain elements of GenAI in order to determine their eligibility for copyright protection will also constitute a monumental challenge, partly because of the scale at which it will need to be done and partly because technical tools to rely on something else than an honest disclosure from the applicant do not currently exist.

Copyright infringement in the form of plagiarism is not a risk that is new or specific to Gen AI systems, but the broad availability of this technology makes it possible on a scale much larger than before. Empirical evidence shows that music's distinctive pitch and rhythm sequences can be memorised by current Gen AI models and are recognizable when re-expressed in new contexts or outputs. In other domains, AI systems tend to paraphrase or re-express learned information from the text and images, providing a defence against infringement claims. The more abstract character of music perception poses a challenge for the makers of Gen AI services, whose original intent is not to infringe copyright, as well as for the rights holders. Gen AI providers may need to establish appropriate guardrails to prevent their systems from plagiarising copyrighted material.

From the rights holders perspective, since Gen AI may be susceptible to produce outputs that infringe copyright in large numbers of small transactions, it may be impractical to discover and take legal action on a case-by-case basis. In response, copyright owners across the space are "opting out" (denying AI companies the right to use their copyrighted material for training). While this is entirely understandable, it could have severe adverse effects in the long term. It may be in the music industry's own interest to avoid another Naptser-like situation, with aggressive use of legal action against any and all training of AI models on copyrighted material [MB03]. An alternative strategy is embracing AI while trying to infuse the space with desired values and offering attractive solutions in terms of technology and access to desirable data. A failure to do so could yield a competitive advantage to actors who, by ignoring rights, would thus be able to train better models at a significantly lower cost. Over time, this could effectively dilute the concept of copyright and take away the abilities of industry, citizens and artists to influence and control the Gen AI space.

Large language models for music, through the application of "disentanglement" of rhythm, style, emotion, timbre, and other musical elements, may lead to an expansion of current concepts of copyrightable content beyond existing copyright coverage of composition, lyrics, and recordings. An important question for Gen AI and for music copyright is: to what extent are generated music elements derivatives of original copyrighted music recordings and compositions vs. original expressions of learned musical knowledge?

Because of the scale and new nature of copyright challenges raised by Gen AI, the development of appropriate and scalable technical tools and solutions may be required to address them. Improved technology could be used to avoid copyright infringement as well as to make copyright infringement detection practical in the new era of mass production by AI. For example, technology performing the following tasks in a tractable, robust and scalable way: identification of GenAI output, training set membership inference, content provenance, authenticity and authorship tracking. As of today, most of these tools either do not exist or are not mature enough for deployment at scale, which undoubtedly warrants more research.

Even though Gen AI raises some profound copyright-related questions, it is undeniable that this technology possesses enormous potential to engender a new generation of creative tools that will empower artists and creatives across domains, including music. The desirable outcome of copyright debates is one where all stakeholders (artists, right holders, Gen AI companies, legislators etc.) establish a model for usage and management of copyrighted material that allows the adoption and continued development of Gen AI technology while preserving the ability of the creative sectors to thrive.

## C. Personality Rights

In addition to the copyright-related challenges discussed in the previous section, Gen AI systems bring additional risks and challenges related to a different type of rights: Personality Rights. Personality Rights, which are sometimes also referred to as right of publicity, give an individual control over the commercial use of their name, image, likeness, and any traits that may allow to identify them. This includes an individual's voice, which is of particular interest and relevance in the music domain.

Voice, image and likeness cloning technology is particularly problematic in the music domain. Owing to recent progress in the development of Gen AI technology, voice cloning systems are now capable of cloning both spoken and singing voice with a level of realism that makes a synthetic output virtually indistinguishable from an authentic recording. Given that these systems are automated, broadly available and that music and sound recordings can be accessed (perhaps illegally) on the internet to train them, this means that voice cloning is already being performed at scale and in a variety of scenarios. Doing so without the authorisation of the individual (artist) being cloned constitutes a profoundly unethical behaviour. Some actors are totally oblivious to ethical considerations. For instance, there exists today a number of commercial organisations that sell services to clone known artists voices, without having received prior authorisation, or proposing any form of compensation. The risks associated with such a questionable use of this technology in the music domain include reputational harm to the artist, misinformation, and deceit of the artist's audience or the general public, to name only a few. Because the likelihood that these risks materialise

is significant, and the consequences can be devastating for the artist and for society at large, it is extremely important that they are acknowledged, considered and addressed jointly by the Gen AI and music communities, and, possibly, also by the legislators.

There is very little debate that using someone's likeness without their prior consent cannot be deemed as an ethical or acceptable behaviour. In most jurisdictions, a legal framework establishing personality rights and their protection already exists (though it may vary depending on the jurisdiction). However, because of the specificity of the threats resulting from the broad availability of likeness cloning technology powered by Gen AI, certain jurisdictions recognise the need to introduce AI-specific regulation to protect personality rights. In the state of Tennessee (USA), the ELVIS act is generally acknowledged as the first piece of legislation voted into law specifically designed to protect musicians against unauthorised use of their vocal likeness in scenarios such as deep fakes, voice cloning, etc. This law has been criticised for its broad reach and potential conflict with other legal rights[48]. With the introduction of any new legal protections comes a risk is that may be used or abused beyond the original intent of the legislation. Hypothetically, a famous person could assert personality rights claims to attempt to block the career of another artist, regardless of how they came to sound alike [Sta94, p. 368]. In the context of AI and voice-cloning, similarity to human artists may become a tunable parameter that intensifies the already existing gray area defining the boundaries of identity. Nevertheless, one may anticipate that other jurisdiction may also consider the introduction of similar legislation, given the growing prominence of voice cloning technology.

While the law provides legal constraints, voluntary initiatives can raise ethical standards beyond (or in spite of) legal requirements. Some providers of music AI technology are already adopting an ethical position that goes beyond strict legal requirements [49]. This type of approach is broadly in line with an open call from musicians and a vast number of music industry organisations for more consideration being given to musicians' rights, livelihood and perspective in the development of Gen AI technology [50]. The development of technical tools may be helpful to not only combat harm when done but also prevent harm being done in the first place. In particular, tools such as the following are identified as possible aids: Automatic identification of a likeness (e.g. image or vocal), detection of cloned vs. authentic voice [DMBHM24], content provenance and authenticity tracking [C2P], and perhaps attribution models to fairly compensate authorised use of a likeness. For the most part, these tools are in their infancy and will require further research and development before they can be effective guardrails at scale.

Another use case that raises arduous ethical questions consists in using likeness cloning technology for deceased artists.

There is already evidence that the public is receptive to this type of use [CDRS20], [Bae], which means we can most likely expect that it is a commercial offering that may experience some growth in the future. Personality rights are fundamentally designed to give a (living) individual rights over their likeness in the general sense of the term. After death, it is often the deceased's estate that is in a position to make decision regarding the potential use of the deceased's likeness. Should there be limits to what can be deemed ethical or acceptable to do with the deceased's likeness? Especially if they did not explicitly express their will during their lifetime? How should those limits be established? Should there be time limits for protection as in copyright? Should artificial personalities be protected, and could, say, 1 billion personalities be detected or distinguished?

Despite a number of risks and pending ethical questions, and much like the case of copyright, provided it is done in an ethical and responsible way, likeness cloning technology can be an enabler of valuable and innovative products or experiences. For example voice cloning alone could be used to: create personalised artist messages to fans, games, localise content (i.e. adapt to local languages), produce avatars to reach audiences that are not able to attend concerts etc. The artist community has already started embracing the technology, showing the exciting path of new creative possibilities for the use of technology that is respectful of artists' personality rights.

## VII. CONCLUSION & DISCUSSION

In conclusion, this survey has articulated the transformative potential of foundation models (FMs), particularly large language models (LLMs) and latent diffusion models (LDMs), in the realm of music. These models exploit self-supervised learning to process complex musical data, significantly enhancing music information retrieval, generation, and multimodal interactions. The integration of these technologies into the music industry promises substantial benefits across various sectors, including entertainment, education, therapy, and cultural preservation. The survey has delineated several core areas where FMs are making an impact: representation of music, downstream tasks, pre-training strategies, adaptation techniques, and the pressing issues of ethics and copyright in music AI applications. The outlined advancements underscore the need for advanced LLMs, visual-language models, and other interdisciplinary approaches to address the challenges of long-sequence modelling and domain-specific knowledge architecture. In addition, we collect the available corpus for music FM development, including a number of large, high-quality and diverse datasets in multiple modalities and evaluation metrics for music understanding and generation. Besides, we suggest that due to the strong socio-cultural and music industry connections, ethical issues and data copyrights need to be taken into account.

Targeted at computer music researchers unfamiliar with LLMs and LDMs and ML researchers interested in music applications yet inexperienced with traditional music processing techniques, this survey serves as a foundational reference. It

---

[48]ELVIS Act Needs to Be 'Returned to Sender' https://www.realclearpolicy.com/articles/2024/02/29/elvis_act_needs_to_be_returned_to_sender_1015123.html

[49]https://aiformusic.info

[50]https://www.humanartistrycampaign.com

aims to guide future research directions, encouraging a deeper engagement with the ethical implications and societal impacts of deploying foundation models in music, thereby shaping a responsible trajectory for the evolution of music technologies.

Our exploration of FMs' **representation in music** has revealed that current models predominantly focus on a limited spectrum of musical representations such as spectrum, waveform, MIDI, and ABC notation. Besides, we observe a considerable gap in fully exploiting multimodal representations that unify symbolic music, audio, text, and music score images within FM development.

On the **music applications** front, music pre-trained models demonstrated robust capabilities across a wide array of music-related tasks—ranging from tonality and harmony analysis to sophisticated music generation and medical applications. These include chord and key detection, melody extraction, pitch detection, and more dynamic tasks like music source separation and emotion recognition. The potency of foundation models in NLP and CV communities, exemplified by systems like chatGPT and stable diffusion, in addressing diverse and previously unseen natural language or image generation tasks suggests a promising future. Products such as Qwen-audio and SunoAI, which offer extensive music understanding and generation capabilities, indicate that multimodal LLMs combined with LDMs or employed as Music Agents could substantially impact music technology. This burgeoning field necessitates exploration of the advanced LLM and LDM techniques in music to enable models' versatility, along with ethical dimensions ensuring that advancements enhance.

Moreover, the survey highlights several significant technical advances and remaining challenges in leveraging LLMs and LDMs for music, including pre-training paradigms, domain adaptation, in-context learning, model architectures, controllability, music agents and scaling law. **Training paradigms**, including contrastive learning, generative pretraining, and mask language modelling, have shown potential in learning from large-scale, unsupervised music data. However, all of these models require fine-tuning on downstream tasks, and models developed with the combination of pre-training and instruction tuning are yet to be fully realised. Besides, combining the music domain to develop pre-training targets has not been fully discovered.

**Domain adaptation techniques** such as adaptor and prefix tuning have proven effective for fine-tuning pre-trained language models on downstream music tasks, enhancing the multimodal capabilities of these models without extensive retraining. These methods facilitate the integration of music-specific knowledge into LLMs, allowing for nuanced understanding and generation of musical content. However, the research field still lacks a comprehensive understanding of the potential of MLP, Q-former or other architecture on music FMs.

**In-context learning (ICL)** techniques have enabled foundation models to perform tasks with little to no task-specific training, leveraging the model's ability to generate contextually appropriate responses based on provided examples. This approach is particularly promising for music applications where training data may be scarce or highly diverse. However, the

development of supervised fine-tuning data, which gives model music ICL capabilities, is still rare. And how music capabilities at different levels (e.g., acoustic information, performance-level information, high-level text description, etc.) can be improved by which type of advanced prompt techniques, such as chain-of-thought (CoT) remain underexplored.

**The architectural** of most current foundation models, primarily based on the Transformer architecture, supports extensive scalability but faces challenges with long-sequence music data. Innovations in model architecture are required to better handle the complexities of musical structure and longer context windows. Although the GPT-4 and Claude-3 can handle more than 100k tokens, there is still much room for improvement in the ability to handle long sequences. Specific architectural adjustments or algorithms may be needed to enhance the modelling and exploitation of long contextual information. Another issue is that current foundation models are mainly decoder-only architectures, and training such models places high demands on the quality of the music-text tokeniser. Whether the Encoder-decoder architecture is an alternative is also worth exploring.

**controllability** remains a critical aspect, with current models needing better mechanisms to ensure music generation adheres to specific stylistic or structural constraints. The integration of explicit control mechanisms and the exploration of interpretable models will enhance the utility and applicability of foundation models in music, making them more transparent and aligned with users' creative intentions.

**Scaling laws** in LLMs have underscored the importance of model size and training data scale for achieving superior performance. Training good foundation models is very challenging due to huge computational consumption and sensitivity to data quality and training techniques. The scaling law solves this problem to some extent. However, the abilities of musical FMs can be inferred from small models, which are emergent abilities that need further exploration to fully leverage their potential in creative and analytical tasks. Further, the tokenisation of the music may greatly influence the ratio of training resources to training data and model parameters. The exact link between the two needs to be explored.

Lastly, **music agents** represent a promising area of development, transitioning from semi-autonomous to fully autonomous systems that can interact with users, manage complex workflows, and integrate various musical tasks and modalities seamlessly. These agents are also poised to revolutionise how we interact with and create music, offering personalised and adaptive music experiences.

Apart from the training methodology mentioned above, the efficacy of foundation models for music heavily depends on **the quality and diversity of the datasets** used for training and evaluation. Challenges related to the size, diversity, and copyright restrictions of music datasets necessitate rigorous efforts in data collection, cleaning, mixing, and curriculum learning to ensure that models can generalise effectively across varied musical backgrounds and cultures. Despite the availability of extensive, open-source, high-quality music audio and scores encompassing various epochs, instruments, and genres, there remains a notable gap in the strategies for data preparation,

such as cleaning and instruction tuning for multimodal data involving diverse musical tasks. Additionally, the deployment of these models requires the development of strategies that appropriately balance data to avoid catastrophic forgetting and continuously adapt to new musical trends and user preferences.

Additionally, the security aspects of models capable of distilling training data require rigorous examination. This scrutiny enables music publishers to assess the legality of training datasets and simultaneously exposes models to the risk of revealing proprietary training data. An illustrative case is a lawsuit initiated by Universal Music Group Publishing, Concord Music Group, and ABKCO against Anthropic, which involved the unauthorized use of copyrighted material in 500 instances for training AI models [51].

On the **evaluation** front, the generalisation of these models, whether in understanding or generation tasks, requires more uniform and comprehensive metrics. Often, models are only tested on a subset of downstream tasks, and many benchmarks use music audio from public datasets that the models might have already been exposed to, leading to potential data contamination and leakage. Furthermore, there is an urgent need to address biases in models that have predominantly been trained on Western music, neglecting the vast array of world music, which raises significant ethical concerns.

Ethically, the use of FMs in music necessitates careful consideration of **ethics issues and copyright issues**, including transparency, accountability, and bias. By fostering collaboration among researchers, practitioners, and policymakers, we hope to better utilise the potential of AI in enriching musical experiences while ensuring its benefits are widely accessible to everyone ethically.

This survey underscores the importance of interdisciplinary collaboration in advancing music foundation models. We call upon researchers in computer music and machine learning to engage with ethical research and development practices rigorously. By doing so, we aim to bring transformative changes to the music industry and human musical culture through the responsible advancement of music foundation models.

## VIII. Acknowledgement

## Appendix

### A. Division of Work in Section Order

- Introduction: Simon Dixon, Yinghao Ma
- II Representations of Music:
  - II-A Music Perception & Notation: Charalampos Saitis
  - II-B Computer Representations of Music
    * II-B1 Acoustic-level Representations of Music: Huan Zhang, Yinghao Ma
    * II-B2 Symbolic Music Format & Their Content: Elona Shatri, György Fazekas
    * II-B3 Transformed Symbolic Music Representations for Foundational Models: György Fazekas, Huan Zhang
  - II-C Multimodal Music Representations: Shangda Wu, Yinghao Ma
- III Applications
  - III-A Music Understanding
    * III-A1 Traditional Music Information Retrieval Tasks: Emmanouil Benetos
    * III-A2 Multimodal Music Understanding Tasks: Shangda Wu, Xingjian Du, Yinghao Ma
    * III-A3 Vocal Music Understanding: Jiawen Huang, Zhizheng Wu
  - III-B Music Generation
    * III-B1 Symbolic Music Generation: Zeyue Tian
    * III-B2 Language Model for Acoustic Music Generation: Shun Lei, Zhiyong Wu
    * III-B3 Audio Diffusion Model: Max W. Y. Lam, Shiyin Kang
    * III-B4 Vocal Music Generation: Jiawen Huang, Zhizheng Wu
  - III-C Music Therapy & Medical Applications: Bleiz MacSen Del Sette, Chenghua Lin, Yizhi Li
- IV Technical Details of Foundation Models
  - IV-A Model Design & Pre-training Strategies: Christos Plachouras, Julien Guinot, Ruibin Yuan, Ziyang Ma, Wenhao Huang,
  - IV-B Music Domain Adaptation for Foundation Models: Ge Zhang, Wenhu Chen, Xingwei Qu, Yinghao Ma
  - IV-C Audio Tokenisers: Qiuqiang Kong
  - IV-D Model Architectures: Qiuqiang Kong
  - IV-E Interpretability & Controllability on Music Generation: Gus Xia, Li-wei Lin, Ziyu Wang
  - IV-F Foundation Models as Music Agents: Xu Tan
  - IV-G Scaling Laws for Music: Xu Tan
  - IV-H Additional Future Improvements: Anton Ragni
- V Datasets & Evaluation
  - V-A1-4 Music Datasets: Yinghao Ma

---

[51]Universal Music Sues AI Company Anthropic for Copyright Infringement - Levi's Sues Coperni for Trade mark Infringement. https://intellectual-property-helpdesk.ec.europa.eu/news-events/news/universal-music-sues-ai-company-anthropic-copyright-infringement-levis-sues-coperni-trade-mark-2023-10-26_en

*B. Author List*

*1) Main Coordinators:*

- Yinghao Ma, *Student Member, IEEE*

*2) Coordinators:*

- Emmanouil Benetos, *Senior Member, IEEE,*
- Fabio Morreale
- Shiyin Kang, *Member, IEEE*

*3) Main Contributors:*

- Anton Ragni
- Bleiz MacSen Del Sette
- Charalampos Saitis
- Chris Donahue
- Christos Plachouras
- Elona Shatri
- Huan Zhang
- Ilaria Manco
- Jiawen Huang
- Julien Guinot
- Liwei Lin
- Luca Marinelli
- Max W. Y. Lam
- Megha Sharma
- Qiuqiang Kong, *Member, IEEE*
- Roger B. Dannenberg, *Member, IEEE,*
- Shangda Wu
- Shih-Lun Wu
- Shun Lei, *Student Member, IEEE*,
- Simon Dixon, *Senior Member, IEEE*
- Wenhu Chen
- Xu Tan, *Senior Member, IEEE*
- Yizhi Li
- Zeyue Tian
- Zhiyong Wu, *Member, IEEE*

*4) Contributors:*

- Anders Øland
- Chenghua Lin
- Ge Zhang
- György Fazekas
- Gus Xia
- Ruibin Yuan
- Shuqi Dai
- Wenhao Huang
- Xingjian Du
- Xingwei Qu
- Zhizheng Wu
- Ziyang Ma, *Student Member, IEEE*
- Ziyu Wang

REFERENCES

[AAA+23] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

[AAAK23] Sara Atito, Muhammed Awais, Tony Alex, and Josef Kittler. Group Masked Model Learning for General Audio Representation. In *2023 IEEE International Conference on Image Processing (ICIP)*, pages 2600–2604, Kuala Lumpur, Malaysia, October 2023. IEEE.

[AANK23] Sara Atito, Muhammed Awais, Srinivasa Nandam, and Josef Kittler. GMML is All You Need. In *2023 IEEE International Conference on Image Processing (ICIP)*, pages 2125–2129, October 2023.

[AAW+22] Sara Atito, Muhammad Awais, Wenwu Wang, Mark D Plumbley, and Josef Kittler. Asit: Audio spectrogram vision transformer for general audio representation. *arXiv preprint arXiv:2211.13189*, 2022.

[AAW+24] Sara Atito, Muhammad Awais, Wenwu Wang, Mark D. Plumbley, and Josef Kittler. ASiT: Local-Global Audio Spectrogram vIsion Transformer for Event Classification, March 2024. arXiv:2211.13189 [cs, eess].

[AC87] Philip E Agre and David Chapman. Pengi: An implementation of a theory of activity. In *Proceedings of the sixth National conference on Artificial intelligence-Volume 1*, pages 268–272, 1987.

[AC11] G Adda and Kevin Bretonnel Cohen. Amazon mechanical turk: Gold mine or coal mine. *Computational Lingustics*, 37(2):2–10, 2011.

[ACSS23] Jonah Anton, Harry Coppock, Pancham Shukla, and Björn W Schuller. Audio barlow twins: Self-supervised audio representation learning. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.

[ADB+23] Andrea Agostinelli, Timo I Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, et al. MusicLM: Generating music from text. *arXiv preprint arXiv:2301.11325*, 2023.

[ADF+23] Rohan Anil, Andrew M Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, et al. Palm 2 technical report. *arXiv preprint arXiv:2305.10403*, 2023.

[ADH+21] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6836–6846, 2021.

[ADL+22] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. *Advances in neural information processing systems*, 35:23716–23736, 2022.

[ADM+23] Mahmoud Assran, Quentin Duval, Ishan Misra, Piotr Bojanowski, Pascal Vincent, Michael Rabbat, Yann LeCun, and Nicolas Ballas. Self-Supervised Learning from Images with a Joint-Embedding Predictive Architecture. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15619–15629, June 2023. ISSN: 2575-7075.

[ADRDS+20] Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bennetot, Siham Tabik, Alberto Barbado, Salvador García, Sergio Gil-López, Daniel Molina, Richard Benjamins, et al. Explainable artificial intelligence (xai): Concepts, taxonomies, opportunities and challenges toward responsible ai. *Information fusion*, 58:82–115, 2020.

[AFL23] Omri Avrahami, Ohad Fried, and Dani Lischinski. Blended latent diffusion. *ACM transactions on graphics (TOG)*, 42(4):1–11, 2023.

[Agr98] Philip E Agre. Toward a Critical Technical Practice: Lessons Learned in Trying to Reform AI. In *Social Science, Technical Systems, and Cooperative Work*, pages 131–157. Psychology Press, 1998.

[AJFF+23] Pablo Alonso-Jiménez, Xavier Favory, Hadrien Foroughmand, Grigoris Bourdalas, Xavier Serra, Thomas Lidy, and Dmitry Bogdanov. Pre-training strategies using contrastive learning and playlist information for music classification and similarity. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.

[AJH+21] Jacob Austin, Daniel D Johnson, Jonathan Ho, Daniel Tarlow, and Rianne Van Den Berg. Structured denoising diffusion models in discrete state-spaces. *Advances in Neural Information Processing Systems*, 34:17981–17993, 2021.

[Alg19] AlgorithmWatch, Bertelsmann Stiftung, Open Society Foundations. Automating Society – Taking Stock of Automated Decision-Making in the EU, 2019. Available at https://algorithmwatch.org/en/automating-society/ (Retrieved 19/06/2024).

[Ano24] Anonymous. Naturalspeech 2: Latent diffusion models are natural and zero-shot speech and singing synthesizers. In *The Twelfth International Conference on Learning Representations*, 2024.

[AP21] Gunjan Aggarwal and Devi Parikh. Dance2music: Automatic dance-driven music generation. *CoRR*, abs/2107.06252, 2021.

[APS05] J-J Aucouturier, François Pachet, and Mark Sandler. "the way it sounds'': Timbre models for analysis and retrieval of music signals. *IEEE Transactions on Multimedia*, 7(6):1028–1035, 2005.

[ARA+22] Reza Yazdani Aminabadi, Samyam Rajbhandari, Ammar Ahmad Awan, Cheng Li, Du Li, Elton Zheng, Olatunji Ruwase, Shaden Smith, Minjia Zhang, Jeff Rasley, et al. Deepspeed-inference: enabling efficient inference of transformer models at unprecedented scale. In *SC22: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–15. IEEE, 2022.

[Arm08] Victoria Armstrong. Hard bargaining on the hard drive: gender bias in the music technology classroom. *Gender and Education*, 20(4):375–386, 2008.

[ASV+21] Kat R Agres, Rebecca S Schaefer, Anja Volk, Susan van Hooren, Andre Holzapfel, Simone Dalla Bella, Meinard Müller, Martina De Witte, Dorien Herremans, Rafael Ramirez Melendez, et al. Music, computing, and health: a roadmap for the current and future roles of music technology for health care and well-being. *Music & Science*, 4:2059204321997709, 2021.

[Ata23] Lilac Atassi. Generating symbolic music using diffusion models. *arXiv preprint arXiv:2303.08385*, 2023.

[ATM21] Haider Al-Tahan and Yalda Mohsenzadeh. Clar: Contrastive learning of auditory representations. In *International Conference on Artificial Intelligence and Statistics*, pages 2530–2538. PMLR, 2021.

[AW18] Luis Aguiar and Joel Waldfogel. Platforms, promotion, and product discovery: Evidence from spotify playlists. Technical report, National Bureau of Economic Research, 2018.

[AWZ+21] Junyi Ao, Rui Wang, Long Zhou, Chengyi Wang, Shuo Ren, Yu Wu, Shujie Liu, Tom Ko, Qing Li, Yu Zhang, et al. Speecht5: Unified-modal encoder-decoder pre-training for spoken language processing. *arXiv preprint arXiv:2110.07205*, 2021.

[AYC+23] Armen Aghajanyan, Lili Yu, Alexis Conneau, Wei-Ning Hsu, Karen Hambardzumyan, Susan Zhang, Stephen Roller, Naman Goyal, Omer Levy, and Luke Zettlemoyer. Scaling laws for generative mixed-modal language models. In *International Conference on Machine Learning*, pages 265–279. PMLR, 2023.

[Bae] Gawon Bae. South korea has used ai to bring a dead superstar's voice back to the stage, but ethical concerns abound.

[BAK23] Lorenzo Betti, Carlo Abrate, and Andreas Kaltenbrunner. Large scale analysis of gender bias and sexism in song lyrics. *EPJ Data Science*, 12(1):10, 2023.

[BAM19] Alexei Baevski, Michael Auli, and Abdelrahman Mohamed. Effectiveness of self-supervised pre-training for speech recognition. *arXiv preprint arXiv:1911.03912*, 2019.

[BAM20] Alexei Baevski, Michael Auli, and Abdelrahman Mohamed. Effectiveness of self-supervised pre-training for speech recognition, May 2020.

[Bar07] Karen Barad. *Meeting the universe halfway: Quantum physics and the entanglement of matter and meaning.* Duke University Press, 2007.

[BB18] Georgina Barton and Georgina Barton. The relationship between music, culture, and society: meaning in music. *Music*

*Learning and Teaching in Culturally and Socially Diverse Contexts: Implications for Classroom Practice*, pages 23–41, 2018.

[BBC+23] Anthony Brohan, Noah Brown, Justice Carbajal, Yevgen Chebotar, Xi Chen, Krzysztof Choromanski, Tianli Ding, Danny Driess, Avinava Dubey, Chelsea Finn, et al. Rt-2: Vision-language-action models transfer web knowledge to robotic control. *arXiv preprint arXiv:2307.15818*, 2023.

[BCA18] Scott Beveridge, Estefanía Cano, and Kat Agres. Rhythmic entrainment for hand rehabilitation using the leap motion controller. In *19th International Society for Music Information Retrieval Conference*, 2018.

[BDA+05] Juan Pablo Bello, Laurent Daudet, Samer Abdallah, Chris Duxbury, Mike Davies, and Mark B Sandler. A tutorial on onset detection in music signals. *IEEE Transactions on speech and audio processing*, 13(5):1035–1047, 2005.

[BDDE19] Emmanouil Benetos, Simon Dixon, Zhiyao Duan, and Sebastian Ewert. Automatic music transcription: An overview. *IEEE Signal Processing Magazine*, 36(1):20–30, 2019.

[Ber24] Adam Eric Berkowitz. Artificial Intelligence and Musicking. *Music Perception: An Interdisciplinary Journal*, 41(5):393–412, June 2024.

[BGMMS21] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? parrot-emoji. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pages 610–623, 2021.

[BHA+21] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.

[BHLP16] Adriano Baratè, Goffredo Haus, Luca Andrea Ludovico, and Giorgio Presti. Advances and perspectives in web technologies for music representation. *DigitCult-Scientific Journal on Digital Cultures*, 1(2):1–18, 2016.

[BHX+22] Alexei Baevski, Wei-Ning Hsu, Qiantong Xu, Arun Babu, Jiatao Gu, and Michael Auli. data2vec: A General Framework for Self-supervised Learning in Speech, Vision and Language. In *Proceedings of Machine Learning Research*, Baltimore, MD, USA, 2022.

[Bie20] Elettra Bietti. From ethics washing to ethics bashing: a view on tech ethics from within moral philosophy. In *Proceedings of the 2020 conference on fairness, accountability, and transparency*, pages 210–219, 2020.

[BIS+23] Randall Balestriero, Mark Ibrahim, Vlad Sobal, Ari Morcos, Shashank Shekhar, Tom Goldstein, Florian Bordes, Adrien Bardes, Gregoire Mialon, Yuandong Tian, et al. A cookbook of self-supervised learning. *arXiv preprint arXiv:2304.12210*, 2023.

[BKKB24] Aydar Bulatov, Yuri Kuratov, Yermek Kapushev, and Mikhail S. Burtsev. Scaling transformer to 1m tokens and beyond with rmt, 2024.

[BKS+16] Sebastian Böck, Filip Korzeniowski, Jan Schlüter, Florian Krebs, and Gerhard Widmer. Madmom: A new python audio and music signal processing library. In *Proc. ACM MM*, 2016.

[BKWW18] Gino Brunner, Andres Konrad, Yuyi Wang, and Roger Wattenhofer. Midi-Vae: Modeling Dynamics and Instrumentation of Music with Applications to Style Transfer. In *Proceedings of the International Society for Music Informaiton Retrieval (ISMIR) Conference*, Paris, France, 2018.

[BL05] Satanjeev Banerjee and Alon Lavie. Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, pages 65–72, 2005.

[BMEWL11] Thierry Bertin-Mahieux, Daniel PW Ellis, Brian Whitman, and Paul Lamere. The million song dataset. In *Proceedings of the 12th International Conference on Music Information Retrieval (ISMIR 2011)*, 2011.

[BMR+20] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901, 2020.

[BMV+23] Zalán Borsos, Raphaël Marinier, Damien Vincent, Eugene Kharitonov, Olivier Pietquin, Matt Sharifi, Dominik Roblek, Olivier Teboul, David Grangier, Marco Tagliasacchi, and Neil Zeghidour. Audiolm: A language modeling approach to audio generation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 31:2523–2533, 2023.

[Bor10] Georgina Born. For a Relational Musicology: Music and Interdisciplinarity, Beyond the Practice Turn: The 2007 Dent Medal Address. *Journal of the Royal Musical Association*, 135(2):205–243, 2010.

[Bor20] Georgina Born. Diversifying MIR: Knowledge and Real-World Challenges, and New Interdisciplinary Futures. *Transactions of the International Society for Music Information Retrieval*, 3(1):193–204, October 2020.

[BPC20] Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer, 2020.

[BPH22] Alan Baade, Puyuan Peng, and David Harwath. MAE-AST: Masked Autoencoding Audio Spectrogram Transformer. In *Interspeech 2022*, pages 2438–2442, Incheon, Korea, September 2022. ISCA.

[Bra16] Rosi Braidotti. Posthuman Critical Theory. In Debashish Banerji and Makarand R. Paranjape, editors, *Critical Posthumanism and Planetary Futures*. Springer, 2016.

[BRGL+23] Roser Batlle-Roca, Emila Gómez, WeiHsiang Liao, Xavier Serra, and Yuki Mitsufuji. Transparency in music-generative ai: A systematic literature review. *Preprint https://doi.org/10.21203/rs.3.rs-3708077/v1*, 2023.

[BS16] Solon Barocas and Andrew D Selbst. Big data's disparate impact. *Calif. L. Rev.*, 104:671, 2016.

[BSV+23] Zalán Borsos, Matt Sharifi, Damien Vincent, Eugene Kharitonov, Neil Zeghidour, and Marco Tagliasacchi. Soundstorm: Efficient parallel audio generation, 2023.

[BWT+19] Dmitry Bogdanov, Minz Won, Philip Tovstogan, Alastair Porter, and Xavier Serra. The mtg-jamendo dataset for automatic music tagging. In *Machine Learning for Music Discovery Workshop, International Conference on Machine Learning (ICML 2019)*, Long Beach, CA, United States, 2019.

[BZMA20a] Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A Framework for Self-Supervised Learning of Speech Representations. In *Advances in Neural Information Processing Systems*, Online, 2020.

[BZMA20b] Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460, 2020.

[C2P] C2PA. Coalition for content provenance and authenticity.

[CB20] Patricia Hill Collins and Sirma Bilge. *Intersectionality, 2nd Edition*. John Wiley & Sons, 2020.

[CBBG19] Pritish Chandna, Merlijn Blaauw, Jordi Bonada, and Emilia Gómez. Wgansing: A multi-voice singing voice synthesizer based on the wasserstein-gan. In *2019 27th European signal processing conference (EUSIPCO)*, pages 1–5. IEEE, 2019.

[CBKH05] Pedro Cano, Eloi Batlle, Ton Kalker, and Jaap Haitsma. A review of audio fingerprinting. *Journal of VLSI signal processing systems for signal, image and video technology*, 41:271–284, 2005.

[CCC+21] Yi-Hui Chou, I Chen, Chin-Jui Chang, Joann Ching, Yi-Hsuan Yang, et al. Midibert-piano: large-scale pre-training for symbolic music understanding. *arXiv preprint arXiv:2107.05223*, 2021.

[CCC+24] Yi-Hui Chou, I.-Chun Chen, Chin-Jui Chang, Joann Ching, and Yi-Hsuan Yang. BERT-like Pre-training for Symbolic Piano Music Classification Tasks, April 2024.

[CCG+23] Qian Chen, Yunfei Chu, Zhifu Gao, Zerui Li, Kai Hu, Xiaohuan Zhou, Jin Xu, Ziyang Ma, Wen Wang, Siqi Zheng, et al. Lauragpt: Listen, attend, understand, and regenerate audio with gpt. *arXiv preprint arXiv:2310.04673*, 2023.

[CCK+22] Kin Wai Cheuk, Keunwoo Choi, Qiuqiang Kong, Bochen Li, Minz Won, Amy Hung, Ju-Chiang Wang, and Dorien Herremans. Jointist: Joint learning for multi-instrument transcription and its applications. *arXiv preprint arXiv:2206.10805*, 2022.

[CDL21]     Rodrigo Castellon, Chris Donahue, and Percy Liang. Codified audio language modeling learns useful representations for music information retrieval. *International Society for Music Information Retrieval*, 2021.

[CDRS20]    Rowland Chen, Roger B. Dannenberg, Bhiksha Raj, and Rita Singh. Artificial creative intelligence: Breaking the imitation barrier. In *International Conference on Innovative Computing and Cloud Computing*, 2020.

[CDZ+22]    Ke Chen, Xingjian Du, Bilei Zhu, Zejun Ma, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. Hts-at: A hierarchical token-semantic audio transformer for sound classification and detection. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 646–650. IEEE, 2022.

[CE21]      Antoine Caillon and Philippe Esling. RAVE: A variational autoencoder for fast and high-quality neural audio synthesis, December 2021. arXiv:2111.05011 [cs, eess].

[CFGH20]    Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020.

[CFL+18]    Estefania Cano, Derry FitzGerald, Antoine Liutkus, Mark D Plumbley, and Fabian-Robert Stöter. Musical source separation: An introduction. *IEEE Signal Processing Magazine*, 36(1):31–40, 2018.

[CFS16]     Keunwoo Choi, George Fazekas, and Mark Sandler. Automatic tagging using deep convolutional neural networks. *arXiv preprint arXiv:1606.00298*, 2016.

[CFSC17]    Keunwoo Choi, George Fazekas, Mark Sandler, and Kyunghyun Cho. Transfer learning for music classification and regression tasks. In *The 18th International Society of Music Information Retrieval (ISMIR) Conference 2017, Suzhou, China*. International Society of Music Information Retrieval, 2017.

[CG18]      Florian Colombo and Wulfram Gerstner. Bachprop: Learning to compose music in multiple styles. *CoRR*, abs/1802.05162, 2018.

[CGRS19]    Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers, 2019.

[CH18]      Ching Hua Chuan and Dorien Herremans. Modeling temporal tonal relations in polyphonic music through deep networks with a novel image-based representation. In *32nd AAAI Conference on Artificial Intelligence, AAAI 2018*, pages 2159–2166, 2018.

[CH19]      Raja Chatila and John C Havens. The ieee global initiative on ethics of autonomous and intelligent systems. *Robotics and well-being*, pages 11–16, 2019.

[CHL+22]    Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022.

[CHL+24]    Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53, 2024.

[CJCC22]    Jeong Choi, Seongwon Jang, Hyunsouk Cho, and Sehee Chung. Towards proper contrastive self-supervised learning strategies for music audio representation. In *2022 IEEE International Conference on Multimedia and Expo (ICME)*, pages 1–6. IEEE, 2022.

[CKG+23a]   Jade Copet, Felix Kreuk, Itai Gat, Tal Remez, David Kant, Gabriel Synnaeve, Yossi Adi, and Alexandre Défossez. Simple and controllable music generation. *arXiv preprint arXiv:2306.05284*, 2023.

[CKG+23b]   Jade Copet, Felix Kreuk, Itai Gat, Tal Remez, David Kant, Gabriel Synnaeve, Yossi Adi, and Alexandre Défossez. Simple and controllable music generation. *arXiv preprint arXiv:2306.05284*, 2023.

[CKG+24]    Jade Copet, Felix Kreuk, Itai Gat, Tal Remez, David Kant, Gabriel Synnaeve, Yossi Adi, and Alexandre Défossez. Simple and controllable music generation. *Advances in Neural Information Processing Systems*, 36, 2024.

[CKM+19]    Wenqin Chen, Jessica Keast, Jordan Moody, Corinne Moriarty, Felicia Villalobos, Virtue Winter, Xueqi Zhang, Xuanqi Lyu, Elizabeth Freeman, Jessie Wang, Sherry Cai, and Katherine Kinnaird. Data Usage in MIR: History & Future Recommendations. In *Proceedings of the 20th International Society for Music Information Retrieval Conference*, pages 25–32, Delft, The Netherlands, November 2019. ISMIR.

[CKNH20]    Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.

[Cla21]     Martin Clancy. *Reflections On The Financial And Ethical Implications Of Music Generated By Artificial Intelligence*. PhD Thesis, Trinity College, Dublin, 2021.

[CLG11]     Wei Cai, Qiang Li, and Xin Guan. Automatic singer identification based on auditory features. In *Seventh International Conference on Natural Computation, ICNC 2011, Shanghai, China, 26-28 July, 2011*, pages 1624–1628. IEEE, 2011.

[CLL+23]    Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot impressing gpt-4 with 90%* chatgpt quality. *See https://vicuna. lmsys. org (accessed 14 April 2023)*, 2(3):6, 2023.

[CLM+24]    Wenxi Chen, Yuzhe Liang, Ziyang Ma, Zhisheng Zheng, and Xie Chen. Eat: Self-supervised pre-training with efficient audio transformer. *arXiv preprint arXiv:2401.03497*, 2024.

[CLPN19]    Jeong Choi, Jongpil Lee, Jiyoung Park, and Juhan Nam. Zero-shot learning and knowledge transfer in music classification and tagging. *arXiv preprint arXiv:1906.08615*, 2019.

[CM24]      Santiago Cuervo and Ricard Marxer. Scaling properties of speech language models. *arXiv preprint arXiv:2404.00685*, 2024.

[CMJG17]    Pritish Chandna, Marius Miron, Jordi Janer, and Emilia Gómez. Monoaural audio source separation using deep convolutional neural networks. In *Latent Variable Analysis and Signal Separation: 13th International Conference, LVA/ICA 2017, Grenoble, France, February 21-23, 2017, Proceedings 13*, pages 258–266. Springer, 2017.

[CND+22]    Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. arxiv 2022. *arXiv preprint arXiv:2204.02311*, 10, 2022.

[CNJ+24]    Sanjoy Chowdhury, Sayan Nag, KJ Joseph, Balaji Vasan Srinivasan, and Dinesh Manocha. Melfusion: Synthesizing music from image and language cues using diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26826–26835, 2024.

[COR30]     Harvy Clyde Carver, AL O'TOOLE, and TE RAIFORD. *The annals of mathematical statistics*. Edwards Bros., 1930.

[CPH+21]    Kyoyun Choi, Jonggwon Park, Wan Heo, Sungwook Jeon, and Jonghun Park. Chord conditioned melody generation with transformer based decoders. *IEEE Access*, 9:42071–42080, 2021.

[CPM+24]    Ruben Ciranni, Emilian Postolache, Giorgio Mariani, Michele Mancusi, Luca Cosmo, and Emanuele Rodolà. Cocola: Coherence-oriented contrastive learning of musical audio representations. *arXiv preprint arXiv:2404.16969*, 2024.

[CPZ+20]    Jiemei Chen, Fan Pan, Ping Zhong, Tiantian He, Leiyu Qi, Jingzhe Lu, Peiyu He, and Yun Zheng. An Automatic Method to Develop Music With Music Segment and Long Short Term Memory for Tinnitus Music Therapy. *IEEE Access*, 8:141860–141871, 2020.

[CQZ+22a]   Chung-Cheng Chiu, James Qin, Yu Zhang, Jiahui Yu, and Yonghui Wu. Self-supervised learning with random-projection quantizer for speech recognition. In *Proceedings of the International Conference on Machine Learning*, pages 3915–3924. PMLR, June 2022. ISSN: 2640-3498.

[CQZ+22b]   Chung-Cheng Chiu, James Qin, Yu Zhang, Jiahui Yu, and Yonghui Wu. Self-supervised learning with random-projection quantizer for speech recognition. In *International Conference on Machine Learning*, pages 3915–3924. PMLR, 2022.

[Cra24]     Kate Crawford. Generative AI's environmental costs are soaring — and mostly secret. *Nature*, 626(8000):693–693, February 2024.

[CS17]      Ralf Gunter Correa Carvalho and Paris Smaragdis. Towards end-to-end polyphonic music transcription: Transforming music audio directly to a score. In *2017 IEEE Workshop on*

*Applications of Signal Processing to Audio and Acoustics (WASPAA)*, pages 151–155. IEEE, 2017.

[CSS19]   Marcelo Caetano, Charalampos Saitis, and Kai Siedenburg. Audio content descriptors of timbre. In Kai Siedenburg, Charalampos Saitis, Stephen McAdams, Arthur N Popper, and Richard R Fay, editors, *Timbre: Acoustics, Perception, and Cognition*, pages 297–333. Springer, Cham, 2019.

[Csu17]   Gabriela Csurka. Domain adaptation for visual applications: A comprehensive survey. *arXiv preprint arXiv:1702.05374*, 2017.

[CSU+23]   Kin Wai Cheuk, Ryosuke Sawata, Toshimitsu Uesaka, Naoki Murata, Naoya Takahashi, Shusuke Takahashi, Dorien Herremans, and Yuki Mitsufuji. DiffRoll: DIFFUSION-BASED GENERATIVE MUSIC TRANSCRIPTION WITH UNSUPERVISED PRETRAINING CAPABILITY. In *Proceedings of the International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Rhodes Island, Greece, 2023.

[CTL+20]   Jiawei Chen, Xu Tan, Jian Luan, Tao Qin, and Tie-Yan Liu. Hifisinger: Towards high-fidelity neural singing voice synthesis. *arXiv preprint arXiv:2009.01776*, 2020.

[CWBKD20]   Ke Chen, Cheng I. Wang, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. Music Sketchnet: Controllable Music Generation Via Factorized Representations of Pitch and Rhythm. In *Proceedings of the 21st International Society for Music Information Retrieval Conference, ISMIR 2020*, pages 77–84, 2020.

[CWC+22a]   Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, Jian Wu, Long Zhou, Shuo Ren, Yanmin Qian, Yao Qian, Jian Wu, Michael Zeng, Xiangzhan Yu, and Furu Wei. Wavlm: Large-scale self-supervised pre-training for full stack speech processing. *IEEE J. Sel. Top. Signal Process.*, 16(6):1505–1518, 2022.

[CWC+22b]   Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, Jian Wu, Long Zhou, Shuo Ren, Yanmin Qian, Yao Qian, Jian Wu, Michael Zeng, Xiangzhan Yu, and Furu Wei. WavLM: Large-Scale Self-Supervised Pre-Training for Full Stack Speech Processing. *IEEE Journal of Selected Topics in Signal Processing*, 16(6):1505–1518, October 2022.

[CWC+24]   Xiaowei Chi, Yatian Wang, Aosong Cheng, Pengjun Fang, Zeyue Tian, Yingqing He, Zhaoyang Liu, Xingqun Qi, Jiahao Pan, Rongyu Zhang, et al. Mmtrail: A multimodal trailer video dataset with language and music descriptions. *arXiv preprint arXiv:2407.20962*, 2024.

[CWL+23a]   Ke Chen, Yusong Wu, Haohe Liu, Marianna Nezhurina, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. Musicldm: Enhancing novelty in text-to-music generation using beat-synchronous mixup strategies. *arXiv preprint arXiv:2308.01546*, 2023.

[CWL+23b]   Ke Chen, Yusong Wu, Haohe Liu, Marianna Nezhurina, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. Musicldm: Enhancing novelty in text-to-music generation using beat-synchronous mixup strategies. *CoRR*, abs/2308.01546, 2023.

[CWW+22]   Sanyuan Chen, Yu Wu, Chengyi Wang, Shujie Liu, Daniel Tompkins, Zhuo Chen, and Furu Wei. Beats: Audio pre-training with acoustic tokenizers. *arXiv preprint arXiv:2212.09058*, 2022.

[CWZZ23]   Dading Chong, Helin Wang, Peilin Zhou, and Qingcheng Zeng. Masked Spectrogram Prediction for Self-Supervised Audio Pre-Training. In *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5, Rhodes Island, Greece, June 2023. IEEE.

[CXH21]   Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9640–9649, 2021.

[CXY+24]   Yunfei Chu, Jin Xu, Qian Yang, Haojie Wei, Xipin Wei, Zhifang Guo, Yichong Leng, Yuanjun Lv, Jinzheng He, Junyang Lin, et al. Qwen2-audio technical report. *arXiv preprint arXiv:2407.10759*, 2024.

[CXZ+23]   Yunfei Chu, Jin Xu, Xiaohuan Zhou, Qian Yang, Shiliang Zhang, Zhijie Yan, Chang Zhou, and Jingren Zhou. Qwen-audio: Advancing universal audio understanding via

unified large-scale audio-language models. *arXiv preprint arXiv:2311.07919*, 2023.

[CYLK23]   Hyeong-Seok Choi, Jinhyeok Yang, Juheon Lee, and Hyeongju Kim. NANSY++: unified voice synthesis with neural analysis and synthesis. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.

[CZH+21a]   Yu-An Chung, Yu Zhang, Wei Han, Chung-Cheng Chiu, James Qin, Ruoming Pang, and Yonghui Wu. W2v-bert: Combining contrastive learning and masked language modeling for self-supervised speech pre-training. In *2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 244–250. IEEE, 2021.

[CZH+21b]   Yu-An Chung, Yu Zhang, Wei Han, Chung-Cheng Chiu, James Qin, Ruoming Pang, and Yonghui Wu. w2v-BERT: Combining Contrastive Learning and Masked Language Modeling for Self-Supervised Speech Pre-Training. In *2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 244–250, December 2021.

[CZJ+22]   Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11315–11325, 2022.

[CZJP20]   Jorge Calvo-Zaragoza, Jan Hajič Jr, and Alexander Pacha. Understanding optical music recognition. *ACM Computing Surveys (CSUR)*, 53(4):1–35, 2020.

[CZL+23]   Xuxin Cheng, Zhihong Zhu, Hongxiang Li, Yaowei Li, and Yuexian Zou. SSVMR: saliency-based self-training for video-music retrieval. In *IEEE International Conference on Acoustics, Speech and Signal Processing ICASSP 2023, Rhodes Island, Greece, June 4-10, 2023*, pages 1–5. IEEE, 2023.

[CZZ+21a]   Nanxin Chen, Yu Zhang, Heiga Zen, Ron J Weiss, Mohammad Norouzi, and William Chan. Wavegrad: Estimating gradients for waveform generation. In *International Conference on Learning Representations*, 2021.

[CZZ+21b]   Nanxin Chen, Yu Zhang, Heiga Zen, Ron J. Weiss, Mohammad Norouzi, Najim Dehak, and William Chan. WaveGrad 2: Iterative Refinement for Text-to-Speech Synthesis. In *Interspeech*, pages 3765–3769, 2021.

[DAD20]   Emir Demirel, Sven Ahlbäck, and Simon Dixon. Automatic lyrics transcription using dilated convolutional neural networks with self-attention. In *2020 International Joint Conference on Neural Networks, IJCNN 2020, Glasgow, United Kingdom, July 19-24, 2020*, pages 1–8. IEEE, 2020.

[DAD21]   Emir Demirel, Sven Ahlbäck, and Simon Dixon. Mstre-net: Multistreaming acoustic modeling for automatic lyrics transcription. In *Proceedings of the 22nd International Society for Music Information Retrieval Conference, ISMIR 2021, Online, November 7-12, 2021*, pages 151–158, 2021.

[DB19]   Gerardo Roa Dabike and Jon Barker. Automatic lyric transcription from karaoke vocal tracks: Resources and a baseline system. In *20th Annual Conference of the International Speech Communication Association, Interspeech 2019, Graz, Austria, September 15-19, 2019*, pages 579–583. ISCA, 2019.

[DB22]   Karlijn Dinnissen and Christine Bauer. Fairness in music recommender systems: A stakeholder-centered mini review. *Frontiers in big Data*, 5:913608, 2022.

[DBK+20]   Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

[DBVB16]   Michaël Defferrard, Kirell Benzi, Pierre Vandergheynst, and Xavier Bresson. Fma: A dataset for music analysis. *arXiv preprint arXiv:1612.01840*, 2016.

[DCD+23]   Hao-Wen Dong, Ke Chen, Shlomo Dubnov, Julian McAuley, and Taylor Berg-Kirkpatrick. Multitrack music transformer. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.

[DCLN23a]   Seungheon Doh, Keunwoo Choi, Jongpil Lee, and Juhan Nam. Lp-musiccaps: Llm-based pseudo music captioning. In Augusto Sarti, Fabio Antonacci, Mark Sandler, Paolo Bestagini, Simon Dixon, Beici Liang, Gaël Richard, and Johan Pauwels, editors, *Proceedings of the 24th International Society for Music Information Retrieval Conference, ISMIR*

*2023, Milan, Italy, November 5-9, 2023*, pages 409–416, 2023.

[DCLN23b] SeungHeon Doh, Keunwoo Choi, Jongpil Lee, and Juhan Nam. Lp-musiccaps: Llm-based pseudo music captioning. *ISMIR*, 2023.

[DCLT18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

[DCR+23] Chris Donahue, Antoine Caillon, Adam Roberts, Ethan Manilow, Philippe Esling, Andrea Agostinelli, Mauro Verzetti, Ian Simon, Olivier Pietquin, Neil Zeghidour, et al. Singsong: Generating musical accompaniments from singing. *arXiv preprint arXiv:2301.12662*, 2023.

[DCSA22] Alexandre Défossez, Jade Copet, Gabriel Synnaeve, and Yossi Adi. High fidelity neural audio compression. *arXiv preprint arXiv:2210.13438*, 2022.

[DDGK18] Roel Dobbe, Sarah Dean, Thomas Gilbert, and Nitin Kohli. A broader view on bias in automated decision-making: Reflecting on epistemology and dynamics. *arXiv preprint arXiv:1807.00553*, 2018.

[DDHB19] Yann Dubois, Gautier Dagan, Dieuwke Hupkes, and Elia Bruni. Location attention for extrapolation to longer sequences. *arXiv preprint arXiv:1911.03872*, 2019.

[DDP09] Matthew EP Davies, Norberto Degara, and Mark D Plumbley. Evaluation methods for musical audio beat tracking algorithms. *Tech. Rep., Queen Mary University of London, Centre for Digital Music,*, 2009.

[DESW23] Soham Deshmukh, Benjamin Elizalde, Rita Singh, and Huaming Wang. Pengi: An audio language model for audio tasks. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023.

[Dev19] Kyle Devine. *Decomposed: the political ecology of music*. The MIT Press, Cambridge, Massachusetts, 2019.

[DGCY22] Chenpeng Du, Yiwei Guo, Xie Chen, and Kai Yu. VQTTS: high-fidelity text-to-speech synthesis with self-supervised VQ acoustic feature. In *23rd Annual Conference of the International Speech Communication Association, Interspeech 2022, Incheon, Korea, September 18-22, 2022*, pages 1596–1600, 2022.

[DHYY18] Hao-Wen Dong, Wen-Yi Hsiao, Li-Chia Yang, and Yi-Hsuan Yang. Musegan: Multi-track sequential generative adversarial networks for symbolic music generation and accompaniment. In Sheila A. McIlraith and Kilian Q. Weinberger, editors, *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana, USA, February 2-7, 2018*, pages 34–41. AAAI Press, 2018.

[DJB+23] Yashar Deldjoo, Dietmar Jannach, Alejandro Bellogin, Alessandro Difonzo, and Dario Zanzonelli. Fairness in recommender systems: research landscape and future directions. *User Modeling and User-Adapted Interaction*, 34(1):59–108, 2023.

[DJGD21] Shuqi Dai, Zeyu Jin, Celso Gomes, and Roger B Dannenberg. Controllable deep melody generation via hierarchical music structure representation. *arXiv preprint arXiv:2109.00663*, 2021.

[DJL+21] Shangzhe Di, Zeren Jiang, Si Liu, Zhaokai Wang, Leyan Zhu, Zexin He, Hongming Liu, and Shuicheng Yan. Video background music generation with controllable music transformer. In Heng Tao Shen, Yueting Zhuang, John R. Smith, Yang Yang, Pablo César, Florian Metze, and Balakrishnan Prabhakaran, editors, *MM '21: ACM Multimedia Conference, Virtual Event, China, October 20 - 24, 2021*, pages 2037–2045. ACM, 2021.

[DJP+20a] Prafulla Dhariwal, Heewoo Jun, Christine Payne, Jong Wook Kim, Alec Radford, and Ilya Sutskever. Jukebox: A generative model for music. *arXiv preprint arXiv:2005.00341*, 2020.

[DJP+20b] Prafulla Dhariwal, Heewoo Jun, Christine Payne, Jong Wook Kim, Alec Radford, and Ilya Sutskever. Jukebox: A generative model for music. *arXiv preprint arXiv:2005.00341*, 2020.

[DL82] DC Dowson and BV Landau. The fréchet distance between multivariate normal distributions. *Journal of multivariate analysis*, 1982.

[DLD+22] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. A survey on in-context learning. *arXiv preprint arXiv:2301.00234*, 2022.

[DLD+24a] Shuangrui Ding, Zihan Liu, Xiaoyi Dong, Pan Zhang, Rui Qian, Conghui He, Dahua Lin, and Jiaqi Wang. Songcomposer: A large language model for lyric and melody composition in song generation. *CoRR*, abs/2402.17645, 2024.

[DLD+24b] Shuangrui Ding, Zihan Liu, Xiaoyi Dong, Pan Zhang, Rui Qian, Conghui He, Dahua Lin, and Jiaqi Wang. Songcomposer: A large language model for lyric and melody composition in song generation. *arXiv preprint arXiv:2402.17645*, 2024.

[dlTPVB24] Pablo González de la Torre, Marta Pérez-Verdugo, and Xabier E Barandiaran. Attention is all they need: Cognitive science and the (techno) political economy of attention in humans and machines. *arXiv preprint arXiv:2405.06478*, 2024.

[DMBHM24] Dorian Desblancs, Gabriel Meseguer-Brocal, Romain Hennequin, and Manuel Moussallam. From real to cloned singer identification, 2024.

[DML+24] Zihao Deng, Yinghao Ma, Yudong Liu, Rongchen Guo, Ge Zhang, Wenhu Chen, Wenhao Huang, and Emmanouil Benetos. Musilingo: Bridging music and text with pre-trained language models for music captioning and query response. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics*, 2024.

[DMM18] Chris Donahue, Huanru Henry Mao, and Julian McAuley. The nes music database: A multi-instrumental dataset with expressive performance attributes. *ISMIR*, 2018.

[DN21] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.

[Dow03] J Stephen Downie. Music information retrieval. *Annual review of information science and technology*, 37(1):295–340, 2003.

[dPS+22] Martina de Witte, Ana da Silva Pinho, Geert-Jan Stams, Xavier Moonen, Arjan E.R. Bos, and Susan van Hooren. Music therapy for stress reduction: A systematic review and meta-analysis. *Health Psychology Review*, 16(1):134–159, January 2022.

[DSCS23] Bleiz Macsen Del Sette, Dawn Carnes, and Charalampos Saitis. Sound of Care: Towards a Co-Operative AI Digital Pain Companion to Support People with Chronic Primary Pain. In *Computer Supported Cooperative Work and Social Computing*, pages 283–288, Minneapolis MN USA, October 2023. ACM.

[dSdLT+21] Maíra Araújo de Santana, Clarisse Lins de Lima, Arianne Sarmento Torcate, Flávio Secco Fonseca, and Wellington Pinheiro dos Santos. Affective computing in the context of music therapy: A systematic review. *Research, Society and Development*, 10(15):e392101522844–e392101522844, November 2021.

[DSE23] Simon Durand, Daniel Stoller, and Sebastian Ewert. Contrastive learning-based audio to lyrics alignment for multiple languages. In *IEEE International Conference on Acoustics, Speech and Signal Processing ICASSP 2023, Rhodes Island, Greece, June 4-10, 2023*, pages 1–5. IEEE, 2023.

[DUBB19] Alexandre Défossez, Nicolas Usunier, Léon Bottou, and Francis Bach. Music source separation in the waveform domain. *arXiv preprint arXiv:1911.13254*, 2019.

[DvdOS18] Sander Dieleman, Aaron van den Oord, and Karen Simonyan. The challenge of realistic music generation: modelling raw audio at scale. *Advances in neural information processing systems*, 31, 2018.

[DVE+23] Grant Davidson, Mark Vinton, Per Ekstrand, Cong Zhou, Lars Villemoes, and Lie Lu. High quality audio coding with mdctnet. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.

[DVK17] Finale Doshi-Velez and Been Kim. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608*, 2017.

[DVVD20] Cedric De Boom, Stephanie Van Laere, Tim Verbelen, and Bart Dhoedt. Rhythm, Chord and Melody Generation for Lead Sheets Using Recurrent Neural Networks. *Communications*

*in Computer and Information Science*, 1168 CCIS:454–461, 2020.

[DWCN23] SeungHeon Doh, Minz Won, Keunwoo Choi, and Juhan Nam. Toward universal text-to-music retrieval. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.

[DYL+24] Xingjian Du, Zhesong Yu, Jiaju Lin, Bilei Zhu, and Qiuqiang Kong. Joint music and language attention models for zero-shot music tagging. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1126–1130. IEEE, 2024.

[DYW+24] Heinrich Dinkel, Zhiyong Yan, Yongqing Wang, Junbo Zhang, Yujun Wang, and Bin Wang. Scaling up masked audio encoder learning for general audio classification, June 2024. arXiv:2406.06992 [cs, eess].

[DYY+19] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G. Carbonell, Quoc Viet Le, and Ruslan Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. In *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019*, pages 2978–2988, 2019.

[DYY+24] Qixin Deng, Qikai Yang, Ruibin Yuan, Yipeng Huang, Yi Wang, Xubo Liu, Zeyue Tian, Jiahao Pan, Ge Zhang, Hanfeng Lin, et al. Composerx: Multi-agent symbolic music composition with llms. In *The 25th International Society of Music Information Retrieval (ISMIR) Conference 20224, USA*. International Society of Music Information Retrieval, 2024.

[ECT+24] Zach Evans, CJ Carr, Josiah Taylor, Scott H Hawley, and Jordi Pons. Fast timing-conditioned latent audio diffusion. *arXiv preprint arXiv:2402.04825*, 2024.

[EDAIW23] Benjamin Elizalde, Soham Deshmukh, Mahmoud Al Ismail, and Huaming Wang. CLAP learning audio concepts from natural language supervision. In *Proc. ICASSP*, 2023.

[EDCB20] Avriel Epps-Darling, Henriette Cramer, and Romain Takeo Bouyer. Artist gender representation in music streaming. In *ISMIR*, pages 248–254, 2020.

[EFJ+19] Maria Eriksson, Rasmus Fleischer, Anna Johansson, Pelle Snickars, and Patrick Vonderau. *Spotify teardown: Inside the black box of streaming music*. Mit Press, 2019.

[Ell07] Daniel P. W. Ellis. Classifying music audio with timbral and chroma features. In *Proceedings of the 8th International Conference on Music Information Retrieval, ISMIR 2007, Vienna, Austria, September 23-27, 2007*, pages 339–340. Austrian Computer Society, 2007.

[Eng77] George L. Engel. The Need for a New Medical Model: A Challenge for Biomedicine. *Science*, 196(4286):129–136, April 1977.

[EP20] Jeff Ens and Philippe Pasquier. Mmm : Exploring conditional multi-track music generation with the transformer, 2020.

[EP21] Jeffrey Ens and Philippe Pasquier. Building the metamidi dataset: Linking symbolic and audio musical data. In *ISMIR*, pages 182–188, 2021.

[EPC+24] Zach Evans, Julian D Parker, CJ Carr, Zack Zukowski, Josiah Taylor, and Jordi Pons. Long-form music generation with latent diffusion. *arXiv preprint arXiv:2404.10301*, 2024.

[ERO21] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12873–12883, 2021.

[ESFC24] Mehmet Hamza Erol, Arda Senocak, Jiu Feng, and Joon Son Chung. Audio mamba: Bidirectional state space model for audio representation learning. *arXiv preprint arXiv:2406.03344*, 2024.

[FBC+21] Nathan Fradet, Jean-Pierre Briot, Fabien Chhel, Amal El Fallah Seghrouchni, and Nicolas Gutowski. Miditok: a Python Package for Midi File Tokenization. In *International Society for Music Information Retrieval (ISMIR) Late Breaking Demo (LBD)*, 2021.

[FCGG19] Michael Fell, Elena Cabrio, Fabien Gandon, and Alain Giboin. Song lyrics summarization inspired by audio thumbnailing. In *Proceedings of the International Conference on Recent Advances in Natural Language Processing, RANLP 2019, Varna, Bulgaria, September 2-4, 2019*, pages 328–337. INCOMA Ltd., 2019.

[FCZ+11] Nizan Friedman, Vicky Chan, Danny Zondervan, Mark Bachman, and David J Reinkensmeyer. Musicglove: Motivating and quantifying hand movement rehabilitation by using functional grips to play music. In *2011 annual international conference of the IEEE engineering in medicine and biology society*, pages 2359–2363. IEEE, 2011.

[FDVS20] Xavier Favory, Konstantinos Drossos, Tuomas Virtanen, and Xavier Serra. Coala: Co-aligned autoencoders for learning semantically enriched audio representations. *arXiv preprint arXiv:2006.08386*, 2020.

[FFH24] Zhengcong Fei, Mingyuan Fan, and Junshi Huang. A-JEPA: Joint-Embedding Predictive Architecture Can Listen, January 2024. arXiv:2311.15830 [cs, eess].

[FGCB23] Nathan Fradet, Nicolas Gutowski, Fabien Chhel, and Jean-Pierre Briot. Byte pair encoding for symbolic music. *arXiv preprint arXiv:2301.11975*, 2023.

[FH19] Batya Friedman and David F. Hendry. *Value sensitive design: shaping technology with moral imagination*. The MIT Press, Cambridge, MA London, England, 2019.

[FLJ18] Zhe-Cheng Fan, Yen-Lin Lai, and Jyh-Shing Roger Jang. SVSGAN: singing voice separation via generative adversarial network. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2018, Calgary, AB, Canada, April 15-20, 2018*, pages 726–730. IEEE, 2018.

[FLTZ10] Zhouyu Fu, Guojun Lu, Kai Ming Ting, and Dengsheng Zhang. A survey of audio-based music classification and annotation. *IEEE transactions on multimedia*, 13(2):303–319, 2010.

[FLW20] Lucas N. Ferreira, Levi H. S. Lelis, and Jim Whitehead. Computer-generated music for tabletop role-playing games. In Levi Lelis and David Thue, editors, *Proceedings of the Sixteenth AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment, AIIDE 2020, virtual, October 19-23, 2020*, pages 59–65. AAAI Press, 2020.

[FM22a] Seth* Forsgren and Hayk* Martiros. Riffusion - Stable diffusion for real-time music generation, 2022.

[FM22b] Seth* Forsgren and Hayk* Martiros. Riffusion - Stable diffusion for real-time music generation, 2022.

[FMG+19] Jörg C Fachner, Clemens Maidhof, Denise Grocke, Inge Nygaard Pedersen, Gro Trondalen, Gerhard Tucek, and Lars O Bonde. "telling me not to worry. . ." hyperscanning and neural dynamics of emotion processing during guided imagery and music. *Frontiers in Psychology*, 10:1561, 2019.

[FN96] Batya Friedman and Helen Nissenbaum. Bias in computer systems. *ACM Transactions on information systems (TOIS)*, 14(3):330–347, 1996.

[FOM+21] Eduardo Fonseca, Diego Ortego, Kevin McGuinness, Noel E O'Connor, and Xavier Serra. Unsupervised contrastive learning of sound event representations. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 371–375. IEEE, 2021.

[FR23] Robert Flynn and Anton Ragni. How much context does my attention-based asr system need?, 2023.

[Fuj99] Takuya Fujishima. Realtime chord recognition of musical sound: Asystem using common lisp music. In *Proc. ICMC*, 1999.

[FVRA21] Abolfazl Farahani, Sahar Voghoei, Khaled Rasheed, and Hamid R Arabnia. A brief review of domain adaptation. *Advances in Data Science and Information Engineering*, page 877, 2021.

[FWGL23] Ruchao Fan, Yiming Wang, Yashesh Gaur, and Jinyu Li. CTCBERT: Advancing hidden-unit bert with ctc objectives. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.

[GCB+22] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10696–10706, 2022.

[GCG21] Yuan Gong, Yu-An Chung, and James Glass. Ast: Audio spectrogram transformer. *arXiv preprint arXiv:2104.01778*, 2021.

[GD23] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces, 2023.

[GDSB23a] Josh Gardner, Simon Durand, Daniel Stoller, and Rachel M. Bittner. Llark: A multimodal foundation model for music. *CoRR*, abs/2310.07160, 2023.

[GDSB23b]   Joshua P Gardner, Simon Durand, Daniel Stoller, and Rachel M Bittner. Llark: A multimodal instruction-following language model for music. In *Forty-first International Conference on Machine Learning*, 2023.

[GEF+17]   Jort F Gemmeke, Daniel PW Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R Channing Moore, Manoj Plakal, and Marvin Ritter. Audio set: An ontology and human-labeled dataset for audio events. In *2017 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 776–780. IEEE, 2017.

[GFI+23]   Mariana-Iuliana Georgescu, Eduardo Fonseca, Radu Tudor Ionescu, Mario Lucic, Cordelia Schmid, and Anurag Arnab. Audiovisual masked autoencoders. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16144–16154, 2023.

[GGDR22]   Karan Goel, Albert Gu, Chris Donahue, and Christopher Ré. It's raw! audio generation with state-space models. In *International Conference on Machine Learning*, pages 7616–7633. PMLR, 2022.

[GGL23]   Xiaoxue Gao, Chitralekha Gupta, and Haizhou Li. Polyscriber: Integrated fine-tuning of extractor and lyrics transcriber for polyphonic music. *IEEE ACM Trans. Audio Speech Lang. Process.*, 31:1968–1981, 2023.

[GGR22]   Albert Gu, Karan Goel, and Christopher Re. Efficiently modeling long sequences with structured state spaces. In *International Conference on Learning Representations*, 2022.

[GGZ+23]   Yuying Ge, Yixiao Ge, Ziyun Zeng, Xintao Wang, and Ying Shan. Planting a SEED of vision in large language model. *CoRR*, abs/2307.08041, 2023.

[GHC+20a]   Chuang Gan, Deng Huang, Peihao Chen, Joshua B. Tenenbaum, and Antonio Torralba. Foley music: Learning to generate music from videos. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XI*, volume 12356 of *Lecture Notes in Computer Science*, pages 758–775. Springer, 2020.

[GHC+20b]   Chuang Gan, Deng Huang, Peihao Chen, Joshua B Tenenbaum, and Antonio Torralba. Foley music: Learning to generate music from videos. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XI 16*, pages 758–775. Springer, 2020.

[GHMS19]   Emilia Gomez, Andre Holzapfel, Marius Miron, and Bob L Sturm. Fairness, accountability and transparency in music information research (fat-mir). In *Tutorial at the International Society for Music Information Retrieval Conference*, 2019.

[GHY+22]   Songwei Ge, Thomas Hayes, Harry Yang, Xi Yin, Guan Pang, David Jacobs, Jia-Bin Huang, and Devi Parikh. Long video generation with time-agnostic vqgan and time-sensitive transformer. In *European Conference on Computer Vision*, pages 102–118. Springer, 2022.

[GJFA21]   Born Georgina, Morris Jeremy, Diaz Fernando, and Anderson Ashton. Artificial intelligence, music recommendation, and the curation of culture. *TSpace: https://tspace. library. utoronto. ca/handle/1807/129105 (University of Toronto, 2021)*, 2021.

[GKH23]   Zixun Guo, Jaeyong Kang, and Dorien Herremans. A domain-knowledge-inspired music embedding space and a novel attention mechanism for symbolic music modeling. In Brian Williams, Yiling Chen, and Jennifer Neville, editors, *Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023*, pages 5070–5077. AAAI Press, 2023.

[GL84]   Daniel Griffin and Jae Lim. Signal estimation from modified short-time fourier transform. *IEEE Transactions on acoustics, speech, and signal processing*, 32(2):236–243, 1984.

[GLCG22a]   Yuan Gong, Cheng-I. Lai, Yu-An Chung, and James Glass. SSAST: Self-Supervised Audio Spectrogram Transformer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 10699–10709, Online, June 2022. Number: 10.

[GLCG22b]   Yuan Gong, Cheng-I Lai, Yu-An Chung, and James R. Glass. SSAST: self-supervised audio spectrogram transformer. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelveth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022*, pages 10699–10709. AAAI Press, 2022.

[gLKS+22]   Sang gil Lee, Heeseung Kim, Chaehun Shin, Xu Tan, Chang Liu, Qi Meng, Tao Qin, Wei Chen, Sungroh Yoon, and Tie-Yan Liu. Priorgrad: Improving conditional denoising diffusion models with data-dependent adaptive prior. In *International Conference on Learning Representations*, 2022.

[GLL+23]   Yuan Gong, Hongyin Luo, Alexander H Liu, Leonid Karlinsky, and James Glass. Listen, think, and understand. *arXiv preprint arXiv:2305.10790*, 2023.

[GMMP23a]   Deepanway Ghosal, Navonil Majumder, Ambuj Mehrish, and Soujanya Poria. Text-to-audio generation using instruction guided latent diffusion model. In *Proceedings of the 31st ACM International Conference on Multimedia*, page 3590–3598, 2023.

[GMMP23b]   Deepanway Ghosal, Navonil Majumder, Ambuj Mehrish, and Soujanya Poria. Text-to-audio generation using instruction guided latent diffusion model. In *Proc. ACM MM*, 2023.

[GMP20]   Carina Geerlings and Albert Merono-Penuela. Interacting with gpt-2 to generate controlled and believable musical sequences in abc notation. In *Proceedings of the 1st Workshop on NLP for Music and Audio (NLP4MusA)*, pages 49–53, 2020.

[Goo01]   Michael Good. Musicxml for notation and analysis. *The virtual score: representation, retrieval, restoration*, 12(113–124):160, 2001.

[Goo12]   Michael D Good. Musicxml. *Structuring Music through Markup Language: Designs and Architectures: Designs and Architectures*, page 187, 2012.

[Gos08]   William Sealy Gosset. The probable error of a mean. *Biometrika*, 1908.

[GOZ+24]   Xiangming Gu, Longshen Ou, Wei Zeng, Jianan Zhang, Nicholas Wong, and Ye Wang. Automatic lyric transcription and automatic music transcription from multimodal singing. *ACM Trans. Multim. Comput. Commun. Appl.*, 20(7):209:1–209:29, 2024.

[GPAM+20]   Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.

[GQL+19]   Qipeng Guo, Xipeng Qiu, Pengfei Liu, Yunfan Shao, Xiangyang Xue, and Zheng Zhang. Star-transformer. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1315–1325, 2019.

[GRHD22]   Andrey Guzhov, Federico Raue, Jörn Hees, and Andreas Dengel. Audioclip: Extending clip to image, text and audio. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 976–980. IEEE, 2022.

[GRL+23]   Yuan Gong, Andrew Rouditchenko, Alexander H. Liu, David Harwath, Leonid Karlinsky, Hilde Kuehne, and James R. Glass. Contrastive Audio-Visual Masked Autoencoder. In *International Conference on Learning Representations*, Kigali, Rwanda, 2023.

[GS23]   Huijie Guo and Lei Shi. Ultimate negative sampling for contrastive learning. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.

[GSA+20]   Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in neural information processing systems*, 33:21271–21284, 2020.

[GSK+22]   Rui Guo, Ivor Simpson, Chris Kiefer, Thor Magnusson, and Dorien Herremans. MusIAC: An Extensible Generative Framework for Music Infilling Applications with Multi-level Control. In *Artificial Intelligence in Music, Sound, Art and Design: 11th International Conference, EvoMUSART 2022, Held as Part of EvoStar 2022, Madrid, Spain, April 20–22, 2022, Proceedings*, pages 341–356, Berlin, Heidelberg, April 2022. Springer-Verlag.

[GSKP23]   Hugo F Flores Garcia, Prem Seetharaman, Rithesh Kumar, and Bryan Pardo. Vampnet: Music generation via masked acoustic token modeling. In *Ismir 2023 Hybrid Conference*, 2023.

[GSL23]    Xin Gu, Yinghua Shen, and Chaohui Lv. A dual-path cross-modal network for video-music retrieval. *Sensors*, 23(2):805, 2023.

[GYL23]    Xiaoxue Gao, Xianghu Yue, and Haizhou Li. Self-transriber: Few-shot lyrics transcription with self-training. In *IEEE International Conference on Acoustics, Speech and Signal Processing ICASSP 2023, Rhodes Island, Greece, June 4-10, 2023*, pages 1–5. IEEE, 2023.

[GYR+21]   Yu Gu, Xiang Yin, Yonghui Rao, Yuan Wan, Benlai Tang, Yang Zhang, Jitong Chen, Yuxuan Wang, and Zejun Ma. Bytesing: A chinese singing voice synthesis system using duration allocated encoder-decoder acoustic models and wavernn vocoders. In *2021 12th International Symposium on Chinese Spoken Language Processing (ISCSLP)*, pages 1–5. IEEE, 2021.

[GYY19]    Faqian Guan, Chunyan Yu, and Suqiong Yang. A gan model with self-attention mechanism to generate multi-instruments symbolic music. In *2019 International joint conference on neural networks (IJCNN)*, pages 1–6. IEEE, 2019.

[GZJ+24]   Hanzhe Guo, Jiawen Zhang, Yueyao Jiang, Yifei Qi, Simeng Chen, Zhen Chen, Weiran Lin, Junwei Cao, and Shuangs hou Li. EMO-Music: Emotion Recognition Based Music Therapy with Deep Learning on Physiological Signals. In *2024 IEEE First International Conference on Artificial Intelligence for Medicine, Health and Care (AIMHC)*, pages 10–13, February 2024.

[GZM23]    Christos Garoufis, Athanasia Zlatintsi, and Petros Maragos. Multi-Source Contrastive Learning from Musical Audio, May 2023.

[HA15]     Elad Hoffer and Nir Ailon. Deep metric learning using triplet network. In *Similarity-Based Pattern Recognition: Third International Workshop, SIMBAD 2015, Copenhagen, Denmark, October 12-14, 2015. Proceedings 3*, pages 84–92. Springer, 2015.

[HAM+18]   Kotaro Hara, Abigail Adams, Kristy Milland, Saiph Savage, Chris Callison-Burch, and Jeffrey P Bigham. A data-driven analysis of workers' earnings on amazon mechanical turk. In *Proceedings of the 2018 CHI conference on human factors in computing systems*, pages 1–14, 2018.

[Han18]    David J Hand. Aspects of data ethics in a changing world: Where are we now? *Big data*, 6(3):176–190, 2018.

[Har88]    Donna Haraway. Situated knowledges: The science question in feminism and the privilege of partial perspective. *Feminist Studies*, 14(3):575–599, 1988.

[Har95]    Sandra Harding. Strong objectivity: A response to the new objectivity question. *Synthese*, 104:331–349, 1995.

[Har15]    Sandra Harding. *Objectivity and diversity: Another logic of scientific research*. University of Chicago Press, 2015.

[HB24]     Jiawen Huang and Emmanouil Benetos. Towards building an end-to-end multilingual automatic lyrics transcription model. *arXiv preprint arXiv:2406.17618*, 2024.

[HBE22]    Jiawen Huang, Emmanouil Benetos, and Sebastian Ewert. Improving lyrics alignment through joint pitch detection. In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2022, Virtual and Singapore, 23-27 May 2022*, pages 451–455. IEEE, 2022.

[HBL12]    Eric J. Humphrey, Juan Pablo Bello, and Yann LeCun. Moving beyond feature design: Deep architectures and automatic feature learning in music informatics. In *Proc. ISMIR*, pages 403–408, 2012.

[HBM+22]   Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*, 2022.

[HBT+21a]  Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:3451–3460, 2021.

[HBT+21b]  Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. Hubert: Self-supervised speech representation learning by masked prediction of hidden units. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:3451–3460, 2021.

[HBT+21c]  Wei-Ning Hsu, Benjamin Bolte, Yao-Hung Hubert Tsai, Kushal Lakhotia, Ruslan Salakhutdinov, and Abdelrahman Mohamed. HuBERT: Self-Supervised Speech Representation Learning by Masked Prediction of Hidden Units. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:3451–3460, 2021.

[HC23]     Jie Huang and Kevin Chen-Chuan Chang. Towards reasoning in large language models: A survey. In *The 61st Annual Meeting Of The Association For Computational Linguistics*, 2023.

[HCC+22]   Rongjie Huang, Chenye Cui, Feiyang Chen, Yi Ren, Jinglin Liu, Zhou Zhao, Baoxing Huai, and Zhefeng Wang. Singgan: Generative adversarial network for high-fidelity singing voice generation. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 2525–2535, 2022.

[HCD21]    Mojtaba Heydari, Frank Cwitkowitz, and Zhiyao Duan. Beat-Net: CRNN and particle filtering for online joint beat down-beat and meter tracking. In *Proc. ISMIR*, 2021.

[HCE+17]   Shawn Hershey, Sourish Chaudhuri, Daniel PW Ellis, Jort F Gemmeke, Aren Jansen, R Channing Moore, Manoj Plakal, Devin Platt, Rif A Saurous, Bryan Seybold, et al. Cnn architectures for large-scale audio classification. In *2017 ieee international conference on acoustics, speech and signal processing (icassp)*, pages 131–135. IEEE, 2017.

[HCF+20]   Tsung-Han Hsieh, Kai-Hsiang Cheng, Zhe-Cheng Fan, Yu-Ching Yang, and Yi-Hsuan Yang. Addressing the confounds of accompaniments in singer identification. In *2020 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2020, Barcelona, Spain, May 4-8, 2020*, pages 1–5. IEEE, 2020.

[HCH+23]   Zhiqing Hong, Chenye Cui, Rongjie Huang, Lichao Zhang, Jinglin Liu, Jinzheng He, and Zhou Zhao. Unisinger: Unified end-to-end singing voice synthesis with cross-modality information matching. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 7569–7579, 2023.

[HCL06]    Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *2006 IEEE computer society conference on computer vision and pattern recognition (CVPR'06)*, volume 2, pages 1735–1742. IEEE, 2006.

[HCS+22]   Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022.

[HCVKL23]  David Hesmondhalgh, Raquel Campos Valverde, D Kaye, and Zhongwei Li. The impact of algorithmically driven recommendation systems on music consumption and production: A literature review. *UK Centre for Data Ethics and Innovation Reports*, 2023.

[HDS+23a]  Bing Han, Junyu Dai, Xuchen Song, Weituo Hao, Xinyan He, Dong Guo, Jitong Chen, Yuxuan Wang, and Yanmin Qian. Instructme: An instruction guided music edit and remix framework with latent diffusion models. *arXiv preprint arXiv:2308.14360*, 2023.

[HDS+23b]  Bing Han, Junyu Dai, Xuchen Song, Weituo Hao, Xinyan He, Dong Guo, Jitong Chen, Yuxuan Wang, and Yanmin Qian. InstructME: An instruction guided music edit and remix framework with latent diffusion models. *arXiv preprint arXiv:2308.14360*, 2023.

[HDSSL20]  Alex Hanna, Emily Denton, Andrew Smart, and Jamila Smith-Loud. Towards a critical race methodology in algorithmic fairness. In *Proceedings of the 2020 conference on fairness, accountability, and transparency*, pages 501–512, 2020.

[Her22]    Erik Hermann. Artificial intelligence and mass personalization of communication content—an ethical and literacy perspective. *New media & society*, 24(5):1258–1277, 2022.

[HFW+20]   Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020.

[HGJ+19] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, pages 2790–2799. PMLR, 2019.

[HGM+23] Shibo Hao, Yi Gu, Haodi Ma, Joshua Hong, Zhen Wang, Daisy Wang, and Zhiting Hu. Reasoning with language model is planning with world model. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 8154–8173, 2023.

[HHIE18] Curtis Hawthorne, Anna Huang, Daphne Ippolito, and Douglas Eck. Transformer-nade for piano performances. In *NIPS 2nd Workshop on Machine Learning for Creativity and Design*, 2018.

[HHL+23] Zihao He, Weituo Hao, Wei Tsung Lu, Changyou Chen, Kristina Lerman, and Xuchen Song. ALCAP: alignment-augmented music captioner. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pages 16501–16512. Association for Computational Linguistics, 2023.

[HHRK98] Holger H. Hoos, Keith Hamel, Kai Renz, and Jürgen Kilian. The GUIDO notation format: A novel approach for adequately representing score-level music. In *International Conference on Mathematics and Computing (ICMC)*, 1998.

[HHRK01] Holger Hoos, Keith Hamel, Kai Renz, and Jürgen Killian. Representing score-level music using the guido music-notation format, 2001.

[HHSK23] Rujing Stacy Huang, Andre Holzapfel, Bob L. T. Sturm, and Anna-Kaisa Kaila. Beyond Diverse Datasets: Responsible MIR, Interdisciplinarity, and the Fractured Worlds of Music. *Transactions of the International Society for Music Information Retrieval*, 6(1):43–59, June 2023.

[HHW+17] Chin-Cheng Hsu, Hsin-Te Hwang, Yi-Chiao Wu, Yu Tsao, and Hsin-Min Wang. Voice conversion from unaligned corpora using variational autoencoding wasserstein generative adversarial networks. In *18th Annual Conference of the International Speech Communication Association, Interspeech 2017, Stockholm, Sweden, August 20-24, 2017*, pages 3364–3368. ISCA, 2017.

[HHY+23] Rongjie Huang, Jiawei Huang, Dongchao Yang, Yi Ren, Luping Liu, Mingze Li, Zhenhui Ye, Jinglin Liu, Xiang Yin, and Zhou Zhao. Make-An-Audio: Text-to-audio generation with latent diffusion models. *Proceedings of the International Conference on Machine Learning*, 2023.

[HIY18] Sungeun Hong, Woobin Im, and Hyun Seung Yang. CBVMR: content-based video-music retrieval using soft intra-modal structure constraint. In Kiyoharu Aizawa, Michael S. Lew, and Shin'ichi Satoh, editors, *Proceedings of the 2018 ACM on International Conference on Multimedia Retrieval, ICMR 2018, Yokohama, Japan, June 11-14, 2018*, pages 353–361. ACM, 2018.

[HJA20] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.

[HJC+22] Curtis Hawthorne, Andrew Jaegle, Cătălina Cangea, Sebastian Borgeaud, Charlie Nash, Mateusz Malinowski, Sander Dieleman, Oriol Vinyals, Matthew Botvinick, Ian Simon, et al. General-purpose, long-context autoregressive modeling with perceiver ar. In *International Conference on Machine Learning*, pages 8535–8558. PMLR, 2022.

[HJL+22] Qingqing Huang, Aren Jansen, Joonseok Lee, Ravi Ganti, Judith Yue Li, and Daniel P. W. Ellis. Mulan: A joint embedding of music audio and natural language. In Preeti Rao, Hema A. Murthy, Ajay Srinivasamurthy, Rachel M. Bittner, Rafael Caro Repetto, Masataka Goto, Xavier Serra, and Marius Miron, editors, *Proceedings of the 23rd International Society for Music Information Retrieval Conference, ISMIR 2022, Bengaluru, India, December 4-8, 2022*, pages 559–566, 2022.

[HJL+23] Zhicheng Huang, Xiaojie Jin, Chengze Lu, Qibin Hou, Ming-Ming Cheng, Dongmei Fu, Xiaohui Shen, and Jiashi Feng. Contrastive masked autoencoders are stronger vision learners. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.

[HKK+20] Tom Henighan, Jared Kaplan, Mor Katz, Mark Chen, Christopher Hesse, Jacob Jackson, Heewoo Jun, Tom B Brown, Prafulla Dhariwal, Scott Gray, et al. Scaling laws for autoregressive generative modeling. *arXiv preprint arXiv:2010.14701*, 2020.

[HKW+22] Tri Huynh, Simon Kornblith, Matthew R Walter, Michael Maire, and Maryam Khademi. Boosting contrastive self-supervised learning with false negative cancellation. In *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, pages 2785–2795, 2022.

[HLA+19] Sicong Huang, Qiyang Li, Cem Anil, Xuchan Bao, Sageev Oore, and Roger B. Grosse. Timbretron: Awavenet(yclegan(Cqt(Audio))) Pipeline for Musical Timbre Transfer. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019.

[HLC+24] Yingqing He, Zhaoyang Liu, Jingye Chen, Zeyue Tian, Hongyu Liu, Xiaowei Chi, Runtao Liu, Ruibin Yuan, Yazhou Xing, Wenhai Wang, et al. Llms meet multimodal generation and editing: A survey. *arXiv preprint arXiv:2405.19334*, 2024.

[HLJK23] Susung Hong, Gyuseong Lee, Wooseok Jang, and Seungryong Kim. Improving sample quality of diffusion models using self-attention guidance. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7462–7471, 2023.

[HLL23] Ji-Sang Hwang, Sang-Hoon Lee, and Seong-Whan Lee. Hiddensinger: High-quality singing voice synthesis via neural audio codec and latent diffusion models. *arXiv preprint arXiv:2306.06814*, 2023.

[HLSS23a] Atin Sakkeer Hussain, Shansong Liu, Chenshuo Sun, and Ying Shan. M$^2$ ugen: Multi-modal music understanding and generation with the power of large language models. *arXiv preprint arXiv:2311.11255*, 2023.

[HLSS23b] Atin Sakkeer Hussain, Shansong Liu, Chenshuo Sun, and Ying Shan. M$^2$ugen: Multi-modal music understanding and generation with the power of large language models. *CoRR*, abs/2311.11255, 2023.

[HLW+22] Rongjie Huang, Max W. Y. Lam, Jun Wang, Dan Su, Dong Yu, Yi Ren, and Zhou Zhao. Fastdiff: A fast conditional diffusion model for high-quality speech synthesis. In *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence, IJCAI-22*, pages 4157–4163, 7 2022.

[HLY+23] Rongjie Huang, Mingze Li, Dongchao Yang, Jiatong Shi, Xuankai Chang, Zhenhui Ye, Yuning Wu, Zhiqing Hong, Jiawei Huang, Jinglin Liu, et al. Audiogpt: Understanding and generating speech, music, sound, and talking head. *arXiv preprint arXiv:2304.12995*, 2023.

[HLYY21a] Wen-Yi Hsiao, Jen-Yu Liu, Yin-Cheng Yeh, and Yi-Hsuan Yang. Compound word transformer: Learning to compose full-song music over dynamic directed hypergraphs. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(1):178–186, May 2021.

[HLYY21b] Wen-Yi Hsiao, Jen-Yu Liu, Yin-Cheng Yeh, and Yi-Hsuan Yang. Compound word transformer: Learning to compose full-song music over dynamic directed hypergraphs. In *Proc. AAAI*, 2021.

[HN17] Gaëtan Hadjeres and Frank Nielsen. Interactive music generation with positional constraints using anticipation-rnns. *arXiv preprint arXiv:1709.06404*, 2017.

[HNM+22] William Harvey, Saeid Naderiparizi, Vaden Masrani, Christian Weilbach, and Frank Wood. Flexible diffusion modeling of long videos. *Advances in Neural Information Processing Systems*, 35:27953–27965, 2022.

[HOB23] Carlos Hernandez-Olivan and José R. Beltrán. Music composition with deep learning: A review. In Anupam Biswas, Emile Wennekes, Alicja Wieczorkowska, and Rabul Hussain Laskar, editors, *Advances in Speech and Music Technology: Computational Aspects and Applications*, pages 25–50. Springer Cham, 2023.

[HOHOB22] Carlos Hernandez-Olivan, Javier Hernandez-Olivan, and Jose R Beltran. A survey on artificial intelligence for music generation: Agents, domains and perspectives. *arXiv e-prints*, pages arXiv–2210, 2022.

[HPN17] Gaëtan Hadjeres, François Pachet, and Frank Nielsen. DeepBach: A steerable model for bach chorales generation. *34th*

International Conference on Machine Learning, ICML 2017, 3:2187–2196, 2017.

[HPW+23a] Qingqing Huang, Daniel S Park, Tao Wang, Timo I Denk, Andy Ly, Nanxin Chen, Zhengdong Zhang, Zhishuai Zhang, Jiahui Yu, Christian Frank, et al. Noise2music: Text-conditioned music generation with diffusion models. *arXiv preprint arXiv:2302.03917*, 2023.

[HPW+23b] Qingqing Huang, Daniel S. Park, Tao Wang, Timo I. Denk, Andy Ly, Nanxin Chen, Zhengdong Zhang, Zhishuai Zhang, Jiahui Yu, Christian Havnø Frank, Jesse H. Engel, Quoc V. Le, William Chan, and Wei Han. Noise2music: Text-conditioned music generation with diffusion models. *CoRR*, abs/2302.03917, 2023.

[HRH+23] Jiawei Huang, Yi Ren, Rongjie Huang, Dongchao Yang, Zhenhui Ye, Chen Zhang, Jinglin Liu, Xiang Yin, Zejun Ma, and Zhou Zhao. Make-An-Audio 2: Temporal-enhanced text-to-audio generation. *arXiv preprint arXiv:2305.18474*, 2023.

[HS21] Jui-Yang Hsu and Li Su. VOCANO: A note transcription framework for singing voice in polyphonic music. In *Proceedings of the 22nd International Society for Music Information Retrieval Conference, ISMIR 2021, Online, November 7-12, 2021*, pages 293–300, 2021.

[HS22a] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.

[HS22b] Wei-Ning Hsu and Bowen Shi. u-HuBERT: Unified Mixed-Modal Speech Pretraining And Zero-Shot Transfer to Unlabeled Modality. In *Advances in Neural Information Processing Systems*, volume 35, pages 21157–21170, New Orleans, LA, USA, December 2022.

[HSC18a] Andre Holzapfel, Bob Sturm, and Mark Coeckelbergh. Ethical dimensions of music information retrieval technology. *Transactions of the International Society for Music Information Retrieval*, 1(1):44–55, 2018.

[HSC18b] Andre Holzapfel, Bob L. Sturm, and Mark Coeckelbergh. Ethical Dimensions of Music Information Retrieval Technology. *Transactions of the International Society for Music Information Retrieval*, 1(1):44–55, September 2018.

[HSC+22] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *The Journal of Machine Learning Research*, 23(1):2249–2281, 2022.

[HSF+24] B Hayes, J Shier, G Fazekas, A McPherson, and C Saitis. A review of differentiable digital signal processing for music and speech synthesis. *Frontiers in Signal Processing*, 2024.

[HSH21] Rujing Stacy Huang, Bob L. T. Sturm, and Andre Holzapfel. De-centering the west: East asian philosophies and the ethics of applying artificial intelligence to music. In Jin Ha Lee, Alexander Lerch, Zhiyao Duan, Juhan Nam, Preeti Rao, Peter van Kranenburg, and Ajay Srinivasamurthy, editors, *Proceedings of the 22nd International Society for Music Information Retrieval Conference, ISMIR 2021, Online, November 7-12, 2021*, pages 301–309, 2021.

[HSR+18] Curtis Hawthorne, Andriy Stasyuk, Adam Roberts, Ian Simon, Cheng-Zhi Anna Huang, Sander Dieleman, Erich Elsen, Jesse Engel, and Douglas Eck. Enabling factorized no music modeling and generation with the maestro dataset. *arXiv preprint arXiv:1810.12247*, 2018.

[HSR+22] Curtis Hawthorne, Ian Simon, Adam Roberts, Neil Zeghidour, Josh Gardner, Ethan Manilow, and Jesse Engel. Multi-instrument Music Synthesis with Spectrogram Diffusion. In *Proceeding of the International Society on Music Information Retrieval (ISMIR)*, Bengaluru, India, 2022.

[HSX+24] Po-Yao Huang, Vasu Sharma, Hu Xu, Chaitanya Ryali, Yanghao Li, Shang-Wen Li, Gargi Ghosh, Jitendra Malik, Christoph Feichtenhofer, et al. Mavil: Masked audio-video learners. *Advances in Neural Information Processing Systems*, 36, 2024.

[Hur02] David Huron. Music information processing using the humdrum toolkit: Concepts, examples, and lessons. *Computer Music Journal*, 26(2):11–26, 2002.

[HVU+18a] Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Noam Shazeer, Ian Simon, Curtis Hawthorne, Andrew M Dai, Matthew D Hoffman, Monica Dinculescu, and Douglas Eck. Music transformer. *arXiv preprint arXiv:1809.04281*, 2018.

[HVU+18b] Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Noam Shazeer, Ian Simon, Curtis Hawthorne, Andrew M Dai, Matthew D Hoffman, Monica Dinculescu, and Douglas Eck. Music transformer. *arXiv preprint arXiv:1809.04281*, 2018.

[HWAZ+21] Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2021.

[HWY+23] Hao Huang, Lin Wang, Jichen Yang, Ying Hu, and Liang He. W2VC: wavlm representation based one-shot voice conversion with gradient reversal distillation and CTC supervision. *EURASIP J. Audio Speech Music. Process.*, 2023(1):45, 2023.

[HXL+22a] Po-Yao Huang, Hu Xu, Juncheng Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian Metze, and Christoph Feichtenhofer. Masked autoencoders that listen. *Advances in Neural Information Processing Systems*, 35:28708–28720, 2022.

[HXL+22b] Po-Yao Huang, Hu Xu, Juncheng B. Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian Metze, and Christoph Feichtenhofer. Masked Autoencoders that Listen. In *Advances in Neural Information Processing Systems*, New Orleans, LA, USA, October 2022.

[HY20a] Yu Siang Huang and Yi Hsuan Yang. Pop music transformer: Beat-based modeling and generation of expressive pop piano compositions. In *Proceedings of the 28th ACM International Conference on Multimedia*, 2020.

[HY20b] Yu-Siang Huang and Yi-Hsuan Yang. Pop music transformer: Beat-based modeling and generation of expressive pop piano compositions. In *Proceedings of the 28th ACM international conference on multimedia*, pages 1180–1188, 2020.

[HZZ+24] Tao Huang, Yuan Zhang, Mingkai Zheng, Shan You, Fei Wang, Chen Qian, and Chang Xu. Knowledge diffusion for distillation. *Advances in Neural Information Processing Systems*, 36, 2024.

[IEE17] IEEE Global Initiative on Ethics of Autonomous and Intelligent Systems. Glossary for discussion of ethics of autonomous and intelligent systems, version 1, 2017. Available at https://standards.ieee.org/wp-content/uploads/import/documents/other/eadv2_glossary.pdf (Retrieved 21/06/2024).

[IMDS20] Stylianos Ioannis Mimilakis, Konstantinos Drossos, and Gerald Schuller. Revisiting representation learning for singing voice separation with sinkhorn distances. *arXiv e-prints*, pages arXiv–2007, 2020.

[JBEW19] Andreas Jansson, Rachel M. Bittner, Sebastian Ewert, and Tillman Weyde. Joint singing voice separation and F0 estimation with deep u-net architectures. In *27th European Signal Processing Conference, EUSIPCO 2019, A Coruña, Spain, September 2-6, 2019*, pages 1–5. IEEE, 2019.

[JCL20] Chang-Bin Jeon, Hyeong-Seok Choi, and Kyogu Lee. Exploring aligned lyrics-informed singing voice separation. In *Proceedings of the 21th International Society for Music Information Retrieval Conference, ISMIR 2020, Montreal, Canada, October 11-16, 2020*, pages 685–692, 2020.

[JdSBTK17] Jim Jones, Diego de Siqueira Braga, Kleber Tertuliano, and Tomi Kauppinen. Musicowl: The music score ontology. In *Proceedings of the International Conference on Web Intelligence*, pages 1222–1229, 2017.

[JETT+23] Jaladi Sam Joel, B. Ernest Thompson, Steve Renny Thomas, T. Revanth Kumar, Shajin Prince, and D Bini. Emotion based Music Recommendation System using Deep Learning Model. In *2023 International Conference on Inventive Computation Technologies (ICICT)*, pages 227–232, April 2023.

[JG82] Biing-Hwang Juang and A. Gray. Multiple stage vector quantization for speech coding. In *ICASSP '82. IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 7, pages 597–600, May 1982.

[JJDZ20] Nan Jiang, Sheng Jin, Zhiyao Duan, and Changshui Zhang. RL-Duet: Online music accompaniment generation using deep reinforcement learning. *AAAI 2020 - 34th AAAI Conference on Artificial Intelligence*, pages 710–718, 2020.

[JKC+21] Myeonghun Jeong, Hyeongju Kim, Sung Jun Cheon, Byoung Jin Choi, and Nam Soo Kim. Diff-tts: A denoising diffusion model for text-to-speech. In *Interspeech*, pages 3605–3609, 2021.

[JKK+19] Dasaem Jeong, Taegyun Kwon, Yoojin Kim, Kyogu Lee, and Juhan Nam. Virtuosonet: A hierarchical rnn-based system for modeling expressive piano performance. In *ISMIR*, pages 908–915, 2019.

[JKKN19] Dasaem Jeong, Taegyun Kwon, Yoojin Kim, and Juhan Nam. Graph neural network for music score data and modeling

expressive piano performance. In *Proceedings of the 36th International Conference on Machine Learning (ICML)*, 2019.

[JLY20] Shulei Ji, Jing Luo, and Xinyu Yang. A comprehensive survey on deep music generation: Multi-level representations, algorithms, evaluations, and future directions. *arXiv preprint arXiv:2011.06801*, 2020.

[JP23] Udeme Samuel Jacob and Jace Pillay. A Systematic Review of Scientific Studies on the Effects of Music Therapy on Individuals with Autism Spectrum Disorder. *Journal for ReAttach Therapy and Developmental Diversities*, 6(6s):545–558, June 2023.

[JPTI21] Ali Jahanian, Xavier Puig, Yonglong Tian, and Phillip Isola. Generative models as a data source for multiview representation learning. *arXiv preprint arXiv:2106.05258*, 2021.

[JSM+23] Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.

[JWS+23] Tejas Jayashankar, Jilong Wu, Leda Sari, David Kant, Vimal Manohar, and Qing He. Self-supervised representations for singing voice conversion. In *IEEE International Conference on Acoustics, Speech and Signal Processing ICASSP 2023, Rhodes Island, Greece, June 4-10, 2023*, pages 1–5. IEEE, 2023.

[JXC+20] Junyan Jiang, Gus G. Xia, Dave B. Carlton, Chris N. Anderson, and Ryan H. Miyakawa. Transformer VAE: A Hierarchical Model for Structure-Aware and Interpretable Music Representation Learning. In *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 516–520, May 2020. ISSN: 2379-190X.

[JYL23] Shulei Ji, Xinyu Yang, and Jing Luo. A survey on deep learning for symbolic music generation: Representations, algorithms, evaluations, and challenges. *ACM Computing Surveys*, 2023.

[KA21] Maximilian Kasy and Rediet Abebe. Fairness, equality, and power in algorithmic decision-making. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 576–586, 2021.

[KCI+20] Qiuqiang Kong, Yin Cao, Turab Iqbal, Yuxuan Wang, Wenwu Wang, and Mark D Plumbley. Panns: Large-scale pretrained audio neural networks for audio pattern recognition. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 28:2880–2894, 2020.

[KFW23] Emmanouil Karystinaios, Francesco Foscarin, and Gerhard Widmer. Musical voice separation as link prediction: Modeling a musical perception task as a multi-trajectory tracking problem. *IJCAI International Joint Conference on Artificial Intelligence*, pages 3866–3874, 2023.

[KG16] Nadine Kroher and Emilia Gómez. Automatic transcription of flamenco singing from polyphonic music recordings. *IEEE ACM Trans. Audio Speech Lang. Process.*, 24(5):901–913, 2016.

[KK23] Arash Vahdat Karsten Kreis, Ruiqi Gao. Latent diffusion models: Is the generative ai revolution happening in latent space?, 2023. https://neurips2023-ldm-tutorial.github.io/.

[KKB20a] Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis. *Advances in neural information processing systems*, 33:17022–17033, 2020.

[KKB20b] Jungil Kong, Jaehyeon Kim, and Jaekyoung Bae. Hifi-gan: Generative adversarial networks for efficient and high fidelity speech synthesis. *Advances in Neural Information Processing Systems*, 33:17022–17033, 2020.

[KKDB+19] Kundan Kumar, Rithesh Kumar, Thibault De Boissiere, Lucas Gestin, Wei Zhen Teoh, Jose Sotelo, Alexandre De Brebisson, Yoshua Bengio, and Aaron C Courville. Melgan: Generative adversarial networks for conditional waveform synthesis. *Advances in neural information processing systems*, 32, 2019.

[KKJK23] Sungjae Kim, Yewon Kim, Jewoo Jun, and Injung Kim. Muse-svs: Multi-singer emotional singing voice synthesizer that controls emotional intensity. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2023.

[KKK24] Chan-Young Kwon, Hyunsu Kim, and Sung-Hee Kim. The Modernization of Oriental Music Therapy: Five-Element Music Therapy Combined with Artificial Intelligence. *Healthcare*, 12(3):411, January 2024.

[KKL20] Nikita Kitaev, Lukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In *International Conference on Learning Representations*, 2020.

[KKS21] Jaehyeon Kim, Jungil Kong, and Juhee Son. Conditional variational autoencoder with adversarial learning for end-to-end text-to-speech. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 5530–5540. PMLR, 2021.

[KKY22a] Heeseung Kim, Sungwon Kim, and Sungroh Yoon. Guided-TTS: A diffusion model for text-to-speech via classifier guidance. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162, pages 11119–11133, 2022.

[KKY22b] Sungwon Kim, Heeseung Kim, and Sungroh Yoon. Guided-tts 2: A diffusion model for high-quality adaptive text-to-speech with untranscribed data. *arXiv preprint arXiv:2205.15370*, 2022.

[KLCW20] Qiuqiang Kong, Bochen Li, Jitong Chen, and Yuxuan Wang. Giantmidi-piano: A large-scale midi dataset for classical piano music. *arXiv preprint arXiv:2010.07061*, 2020.

[KLK+22] Sangeun Kum, Jongpil Lee, Keunhyoung Luke Kim, Taehyoung Kim, and Juhan Nam. Pseudo-label transfer from frame-level to note-level in a teacher-student framework for singing transcription from polyphonic music. In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2022, Virtual and Singapore, 23-27 May 2022*, pages 796–800. IEEE, 2022.

[KLL+22] Jean Kaddour, Aengus Lynch, Qi Liu, Matt J. Kusner, and Ricardo Silva. Causal machine learning: A survey and open problems, 2022.

[KLP+22] Eugene Kharitonov, Ann Lee, Adam Polyak, Yossi Adi, Jade Copet, Kushal Lakhotia, Tu Anh Nguyen, Morgane Riviere, Abdelrahman Mohamed, Emmanuel Dupoux, and Wei-Ning Hsu. Text-free prosody-aware generative spoken language modeling. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 8666–8681, 2022.

[KLS21a] Otto Kässi, Vili Lehdonvirta, and Fabian Stephany. How many online workers are there in the world? a data-driven assessment. *Open Research Europe*, 1, 2021.

[KLS+21b] Qiuqiang Kong, Bochen Li, Xuchen Song, Yuan Wan, and Yuxuan Wang. High-resolution piano transcription with pedals by regressing onset and offset times. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:3707–3717, 2021.

[KMH+20] Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.

[KPC+22] Felix Kreuk, Adam Polyak, Jade Copet, Eugene Kharitonov, Tu Anh Nguyen, Morgan Rivière, Wei-Ning Hsu, Abdelrahman Mohamed, Emmanuel Dupoux, and Yossi Adi. Textless speech emotion conversion using discrete & decomposed representations. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11200–11214, 2022.

[KPH+21] Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion model for audio synthesis. In *International Conference on Learning Representations*, 2021.

[KPH23] Jaeyong Kang, Soujanya Poria, and Dorien Herremans. Video2music: Suitable music generation from videos using an affective multimodal transformer model. *CoRR*, abs/2311.00968, 2023.

[Kra91] Mark A. Kramer. Nonlinear principal component analysis using autoassociative neural networks. *AIChE Journal*, 37(2):233–243, 1991.

[Kru14] Anna M. Kruspe. Improving singing language identification through i-vector extraction. In *Proceedings of the 17th International Conference on Digital Audio Effects, DAFx-14, Erlangen, Germany, September 1-5, 2014*, pages 227–233, 2014.

[KSEZW21] Khaled Koutini, Jan Schlüter, Hamid Eghbal-Zadeh, and

Gerhard Widmer. Efficient training of audio transformers with patchout. *arXiv preprint arXiv:2110.05069*, 2021.

[KSL+23] Rithesh Kumar, Prem Seetharaman, Alejandro Luebs, Ishaan Kumar, and Kundan Kumar. High-fidelity audio compression with improved rvqgan. *Advances in neural information processing systems*, 2023.

[KSL+24] Rithesh Kumar, Prem Seetharaman, Alejandro Luebs, Ishaan Kumar, and Kundan Kumar. High-fidelity audio compression with improved rvqgan. *Advances in Neural Information Processing Systems*, 36, 2024.

[KSP+22] Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre Défossez, Jade Copet, Devi Parikh, Yaniv Taigman, and Yossi Adi. Audiogen: Textually guided audio generation. *arXiv preprint arXiv:2209.15352*, 2022.

[KSP+23] Felix Kreuk, Gabriel Synnaeve, Adam Polyak, Uriel Singer, Alexandre Défossez, Jade Copet, Devi Parikh, Yaniv Taigman, and Yossi Adi. Audiogen: Textually guided audio generation. In *The Eleventh International Conference on Learning Representations*, 2023.

[KSPH21] Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Advances in neural information processing systems*, 34:21696–21707, 2021.

[KT19] Jacob Devlin Ming-Wei Chang Kenton and Lee Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of naacL-HLT*, volume 1, page 2, 2019.

[KVB+23] Eugene Kharitonov, Damien Vincent, Zalán Borsos, Raphaël Marinier, Sertan Girgin, Olivier Pietquin, Matt Sharifi, Marco Tagliasacchi, and Neil Zeghidour. Speak, read and prompt: High-fidelity text-to-speech with minimal supervision. *arXiv preprint arXiv:2302.03540*, 2023.

[KW14] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *International Conference on Learning Representations*, Banff, Canada, 2014.

[KWB+10] Gregory Kramer, Bruce Walker, Terri Bonebright, Perry Cook, John H. Flowers, Nadine Miner, and John Neuhoff. Sonification Report: Status of the Field and Research Agenda. Technical Report 444, Faculty Publications, Department of Psychology, 2010.

[KWMZ20] A. Sophia Koepke, Olivia Wiles, Yael Moses, and Andrew Zisserman. Sight to sound: An end-to-end approach for visual piano transcription. In *2020 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2020, Barcelona, Spain, May 4-8, 2020*, pages 1838–1842. IEEE, 2020.

[KZRS19] Kevin Kilgour, Mauricio Zuluaga, Dominik Roblek, and Matthew Sharifi. Fréchet audio distance: A reference-free metric for evaluating music enhancement algorithms. In *Proc. Interspeech*, 2019.

[LAD+22] Aitor Lewkowycz, Anders Johan Andreassen, David Dohan, Ethan Dyer, Henryk Michalewski, Vinay Venkatesh Ramasesh, Ambrose Slone, Cem Anil, Imanol Schlag, Theo Gutman-Solo, et al. Solving quantitative reasoning problems with language models. In *Advances in Neural Information Processing Systems*, 2022.

[LAH19] Yin-Jyun Luo, Kat Agres, and Dorien Herremans. Learning Disentangled Representations of Timbre and Pitch for Musical Instrument Sounds Using Gaussian Mixture Variational Autoencoders. In *Proceedings of the International Society for Music Information Retrieval (ISMIR) Conference*, Delft, The Netherlands, 2019.

[LARC21] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059, 2021.

[LCKL20] Juheon Lee, Hyeong-Seok Choi, Junghyun Koo, and Kyogu Lee. Disentangling timbre and singing style with multi-singer singing synthesis system. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7224–7228. IEEE, 2020.

[LCL22] Juheon Lee, Hyeong-Seok Choi, and Kyogu Lee. Expressive singing synthesis using local style token and dual-path pitch encoder. *arXiv preprint arXiv:2204.03249*, 2022.

[LCP+23] Siddique Latif, Heriberto Cuayáhuitl, Farrukh Pervez, Fahad Shamshad, Hafiz Shehbaz Ali, and Erik Cambria. A survey on deep reinforcement learning for audio-based applications. *Artificial Intelligence Review*, 56(3):2193–2240, 2023.

[LCSM21] Songxiang Liu, Yuewen Cao, Dan Su, and Helen Meng. Diffsvc: A diffusion probabilistic model for singing voice conversion. In *IEEE Automatic Speech Recognition and Understanding Workshop, ASRU 2021, Cartagena, Colombia, December 13-17, 2021*, pages 741–748. IEEE, 2021.

[LCW+22] Ann Lee, Peng-Jen Chen, Changhan Wang, Jiatao Gu, Sravya Popuri, Xutai Ma, Adam Polyak, Yossi Adi, Qing He, Yun Tang, Juan Pino, and Wei-Ning Hsu. Direct speech-to-speech translation with discrete units. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3327–3339, 2022.

[LCY+23a] Peike Li, Boyu Chen, Yao Yao, Yikai Wang, Allen Wang, and Alex Wang. JEN-1: text-guided universal music generation with omnidirectional diffusion models. *CoRR*, abs/2308.04729, 2023.

[LCY+23b] Peike Li, Boyu Chen, Yao Yao, Yikai Wang, Allen Wang, and Alex Wang. Jen-1: Text-guided universal music generation with omnidirectional diffusion models. *arXiv preprint arXiv:2308.04729*, 2023.

[LCY+23c] Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. AudioLDM: Text-to-audio generation with latent diffusion models. *Proceedings of the International Conference on Machine Learning*, 2023.

[LCY+23d] Haohe Liu, Zehua Chen, Yi Yuan, Xinhao Mei, Xubo Liu, Danilo Mandic, Wenwu Wang, and Mark D Plumbley. AudioLDM: Text-to-audio generation with latent diffusion models. In *Proc. ICML*, 2023.

[LDC+22] Jiafeng Liu, Yuanliang Dong, Zehua Cheng, Xinran Zhang, Xiaobing Li, Feng Yu, and Maosong Sun. Symphony generation with permutation invariant language model. *arXiv preprint arXiv:2205.05448*, 2022.

[LDJ23] Minhee Lee, Seungheon Doh, and Dasaem Jeong. Annotator subjectivity in the musiccaps dataset. In Lorenzo Porcaro, Roser Batlle-Roca, and Emilia Gómez, editors, *Proceedings of the 2nd Workshop on Human-Centric Music Information Retrieval 2023 co-located with the 24th International Society for Music Information Retrieval Conference (ISMIR 2023), Milan, Italy, November 10, 2023*, volume 3528 of *CEUR Workshop Proceedings*. CEUR-WS.org, 2023.

[Lee16] Jin Hyung Lee. The Effects of Music on Pain: A Meta-Analysis. *Journal of Music Therapy*, 53(4):430–477, December 2016.

[Ler22] Alexander Lerch. *An introduction to audio content analysis: Music Information Retrieval tasks and applications*. John Wiley & Sons, 2022.

[Les01] Lawrence Lessig. *The Future of Ideas: The Fate of the Commons in a Connected World*. Random House Inc., USA, 2001.

[LGFW24] Luca Lanzendörfer, Florian Grötschla, Emil Funke, and Roger Wattenhofer. Disco-10m: A large-scale music dataset. *Advances in Neural Information Processing Systems*, 36, 2024.

[LGJS17] Feynman Liang, Mark Gotham, Matthew Johnson, and Jamie Shotton. Automatic Stylistic Composition of Bach Chorales with Deep LSTM. In *Proceeding of the 18th International Society on Music Information Retrieval (ISMIR)*, 2017.

[LGW18] Stefan Lattner, Maarten Grachten, and Gerhard Widmer. Imposing higher-level structure in polyphonic music generation using convolutional restricted boltzmann machines and constraints. *Journal of Creative Music Systems*, 2:1–31, 2018.

[LHAH20] Yin-Jyun Luo, Chin-Cheng Hsu, Kat Agres, and Dorien Herremans. Singing voice conversion with disentangled representations of singer and vocal technique using variational autoencoders. In *2020 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2020, Barcelona, Spain, May 4-8, 2020*, pages 3277–3281. IEEE, 2020.

[LHSS23a] Shansong Liu, Atin Sakkeer Hussain, Chenshuo Sun, and Ying Shan. Music understanding llama: Advancing text-to-music generation with question answering and captioning. *CoRR*, abs/2308.11276, 2023.

[LHSS23b] Shansong Liu, Atin Sakkeer Hussain, Chenshuo Sun, and Ying Shan. Music understanding llama: Advancing text-to-music generation with question answering and captioning. *arXiv preprint arXiv:2308.11276*, 2023.

[LHSS24] Shansong Liu, Atin Sakkeer Hussain, Chenshuo Sun, and Ying Shan. Music understanding llama: Advancing text-to-

music generation with question answering and captioning. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 286–290. IEEE, 2024.

[Lin04] Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain, July 2004. Association for Computational Linguistics.

[LK19] Bochen Li and Aparna Kumar. Query by video: Cross-modal music retrieval. In *ISMIR*, pages 604–611, 2019.

[LKH+21] Kushal Lakhotia, Eugene Kharitonov, Wei-Ning Hsu, Yossi Adi, Adam Polyak, Benjamin Bolte, Tu-Anh Nguyen, Jade Copet, Alexei Baevski, Abdelrahman Mohamed, and Emmanuel Dupoux. On generative spoken language modeling from raw audio. *Transactions of the Association for Computational Linguistics*, 9:1336–1354, 2021.

[LKHS20] Phuc H Le-Khac, Graham Healy, and Alan F Smeaton. Contrastive representation learning: A framework and review. *Ieee Access*, 8:193907–193934, 2020.

[LKJX21] Liwei Lin, Qiuqiang Kong, Junyan Jiang, and Gus Xia. A unified model for zero-shot music source separation, transcription and synthesis. *arXiv preprint arXiv:2108.03456*, 2021.

[LKL+22] Sang-Hoon Lee, Seung-Bin Kim, Ji-Hyun Lee, Eunwoo Song, Min-Jae Hwang, and Seong-Whan Lee. Hierspeech: Bridging the gap between text and speech by hierarchical variational inference using self-supervised representations for speech synthesis. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 - December 9, 2022*, 2022.

[LL21] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4582–4597, 2021.

[LLD+18] Bochen Li, Xinzhao Liu, Karthik Dinesh, Zhiyao Duan, and Gaurav Sharma. Creating a multitrack classical music performance dataset for multimodal music analysis: Challenges, insights, and applications. *IEEE Transactions on Multimedia*, 21(2):522–535, 2018.

[LLD24] Xia Liang, Jiaju Lin, and Xinjian Du. Bytecomposer: a human-like melody composition method based on language model agent. *arXiv preprint arXiv:2402.17785*, 2024.

[LLLC22] Ya Li, Xiulai Li, Zheng Lou, and Chaofan Chen. Long Short-Term Memory-Based Music Analysis System for Music Therapy. *Frontiers in Psychology*, 13, June 2022.

[LLR+22a] Jinglin Liu, Chengxi Li, Yi Ren, Feiyang Chen, and Zhou Zhao. Diffsinger: Singing voice synthesis via shallow diffusion mechanism. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelveth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022*, pages 11020–11028. AAAI Press, 2022.

[LLR+22b] Jinglin Liu, Chengxi Li, Yi Ren, Feiyang Chen, and Zhou Zhao. Diffsinger: Singing voice synthesis via shallow diffusion mechanism. In *AAAI*, pages 11020–11028. AAAI Press, 2022.

[LLSH23] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023.

[LLY+20] Zheng Liu, Jianxun Lian, Junhan Yang, Defu Lian, and Xing Xie. Octopus: Comprehensive and elastic user representation for the generation of recommendation candidates. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 289–298, 2020.

[LM19] Yi Luo and Nima Mesgarani. Conv-tasnet: Surpassing ideal time–frequency magnitude masking for speech separation. *IEEE/ACM transactions on audio, speech, and language processing*, 27(8):1256–1266, 2019.

[LMW+24] Dichucheng Li, Yinghao Ma, Weixing Wei, Qiuqiang Kong, Yulun Wu, Mingjin Che, Fan Xia, Emmanouil Benetos, and Wei Li. Mertech: Instrument playing technique detection us-

ing self-supervised pretrained model with multi-task finetuning. In *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 521–525, 2024.

[LN19] Kyungyun Lee and Juhan Nam. Learning a joint embedding space of monophonic and mixed music signals for singing voice. In *Proceedings of the 20th International Society for Music Information Retrieval Conference, ISMIR 2019, Delft, The Netherlands, November 4-8, 2019*, pages 295–302, 2019.

[Lon23] Mitchell Longen. A System Out of Balance: A Critical Analysis of Philosophical Justifications for Copyright Law Through the Lenz of Users' Rights. *University of Michigan Journal of Law Reform*, 56:779–826, 2023.

[LPKN17] Jongpil Lee, Jiyoung Park, Keunhyoung Luke Kim, and Juhan Nam. Sample-level deep convolutional neural networks for music auto-tagging using raw waveforms. *arXiv preprint arXiv:1703.01789*, 2017.

[LPP+20] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474, 2020.

[LPV+23] Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time intervention: Eliciting truthful answers from a language model. In *Advances in Neural Information Processing Systems*, volume 36, pages 41451–41530, 2023.

[LQZ+24a] Sizhe Li, Yiming Qin, Minghang Zheng, Xin Jin, and Yang Liu. Diff-BGM: A Diffusion Model for Video Background Music Generation. In *Proceedings of the Conference on Computer Vision and Pattern Recognition*, 2024.

[LQZ+24b] Sizhe Li, Yiming Qin, Minghang Zheng, Xin Jin, and Yang Liu. Diff-bgm: A diffusion model for video background music generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 27348–27357, 2024.

[LR18] Steven R Livingstone and Frank A Russo. The ryerson audio-visual database of emotional speech and song (ravdess): A dynamic, multimodal set of facial and vocal expressions in north american english. *PloS one*, 13(5):e0196391, 2018.

[LRL17] Hyungui Lim, Seungyeon Rhyu, and Kyogu Lee. Chord generation from symbolic melody using blstm networks. *arXiv preprint arXiv:1712.01011*, 2017.

[LS23a] Shuyu Li and Yunsick Sung. Melodydiffusion: chord-conditioned melody generation using a transformer-based diffusion model. *Mathematics*, 11(8):1915, 2023.

[LS23b] Shuyu Li and Yunsick Sung. MRBERT: Pre-Training of Melody and Rhythm for Automatic Music Generation. *Mathematics*, 11(4):798, January 2023. Number: 4 Publisher: Multidisciplinary Digital Publishing Institute.

[LSS+23] Siddique Latif, Moazzam Shoukat, Fahad Shamshad, Muhammad Usama, Heriberto Cuayáhuitl, and Björn W Schuller. Sparks of large audio models: A survey and outlook. *arXiv preprint arXiv:2308.12792*, 2023.

[LSZ+21] Tingtian Li, Zixun Sun, Haoruo Zhang, Jin Li, Ziming Wu, Hui Zhan, Yipeng Yu, and Hengcan Shi. Deep music retrieval for fine-grained videos by exploiting cross-modal-encoded voice-overs. In Fernando Diaz, Chirag Shah, Torsten Suel, Pablo Castells, Rosie Jones, and Tetsuya Sakai, editors, *SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021*, pages 1880–1884. ACM, 2021.

[LTL+23a] Max WY Lam, Qiao Tian, Tang Li, Zongyu Yin, Siyuan Feng, Ming Tu, Yuliang Ji, Rui Xia, Mingbo Ma, Xuchen Song, et al. Efficient neural music generation. *arXiv preprint arXiv:2305.15719*, 2023.

[LTL+23b] Ang Lv, Xu Tan, Peiling Lu, Wei Ye, Shikun Zhang, Jiang Bian, and Rui Yan. Getmusic: Generating any music tracks with a unified representation and diffusion framework. *arXiv preprint arXiv:2305.10841*, 2023.

[LTY+21] Zhonghao Li, Benlai Tang, Xiang Yin, Yuan Wan, Ling Xu, Chen Shen, and Zejun Ma. Ppg-based singing voice conversion with adversarial representation learning. In *IEEE International Conference on Acoustics, Speech and Signal*

*Processing, ICASSP 2021, Toronto, ON, Canada, June 6-11, 2021*, pages 7073–7077. IEEE, 2021.

[LTY+23a] Haohe Liu, Qiao Tian, Yi Yuan, Xubo Liu, Xinhao Mei, Qiuqiang Kong, Yuping Wang, Wenwu Wang, Yuxuan Wang, and Mark D. Plumbley. AudioLDM 2: Learning holistic audio generation with self-supervised pretraining. *arXiv preprint arXiv:2308.05734*, 2023.

[LTY+23b] Haohe Liu, Qiao Tian, Yi Yuan, Xubo Liu, Xinhao Mei, Qiuqiang Kong, Yuping Wang, Wenwu Wang, Yuxuan Wang, and Mark D Plumbley. AudioLDM 2: Learning holistic audio generation with self-supervised pretraining. *arXiv preprint arXiv:2308.05734*, 2023.

[LWD21] Bochen Li, Yuxuan Wang, and Zhiyao Duan. Audiovisual singing voice separation. *arXiv preprint arXiv:2107.00231*, 2021.

[LWKH24] Wei-Tsung Lu, Ju-Chiang Wang, Qiuqiang Kong, and Yun-Ning Hung. Music source separation with band-split rope transformer. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 481–485. IEEE, 2024.

[LWL+20] Peiling Lu, Jie Wu, Jian Luan, Xu Tan, and Li Zhou. Xiaoicesing: A high-quality and integrated singing voice synthesis system. *arXiv preprint arXiv:2006.06261*, 2020.

[LWSY22] Max W. Y. Lam, Jun Wang, Dan Su, and Dong Yu. BDDM: Bilateral denoising diffusion models for fast and high-quality speech synthesis. In *International Conference on Learning Representations*, 2022.

[LWW+21] Wei-Tsung Lu, Ju-Chiang Wang, Minz Won, Keunwoo Choi, and Xuchen Song. Spectnt: A time-frequency transformer for music audio. *arXiv preprint arXiv:2110.09127*, 2021.

[LXJZ23] Liwei Lin, Gus Xia, Junyan Jiang, and Yixiao Zhang. Content-based controls for music large language modeling. *arXiv preprint arXiv:2310.17162*, 2023.

[LXK+23a] Peiling Lu, Xin Xu, Chenfei Kang, Botao Yu, Chengyi Xing, Xu Tan, and Jiang Bian. Musecoco: Generating symbolic music from text. *arXiv preprint arXiv:2306.00110*, 2023.

[LXK+23b] Peiling Lu, Xin Xu, Chenfei Kang, Botao Yu, Chengyi Xing, Xu Tan, and Jiang Bian. MuseCoco: Generating Symbolic Music from Text, May 2023. arXiv:2306.00110 [cs, eess].

[LXK+23c] Peiling Lu, Xin Xu, Chenfei Kang, Botao Yu, Chengyi Xing, Xu Tan, and Jiang Bian. Musecoco: Generating symbolic music from text. *arXiv preprint arXiv:2306.00110*, 2023.

[LXY+22] Zhaoyang Lyu, Xudong Xu, Ceyuan Yang, Dahua Lin, and Bo Dai. Accelerating diffusion models via early stop of the diffusion process. *arXiv preprint arXiv:2205.12524*, 2022.

[LXY+24] Haohe Liu, Xuenan Xu, Yi Yuan, Mengyue Wu, Wenwu Wang, and Mark D Plumbley. Semanticodec: An ultra low bitrate semantic audio codec for general sound. *arXiv preprint arXiv:2405.00233*, 2024.

[LXZJ24] Liwei Lin, Gus Xia, Yixiao Zhang, and Junyan Jiang. Arrange, inpaint, and refine: Steerable long-term music audio generation and editing via content-based controls. *arXiv preprint arXiv:2402.09508*, 2024.

[LYC+22] Haoying Li, Yifan Yang, Meng Chang, Shiqi Chen, Huajun Feng, Zhihai Xu, Qi Li, and Yueting Chen. Srdiff: Single image super-resolution with diffusion probabilistic models. *Neurocomputing*, 479:47–59, 2022.

[LYDB17] Wang Ling, Dani Yogatama, Chris Dyer, and Phil Blunsom. Program induction by rationale generation: Learning to solve and explain algebraic word problems. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 158–167, 2017.

[Lyo17] Richard F Lyon. *Human and Machine Hearing: Extracting Meaning from Sound*. Cambridge University Press, Cambridge, 2017.

[LYT+24] Jiajia Li, Lu Yang, Mingni Tang, Cong Chen, Zuchao Li, Ping Wang, and Hai Zhao. The music maestro or the musically challenged, a massive music evaluation benchmark for large language models. *arXiv preprint arXiv:2406.15885*, 2024.

[LYX+24] Yiwen Lu, Zhen Ye, Wei Xue, Xu Tan, Qifeng Liu, and Yike Guo. Comosvc: Consistency model-based singing voice conversion. *arXiv preprint arXiv:2401.01792*, 2024.

[LYZ+22] Yizhi Li, Ruibin Yuan, Ge Zhang, Yinghao Ma, Chenghua Lin, Xingran Chen, Anton Ragni, Hanzhi Yin, Zhijie Hu, Haoyu He, et al. Map-music2vec: A simple and effective baseline for self-supervised music audio representation learn-

ing. *International Society for Music Information Retrieval, late-breaking demo*, 2022.

[LYZ+24] Yizhi Li, Ruibin Yuan, Ge Zhang, Yinghao Ma, Xingran Chen, Hanzhi Yin, Chenghua Lin, Anton Ragni, Emmanouil Benetos, Norbert Gyenge, et al. Mert: Acoustic music understanding model with large-scale self-supervised training. *International Conference on Learning Representations*, 2024.

[LZL+23] Xubo Liu, Zhongkai Zhu, Haohe Liu, Yi Yuan, Meng Cui, Qiushi Huang, Jinhua Liang, Yin Cao, Qiuqiang Kong, Mark D Plumbley, et al. Wavjourney: Compositional audio creation with large language models. *arXiv preprint arXiv:2307.14335*, 2023.

[LZL+24] Jinhua Liang, Huan Zhang, Haohe Liu, Yin Cao, Qiuqiang Kong, Xubo Liu, Wenwu Wang, Mark D. Plumbley, Huy Phan, and Emmanouil Benetos. Wavcraft: Audio editing and generation with large language models. *arXiv preprint arXiv:2403.09527*, 2024.

[LZSL20] Junchen Lu, Kun Zhou, Berrak Sisman, and Haizhou Li. VAW-GAN for singing voice conversion with non-parallel training data. In *Asia-Pacific Signal and Information Processing Association Annual Summit and Conference, APSIPA 2020, Auckland, New Zealand, December 7-10, 2020*, pages 514–519. IEEE, 2020.

[MÏ5] Meinard Müller. *Fundamentals of Music Processing: Audio, Analysis, Algorithms, Applications*. Springer, 2015.

[MAB+19] Meinard Müller, Andreas Arzt, Stefan Balke, Matthias Dorfer, and Gerhard Widmer. Cross-modal music retrieval and applications: An overview of key methodologies. *IEEE Signal Process. Mag.*, 36(1):52–62, 2019.

[Mac24] Knowing Machines. Models All The Way Down, 2024.

[MAP+20] Masoud Mansoury, Himan Abdollahpouri, Mykola Pechenizkiy, Bamshad Mobasher, and Robin Burke. Feedback loop and bias amplification in recommender systems. In *Proceedings of the 29th ACM international conference on information & knowledge management*, pages 2145–2148, 2020.

[MB03] Tom McCourt and Patrick Burkart. When creators, corporations and consumers collide: Napster and the development of on-line music distribution. *Media, Culture & Society*, 25(3):333–350, 2003.

[MB23] Ted Moore and Jean Brazeau. Serge modular archive instrument (smai): Bridging skeuomorphic & machine learning enabled interfaces. In *Proceedings of the International Conference on New Interfaces for Musical Expression*, 2023.

[MBAB22] Kevin Meng, David Bau, Alex J Andonian, and Yonatan Belinkov. Locating and editing factual associations in GPT. In *Advances in Neural Information Processing Systems*, 2022.

[MBB+23] Fabio Morreale, Elham Bahmanteymouri, Brent Burmester, Andrew Chen, and Michelle Thorp. The unwitting labourer: extracting humanness in AI training. *AI & SOCIETY*, May 2023.

[MBM+20] Margaret Mitchell, Dylan Baker, Nyalleng Moorosi, Emily Denton, Ben Hutchinson, Alex Hanna, Timnit Gebru, and Jamie Morgenstern. Diversity and inclusion metrics in subset selection. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*, pages 117–123, 2020.

[MBQF21] Ilaria Manco, Emmanouil Benetos, Elio Quinton, and György Fazekas. Muscaps: Generating captions for music audio. In *International Joint Conference on Neural Networks, IJCNN 2021, Shenzhen, China, July 18-22, 2021*, pages 1–8. IEEE, 2021.

[MBQF22a] Ilaria Manco, Emmanouil Benetos, Elio Quinton, and George Fazekas. Contrastive audio-language learning for music. In Preeti Rao, Hema A. Murthy, Ajay Srinivasamurthy, Rachel M. Bittner, Rafael Caro Repetto, Masataka Goto, Xavier Serra, and Marius Miron, editors, *Proceedings of the 23rd International Society for Music Information Retrieval Conference, ISMIR 2022, Bengaluru, India, December 4-8, 2022*, pages 640–649, 2022.

[MBQF22b] Ilaria Manco, Emmanouil Benetos, Elio Quinton, and György Fazekas. Learning music audio representations via weak language supervision. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2022)*, pages 456–460. IEEE, 2022.

[MCA24] Tamsin Mains, Victoria Clarke, and Luke Annesley. "Music therapy is the very definition of white privilege": Music therapists' perspectives on race and class in UK music therapy.

*Approaches: An Interdisciplinary Journal of Music Therapy*, January 2024.

[MCMP20] Celia Moreno-Morales, Raul Calero, Pedro Moreno-Morales, and Cristina Pintado. Music Therapy in the Treatment of Dementia: A Systematic Review and Meta-Analysis. *Frontiers in Medicine*, 7, May 2020.

[MDH+24] Matthew C McCallum, Matthew EP Davies, Florian Henkel, Jaehun Kim, and Samuel E Sandberg. On the effect of data-augmentation on local embedding properties in the contrastive learning of music audio representations. *arXiv preprint arXiv:2401.08889*, 2024.

[MDS20] Stylianos I. Mimilakis, Konstantinos Drossos, and Gerald Schuller. Unsupervised interpretable representation learning for singing voice separation. In *28th European Signal Processing Conference, EUSIPCO 2020, Amsterdam, Netherlands, January 18-21, 2021*, pages 1412–1416. IEEE, 2020.

[MEHS21] Gautam Mittal, Jesse Engel, Curtis Hawthorne, and Ian Simon. Symbolic music generation with diffusion models. *arXiv preprint arXiv:2103.16091*, 2021.

[MGG+23] Jan Melechovsky, Zixun Guo, Deepanway Ghosal, Navonil Majumder, Dorien Herremans, and Soujanya Poria. MusTango: Toward controllable text-to-music generation. *arXiv preprint arXiv:2311.08355*, 2023.

[Mil21] Boaz Miller. Is technology value-neutral? *Science, Technology, & Human Values*, 46(1):53–80, 2021.

[Mir21] Eduardo Reck Miranda. *Handbook of artificial intelligence for music*. Springer, 2021.

[MJ12] Meinard Müller and Nanzhu Jiang. A scape plot representation for visualizing repetitive structures of music recordings. In *Proc. ISMIR*, 2012.

[MJXZ23] Lejun Min, Junyan Jiang, Gus Xia, and Jingwei Zhao. Polyffusion: A Diffusion Model for Polyphonic Score Generation with Internal and External Controls. In *Proceeding of the 24th International Society on Music Information Retrieval (ISMIR)*, Milan, Italy, 2023.

[MK22] Anna Maria Matziorinis and Stefan Koelsch. The promise of music therapy for Alzheimer's disease: A review. *Annals of the New York Academy of Sciences*, 1516(1):11–17, 2022.

[MKG+16] Soroush Mehri, Kundan Kumar, Ishaan Gulrajani, Rithesh Kumar, Shubham Jain, Jose Sotelo, Aaron Courville, and Yoshua Bengio. Samplernn: An unconditional end-to-end neural audio generation model. *arXiv preprint arXiv:1612.07837*, 2016.

[MKKW19] Deirdre K Mulligan, Joshua A Kroll, Nitin Kohli, and Richmond Y Wong. This thing called fairness: Disciplinary confusion realizing a value in technology. *Proceedings of the ACM on Human-Computer Interaction*, 3(CSCW):1–36, 2019.

[MKO+22] Matthew C McCallum, Filip Korzeniowski, Sergio Oramas, Fabien Gouyon, and Andreas F Ehmann. Supervised and unsupervised learning of audio representations for music understanding. *arXiv preprint arXiv:2210.03799*, 2022.

[MLL+22] Margaret Mitchell, Alexandra Sasha Luccioni, Nathan Lambert, Marissa Gerchick, Angelina McMillan-Major, Ezinwanne Ozoani, Nazneen Rajani, Tristan Thrush, Yacine Jernite, and Douwe Kiela. Measuring data. *arXiv preprint arXiv:2212.05129*, 2022.

[MLPW22] Xinhao Mei, Xubo Liu, Mark D Plumbley, and Wenwu Wang. Automated audio captioning: an overview of recent progress and new challenges. *EURASIP journal on audio, speech, and music processing*, 2022(1):1–18, 2022.

[MMB+23] Ambuj Mehrish, Navonil Majumder, Rishabh Bharadwaj, Rada Mihalcea, and Soujanya Poria. A review of deep learning techniques for speech processing. *Information Fusion*, page 101869, 2023.

[MMKI19] Takeru Miyato, Shin-ichi Maeda, Masanori Koyama, and Shin Ishii. Virtual adversarial training: A regularization method for supervised and semi-supervised learning. *IEEE Trans. Pattern Anal. Mach. Intell.*, 41(8):1979–1993, 2019.

[MMR+20] Foad Moradi, Hiwa Mohammadi, Mohammad Rezaei, Payam Sariaslani, Nazanin Razazian, Habibolah Khazaie, and Hojjat Adeli. A Novel Method for Sleep-Stage Classification Based on Sonification of Sleep Electroencephalogram Signals Using Wavelet Transform and Recurrent Neural Network. *European Neurology*, 83(5):468–486, 2020.

[Mog16] Olof Mogren. C-rnn-gan: Continuous recurrent neural networks with adversarial training. *arXiv preprint arXiv:1611.09904*, 2016.

[Mor21] Fabio Morreale. Where Does the Buck Stop? Ethical and Political Issues with AI in Music Creation. *Transactions of the International Society for Music Information Retrieval*, 4(1):105–113, July 2021.

[MP20] Gabriel Meseguer-Brocal and Geoffroy Peeters. Content based singing voice source separation via strong conditioning using aligned phonemes. In *Proceedings of the 21th International Society for Music Information Retrieval Conference, ISMIR 2020, Montreal, Canada, October 11-16, 2020*, pages 819–827, 2020.

[MP22] Milagros Miceli and Julian Posada. The data-production dispositif. *Proceedings of the ACM on Human-Computer Interaction*, 6(CSCW2):1–37, 2022.

[MRB+24] Niklas Muennighoff, Alexander Rush, Boaz Barak, Teven Le Scao, Nouamane Tazi, Aleksandra Piktus, Sampo Pyysalo, Thomas Wolf, and Colin A Raffel. Scaling data-constrained language models. *Advances in Neural Information Processing Systems*, 36, 2024.

[MRW19] Brent Mittelstadt, Chris Russell, and Sandra Wachter. Explaining explanations in ai. In *Proceedings of the conference on fairness, accountability, and transparency*, pages 279–288, 2019.

[MSRNDB14] Matt McVicar, Raúl Santos-Rodríguez, Yizhao Ni, and Tijl De Bie. Automatic chord estimation from audio: A review of the state of the art. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 22(2):556–575, 2014.

[MSSR23] Daniel McKee, Justin Salamon, Josef Sivic, and Bryan C. Russell. Language-guided music recommendation for video via prompt analogies. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pages 14784–14793. IEEE, 2023.

[MSW23] Fabio Morreale, Megha Sharma, and I-Chieh Wei. Data collection in music generation training sets: A critical analysis. In Augusto Sarti, Fabio Antonacci, Mark Sandler, Paolo Bestagini, Simon Dixon, Beici Liang, Gaël Richard, and Johan Pauwels, editors, *Proceedings of the 24th International Society for Music Information Retrieval Conference, ISMIR 2023, Milan, Italy, November 5-9, 2023*, pages 37–46, 2023.

[MTBB15] Emilio Molina, Lorenzo J. Tardón, Ana M. Barbancho, and Isabel Barbancho. Sipth: Singing transcription based on hysteresis defined on the pitch-time curve. *IEEE ACM Trans. Audio Speech Lang. Process.*, 23(2):252–263, 2015.

[Mül15] Meinard Müller. *Fundamentals of music processing: Audio, analysis, algorithms, applications*, volume 5. Springer, 2015.

[MWD+23a] I Manco, B Weck, S Doh, M Won, Y Zhang, D Bogdanov, Y Wu, K Chen, P Tovstogan, E Benetos, et al. The song describer dataset: a corpus of audio captions for music-and-language evaluation, 2023.

[MWD+23b] Ilaria Manco, Benno Weck, Seungheon Doh, Minz Won, Yixiao Zhang, Dmitry Bogdanov, Yusong Wu, Ke Chen, Philip Tovstogan, Emmanouil Benetos, Elio Quinton, György Fazekas, and Juhan Nam. The song describer dataset: A corpus of audio captions for music-and-language evaluation. *CoRR*, abs/2311.10057, 2023.

[MYL+23] Yinghao Ma, Ruibin Yuan, Yizhi Li, Ge Zhang, Chenghua Lin, Xingran Chen, Anton Ragni, Hanzhi Yin, Emmanouil Benetos, Norbert Gyenge, Ruibo Liu, Gus Xia, Roger B. Dannenberg, Yike Guo, and Jie Fu. On the effectiveness of speech self-supervised learning for music. In *Proceedings of the 24th International Society for Music Information Retrieval Conference, ISMIR 2023, Milan, Italy, November 5-9, 2023*, pages 457–465, 2023.

[MYO16] Masanori Morise, Fumiya Yokomori, and Kenji Ozawa. WORLD: A vocoder-based high-quality speech synthesis system for real-time applications. *IEICE Trans. Inf. Syst.*, 99-D(7):1877–1884, 2016.

[MZT+23] Ziyang Ma, Zhisheng Zheng, Changli Tang, Yujin Wang, and Xie Chen. MT4SSL: Boosting self-supervised speech representation learning by integrating multiple targets. In *Proc. Interspeech*, 2023.

[MZY+23] Ziyang Ma, Zhisheng Zheng, Guanrou Yang, Yu Wang, Chao Zhang, and Xie Chen. Pushing the limits of unsupervised unit discovery for ssl speech representation. In *Proc. Interspeech*, 2023.

[NBG+16] Joseph W. Newbold, Nadia Bianchi-Berthouze, Nicolas E. Gold, Ana Tajadura-Jiménez, and Amanda CdC Williams. Musically Informed Sonification for Chronic Pain Rehabil-

itation: Facilitating Progress & Avoiding Over-Doing. In *Proceedings of the 2016 CHI Conference on Human Factors in Computing Systems*, CHI '16, pages 5698–5703, New York, NY, USA, May 2016. Association for Computing Machinery.

[NdSC+24] Ingrid Bruno Nunes, Maíra Araújo de Santana, Nicole Charron, Hyngrid Souza e Silva, Caylane Mayssa de Lima Simões, Camila Lins, Ana Beatriz de Souza Sampaio, Arthur Moreira Nogueira de Melo, Thailson Caetano Valdeci da Silva, Camila Tiodista, et al. Automatic identification of preferred music genres: an exploratory machine learning approach to support personalized music therapy. *Multimedia Tools and Applications*, pages 1–17, 2024.

[Ner20] Shahan Nercessian. Zero-shot singing voice conversion. In *Proceedings of the 21th International Society for Music Information Retrieval Conference, ISMIR 2020, Montreal, Canada, October 11-16, 2020*, pages 70–76, 2020.

[NJR+23] Elvis Nunez, Yanzi Jin, Mohammad Rastegari, Sachin Mehta, and Maxwell Horton. Diffusion models as masked audio-video learners. *arXiv preprint arXiv:2310.03937*, 2023.

[NJW+23] Ziqian Ning, Yuepeng Jiang, Zhichao Wang, Bin Zhang, and Lei Xie. Vits-based singing voice conversion leveraging whisper and multi-scale F0 modeling. In *IEEE Automatic Speech Recognition and Understanding Workshop, ASRU 2023, Taipei, Taiwan, December 16-20, 2023*, pages 1–8. IEEE, 2023.

[NKC+23] Tu Anh Nguyen, Eugene Kharitonov, Jade Copet, Yossi Adi, Wei-Ning Hsu, Ali Elkahky, Paden Tomasello, Robin Algayres, Benoit Sagot, Abdelrahman Mohamed, et al. Generative spoken dialogue language modeling. *Transactions of the Association for Computational Linguistics*, 11:250–266, 2023.

[NKKK15] Takashi Nose, Misa Kanemoto, Tomoki Koriyama, and Takao Kobayashi. Hmm-based expressive singing voice synthesis with singing style control and robust pitch modeling. *Computer Speech & Language*, 34(1):308–322, 2015.

[NLR20] Javier Nistal, Stefan Lattner, and Gaël Richard. DRUMGAN: synthesis of drum sounds with timbral feature conditioning using generative adversarial networks. In *Proceedings of the 21th International Society for Music Information Retrieval Conference, ISMIR 2020, Montreal, Canada, October 11-16, 2020*, pages 590–597, 2020.

[NML23] Michele Newman, Lidia Morris, and Jin Ha Lee. Human-ai music creation: Understanding the perceptions and experiences of music creators for ethical and productive collaboration. In Augusto Sarti, Fabio Antonacci, Mark Sandler, Paolo Bestagini, Simon Dixon, Beici Liang, Gaël Richard, and Johan Pauwels, editors, *Proceedings of the 24th International Society for Music Information Retrieval Conference, ISMIR 2023, Milan, Italy, November 5-9, 2023*, pages 80–88, 2023.

[NN03] Han-Wen Nienhuys and Jan Nieuwenhuizen. Lilypond, a system for automated music engraving. In *Proceedings of the XIV Colloquium on Musical Informatics (XIV CIM 2003)*, volume 1, pages 167–171. Citeseer, 2003.

[NNF+19] Ryo Nishikimi, Eita Nakamura, Satoru Fukayama, Masataka Goto, and Kazuyoshi Yoshii. Automatic singing transcription based on encoder-decoder recurrent neural networks with a weakly-supervised attention mechanism. In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2019, Brighton, United Kingdom, May 12-17, 2019*, pages 161–165. IEEE, 2019.

[NNGY19] Ryo Nishikimi, Eita Nakamura, Masataka Goto, and Kazuyoshi Yoshii. End-to-end melody note transcription based on a beat-synchronous attention mechanism. In *2019 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics, WASPAA 2019, New Paltz, NY, USA, October 20-23, 2019*, pages 26–30. IEEE, 2019.

[Noc80] Jorge Nocedal. Updating quasi-newton matrices with limited storage. *Mathematics of computation*, 35(151):773–782, 1980.

[NPA+24] Javier Nistal, Marco Pasini, Cyran Aouameur, Maarten Grachten, and Stefan Lattner. Diff-a-riff: Musical accompaniment co-creation via latent diffusion models. *CoRR*, abs/2406.08384, 2024.

[NSS59] A. Newell, J.C. Shaw, and H.A. Simon. Report on a general problem-solving program. In *Proceedings of the International Conference on Information Processing*, pages 256–264, 1959.

[NTO+21] Daisuke Niizumi, Daiki Takeuchi, Yasunori Ohishi, Noboru Harada, and Kunio Kashino. Byol for audio: Self-supervised learning for general-purpose audio representation. In *2021 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8. IEEE, 2021.

[NTO+22a] Daisuke Niizumi, Daiki Takeuchi, Yasunori Ohishi, Noboru Harada, and Kunio Kashino. Masked Spectrogram Modeling using Masked Autoencoders for Learning General-purpose Audio Representation. In *Proceedings of Machine Learning Research*, Baltimore, MD, USA, 2022.

[NTO+22b] Daisuke Niizumi, Daiki Takeuchi, Yasunori Ohishi, Noboru Harada, and Kunio Kashino. Masked spectrogram modeling using masked autoencoders for learning general-purpose audio representation. In *HEAR: Holistic Evaluation of Audio Representations*, pages 1–24. PMLR, 2022.

[NTO+24] Daisuke Niizumi, Daiki Takeuchi, Yasunori Ohishi, Noboru Harada, and Kunio Kashino. Masked Modeling Duo: Towards a Universal Audio Pre-Training Framework. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 32:2391–2406, 2024.

[OD23] Brendan O'Connor and Simon Dixon. A comparative analysis of latent regressor losses for singing voice conversion. In *20th Sound and Music Computing Conference (SMC 2023), Stockholm, Sweden, June, 2023*, pages 289–295, 2023.

[OEAL+16] Sergio Oramas, Luis Espinosa-Anke, Aonghus Lawlor, et al. Exploring customer reviews for music genre classification and evolutionary studies. *International Society for Music Information Retrieval*, 2016.

[OGW22] Longshen Ou, Xiangming Gu, and Ye Wang. Transfer learning of wav2vec 2.0 for automatic lyric transcription. In *Proceedings of the 23rd International Society for Music Information Retrieval Conference, ISMIR 2022, Bengaluru, India, December 4-8, 2022*, pages 891–899, 2022.

[OLS03] Nicola Orio, Serge Lemouton, and Diemo Schwarz. Score following: State of the art and new developments. *New Interfaces for Musical Expression (NIME)*, 2003.

[OMW23] Longshen Ou, Xichu Ma, and Ye Wang. Loaf-m2l: Joint learning of wording and formatting for singable melody-to-lyric generation. *arXiv preprint arXiv:2307.02146*, 2023.

[OVK17] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural Discrete Representation Learning. In *Advances in Neural Information Processing Systems*, Long Beach, CA, USA, 2017.

[OWJ+22] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, volume 35, pages 27730–27744, 2022.

[OWP+21] Xuan Ouyang, Shuohuan Wang, Chao Pang, Yu Sun, Hao Tian, Hua Wu, and Haifeng Wang. Ernie-m: Enhanced multilingual representation by aligning cross-lingual semantics with monolingual corpora. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 27–38, 2021.

[PAC+21] Adam Polyak, Yossi Adi, Jade Copet, Eugene Kharitonov, Kushal Lakhotia, Wei-Ning Hsu, Abdelrahman Mohamed, and Emmanuel Dupoux. Speech Resynthesis from Discrete Disentangled Self-Supervised Representations. In *Proc. Interspeech 2021*, pages 3615–3619, 2021.

[PAJB23] Christos Plachouras, Pablo Alonso-Jiménez, and Dmitry Bogdanov. mir_ref: A representation evaluation framework for music information retrieval tasks. In *37th Conference on Neural Information Processing Systems (NeurIPS), Machine Learning for Audio Workshop*, New Orleans, LA, USA, 2023.

[PART13] José Portelo, Alberto Abad, Bhiksha Raj, and Isabel Trancoso. Secure binary embeddings of front-end factor analysis for privacy preserving speaker verification. In *14th Annual Conference of the International Speech Communication Association, INTERSPEECH 2013, Lyon, France, August 25-29, 2013*, pages 2494–2498. ISCA, 2013.

[PBK+15] Alastair Porter, Dmitry Bogdanov, Robert Kaye, Roman Tsukanov, and Xavier Serra. Acousticbrainz: a community platform for gathering music information obtained from audio. In *Müller M, Wiering F, editors. ISMIR 2015. 16th*

*International Society for Music Information Retrieval Conference; 2015 Oct 26-30; Málaga, Spain. Canada: ISMIR; 2015.* International Society for Music Information Retrieval (ISMIR), 2015.

[PCC⁺20] Darius Petermann, Pritish Chandna, Helena Cuesta, Jordi Bonada, and Emilia Gómez. Deep learning based source separation applied to choir ensembles. In *Proceedings of the 21th International Society for Music Information Retrieval Conference, ISMIR 2020, Montreal, Canada, October 11-16, 2020*, pages 733–739, 2020.

[PCCF⁺23] Silvan David Peter, Carlos Eduardo Cancino-Chacón, Francesco Foscarin, Andrew Philip McLeod, Florian Henkel, Emmanouil Karystinaios, and Gerhard Widmer. Automatic note-level score-to-performance alignments in the asap dataset. *Transactions of the International Society for Music Information Retrieval (TISMIR)*, 2023.

[PCCKW23] Silvan David Peter, Carlos Eduardo Cancino-Chacón, Emmanouil Karystinaios, and Gerhard Widmer. Sounding out reconstruction error-based evaluation of generative models of expressive performance. In *Proceedings of the 10th International Conference on Digital Libraries for Musicology*, DLfM '23, page 58–66, New York, NY, USA, 2023. Association for Computing Machinery.

[PCGV23] Lorenzo Porcaro, Carlos Castillo, Emilia Gómez, and João Vinagre. Fairness and diversity in information access systems. *arXiv preprint arXiv:2305.09319*, 2023.

[PCW⁺22] Sravya Popuri, Peng-Jen Chen, Changhan Wang, Juan Pino, Yossi Adi, Jiatao Gu, Wei-Ning Hsu, and Ann Lee. Enhanced Direct Speech-to-Speech Translation Using Self-supervised Pre-training and Data Augmentation. In *Proc. Interspeech 2022*, pages 5195–5199, 2022.

[Pea09] Judea Pearl. Causal inference in statistics: An overview. *Statistics Surveys*, 3(none):96 – 146, 2009.

[PG21] Andrea Poltronieri and Aldo Gangemi. The music note ontology. *CEUR-WS*, 11 2021.

[Plo76] R Plomp. *Aspects of Tone Perception*. Academic Press, 1976.

[PMS22] Genís Plaja-Roglans, Marius Miron, and Xavier Serra. A diffusion-inspired training strategy for singing voice extraction in the waveform domain. In *Proceedings of the 23rd International Society for Music Information Retrieval Conference, ISMIR 2022, Bengaluru, India, December 4-8, 2022*, pages 685–693, 2022.

[PMSS23] Genís Plaja-Roglans, Marius Miron, Adithi Shankar, and Xavier Serra. Carnatic singing voice separation using cold diffusion on training data with bleeding. In *Proceedings of the 24th International Society for Music Information Retrieval Conference, ISMIR 2023, Milan, Italy, November 5-9, 2023*, pages 553–560, 2023.

[PR24] Yogesh Prabhakar Pingle and Lakshmappa K. Ragha. Harmonic Healing and Neural Networks: Enhancing Music Therapy Through AI Integration. In Mohammad Shorif Uddin and Jagdish Chand Bansal, editors, *Proceedings of International Joint Conference on Advances in Computational Intelligence*, pages 567–581, Singapore, 2024. Springer Nature.

[PRF24] Leonardo Pepino, Pablo Riera, and Luciana Ferrer. EnCodecMAE: Leveraging neural codecs for universal audio representation learning, May 2024. arXiv:2309.07391 [cs, eess].

[Pro91] Robert Proctor. *Value-free science?: Purity and power in modern knowledge*. Harvard University Press, 1991.

[PRP21] Laure Prétet, Gaël Richard, and Geoffroy Peeters. Cross-modal music-video recommendation: A study of design choices. In *International Joint Conference on Neural Networks, IJCNN 2021, Shenzhen, China, July 18-22, 2021*, pages 1–9. IEEE, 2021.

[PRWZ02] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, 2002.

[PS16] Matthew Pittman and Kim Sheehan. Amazon's mechanical turk a digital sweatshop? transparency and accountability in crowdsourced online research. *Journal of media ethics*, 31(4):260–262, 2016.

[PVG⁺21] Vadim Popov, Ivan Vovk, Vladimir Gogoryan, Tasnima Sadekova, and Mikhail Kudinov. Grad-tts: A diffusion probabilistic model for text-to-speech. In *International Conference on Machine Learning*, pages 8599–8608. PMLR, 2021.

[PVG⁺22] Charilaos Papaioannou, Ioannis Valiantzas, Theodoros Giannakopoulos, Maximos Kaliakatsos-Papakostas, and Alexandros Potamianos. A dataset for greek traditional and folk music: Lyra. *arXiv preprint arXiv:2211.11479*, 2022.

[PX23] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4195–4205, 2023.

[PZV⁺19] Raghavendra Pappagari, Piotr Zelasko, Jesús Villalba, Yishay Carmiel, and Najim Dehak. Hierarchical transformers for long document classification. In *2019 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 838–844, 2019.

[QBM⁺24] Xingwei Qu, Yuelin Bai, Yinghao Ma, Ziya Zhou, Ka Man Lo, Jiaheng Liu, Ruibin Yuan, Lejun Min, Xueling Liu, Tianyu Zhang, et al. Mupt: A generative symbolic music pretrained transformer. *arXiv preprint arXiv:2404.06393*, 2024.

[QML⁺19] Jiezhong Qiu, Hao Ma, Omer Levy, Scott Wen-tau Yih, Sinong Wang, and Jie Tang. Blockwise self-attention for long document understanding. *arXiv preprint arXiv:1911.02972*, 2019.

[QML⁺20] Jiezhong Qiu, Hao Ma, Omer Levy, Wen-tau Yih, Sinong Wang, and Jie Tang. Blockwise self-attention for long document understanding. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2555–2565, 2020.

[QZC⁺19] Kaizhi Qian, Yang Zhang, Shiyu Chang, Xuesong Yang, and Mark Hasegawa-Johnson. Autovc: Zero-shot voice style transfer with only autoencoder loss. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 5210–5219. PMLR, 2019.

[Rab89] L.R. Rabiner. A tutorial on hidden Markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–286, 1989.

[Raf16] Colin Raffel. *Learning-based methods for comparing sequences, with applications to audio-to-midi alignment and matching*. Columbia University, 2016.

[RAN⁺23] Paul K Rubenstein, Chulayuth Asawaroengchai, Duc Dung Nguyen, Ankur Bapna, Zalán Borsos, Félix de Chaumont Quitry, Peter Chen, Dalia El Badawy, Wei Han, Eugene Kharitonov, et al. Audiopalm: A large language model that can speak and listen. *arXiv preprint arXiv:2306.12925*, 2023.

[RBL⁺21] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10674–10685, 2021.

[RBL⁺22] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022.

[RBR24] Madalina Dana Rucsanda, Alexandra Belibou, and Alexandra Ioana Rucsanda. Exploring the relationship between music, medicine and physics. Why pluralism is necessary in music therapy? In *Proceedings of the 2023 12th International Conference on Software and Information Engineering*, ICSIE '23, pages 78–84, New York, NY, USA, January 2024. Association for Computing Machinery.

[RCT24] Gael Richard, Pierre Chouteau, and Bernardo Torres. A fully differentiable model for unsupervised singing voice separation. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 946–950. IEEE, 2024.

[RDRM18] Sabbir M Rashid, David De Roure, and Deborah L McGuinness. A music theory ontology. In *Proceedings of the 1st International Workshop on Semantic Applications for Audio and Music*, pages 6–14, 2018.

[RER⁺18a] Adam Roberts, Jesse Engel, Colin Raffel, Curtis Hawthorne, and Douglas Eck. A hierarchical latent vector model for learning long-term structure in music. In *International conference on machine learning*, pages 4364–4373. PMLR, 2018.

[RER+18b] Adam Roberts, Jesse Engel, Colin Raffel, Curtis Hawthorne, and Douglas Eck. A Hierarchical Latent Vector Model for Learning Long-Term Structure in Music. In *Proceedings of the 35th International Conference on Machine Learning*, pages 4364–4373. PMLR, July 2018. ISSN: 2640-3498.

[RFB15] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In Nassir Navab, Joachim Hornegger, William M. Wells III, and Alejandro F. Frangi, editors, *Medical Image Computing and Computer-Assisted Intervention - MICCAI 2015 - 18th International Conference Munich, Germany, October 5 - 9, 2015, Proceedings, Part III*, volume 9351 of *Lecture Notes in Computer Science*, pages 234–241. Springer, 2015.

[RGC+21] Jessica Sharmin Rahman, Tom Gedeon, Sabrina Caldwell, Richard Jones, and Zi Jin. Towards Effective Music Therapy for Mental Health Care Using Machine Learning Tools: Human Affective Reasoning and Music Genres. *Journal of Artificial Intelligence and Soft Computing Research*, 11(1):5–20, January 2021.

[RH21] Simon Rouard and Gaëtan Hadjeres. CRASH: raw audio score-based generative modeling for controllable high-resolution drum sound synthesis. In *Proceedings of the 22nd International Society for Music Information Retrieval Conference, ISMIR 2021, Online, November 7-12, 2021*, pages 579–585, 2021.

[RHT+20a] Yi Ren, Jinzheng He, Xu Tan, Tao Qin, Zhou Zhao, and Tie-Yan Liu. Popmag: Pop music accompaniment generation. In *Proceedings of the 28th ACM International Conference on Multimedia*, page 1198–1206. Association for Computing Machinery, 2020.

[RHT+20b] Yi Ren, Chenxu Hu, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. Fastspeech 2: Fast and high-quality end-to-end text to speech. *arXiv preprint arXiv:2006.04558*, 2020.

[RHW86] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Learning internal representations by error propagation. In *Parallel distributed processing: explorations in the microstructure of cognition, vol. 1: foundations*, pages 318–362. MIT Press, Cambridge, MA, USA, January 1986.

[RKH+21a] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.

[RKH+21b] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021.

[RKX+23a] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *International Conference on Machine Learning*, pages 28492–28518. PMLR, 2023.

[RKX+23b] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pages 28492–28518. PMLR, 2023.

[RLW+22] Shuo Ren, Shujie Liu, Yu Wu, Long Zhou, and Furu Wei. Speech pre-training with acoustic piece. *arXiv preprint arXiv:2204.03240*, 2022.

[RMY+23] Ludan Ruan, Yiyang Ma, Huan Yang, Huiguo He, Bei Liu, Jianlong Fu, Nicholas Jing Yuan, Qin Jin, and Baining Guo. Mm-diffusion: Learning multi-modal diffusion models for joint audio and video generation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pages 10219–10228. IEEE, 2023.

[RNS+18] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training, 2018.

[Rol02] Perry Roland. The music encoding initiative (mei). In *Proceedings of the First International Conference on Musical Applications Using XML*, volume 1060, pages 55–59. Citeseer, 2002.

[RPG+21] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine learning*, pages 8821–8831. Pmlr, 2021.

[RPJ+20] Jack W. Rae, Anna Potapenko, Siddhant M. Jayakumar, Chloe Hillier, and Timothy P. Lillicrap. Compressive transformers for long-range sequence modelling. In *International Conference on Learning Representations*, 2020.

[RPT15] Bernardino Romera-Paredes and Philip Torr. An embarrassingly simple approach to zero-shot learning. In *International conference on machine learning*, pages 2152–2161. PMLR, 2015.

[RRT+19] Yi Ren, Yangjun Ruan, Xu Tan, Tao Qin, Sheng Zhao, Zhou Zhao, and Tie-Yan Liu. Fastspeech: Fast, robust and controllable text to speech. *Advances in neural information processing systems*, 32, 2019.

[RSIM21] Imad Rahal, Ryan Strelow, Jeremy Iverson, and Katherine Mendel. Separated feature learning for music composition using memory-based neural networks. In *Proceedings of ISCA 30th International Confer*, volume 77, pages 41–51, 2021.

[RSR+20a] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of machine learning research*, 21(140):1–67, 2020.

[RSR+20b] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.

[RST+24] Machel Reid, Nikolay Savinov, Denis Teplyashin, Dmitry Lepikhin, Timothy Lillicrap, Jean-baptiste Alayrac, Radu Soricut, Angeliki Lazaridou, Orhan Firat, Julian Schrittwieser, et al. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. *arXiv preprint arXiv:2403.05530*, 2024.

[RVdOV19] Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems*, 32, 2019.

[RZP+20] Mirco Ravanelli, Jianyuan Zhong, Santiago Pascual, Pawel Swietojanski, Joao Monteiro, Jan Trmal, and Yoshua Bengio. Multi-task self-supervised learning for robust speech recognition. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6989–6993. IEEE, 2020.

[SAFN24] Shanti Stewart, Kleanthis Avramidis, Tiantian Feng, and Shrikanth Narayanan. Emotion-aligned contrastive learning between images and music. In *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8135–8139, 2024.

[SAL+24] Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.

[Sam23] Pamela Samuelson. "legal challenges to generative ai, part ii: Deliberating on inconclusive ai-generated policy questions. *Communications of the ACM*, 66(11):16–19, November 2023.

[SB21] Janne Spijkervet and John Ashley Burgoyne. Contrastive learning of musical representations. *arXiv preprint arXiv:2103.09410*, 2021.

[SBCA19] Steffen Schneider, Alexei Baevski, Ronan Collobert, and Michael Auli. wav2vec: Unsupervised pre-training for speech recognition. *arXiv preprint arXiv:1904.05862*, 2019.

[SDCS23] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. *arXiv preprint arXiv:2303.01469*, 2023.

[SDE19] Daniel Stoller, Simon Durand, and Sebastian Ewert. End-to-end lyrics alignment for polyphonic music using an audio-to-character recognition model. In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2019, Brighton, United Kingdom, May 12-17, 2019*, pages 181–185. IEEE, 2019.

[SDJM24] Siavash Shams, Sukru Samet Dindar, Xilin Jiang, and Nima Mesgarani. Ssamba: Self-supervised audio representation learning with mamba state space model. *arXiv preprint arXiv:2405.11831*, 2024.

[SDL19] Bidisha Sharma, Rohan Kumar Das, and Haizhou Li. On the importance of audio-source separation for singer identification in polyphonic music. In *20th Annual Conference of the International Speech Communication Association, Interspeech 2019, Graz, Austria, September 15-19, 2019*, pages 2020–2024. ISCA, 2019.

[SDP+23] Yutao Sun, Li Dong, Barun Patra, Shuming Ma, Shaohan Huang, Alon Benhaim, Vishrav Chaudhary, Xia Song, and Furu Wei. A length-extrapolatable transformer. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14590–14604, 2023.

[SDvRJ22] Hubert Siuzdak, Piotr Dura, Pol van Rijn, and Nori Jacoby. Wavthruvec: Latent speech representation as intermediate features for neural speech synthesis. In *23rd Annual Conference of the International Speech Communication Association, Interspeech 2022, Incheon, Korea, September 18-22, 2022*, pages 833–837, 2022.

[SDWMG15] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015.

[SDYD+23] Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. *arXiv preprint arXiv:2302.04761*, 2023.

[SED18a] Daniel Stoller, Sebastian Ewert, and Simon Dixon. Adversarial semi-supervised audio source separation applied to singing voice extraction. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2018, Calgary, AB, Canada, April 15-20, 2018*, pages 2391–2395. IEEE, 2018.

[SED18b] Daniel Stoller, Sebastian Ewert, and Simon Dixon. Jointly detecting and separating singing voice: A multi-task approach. In *Latent Variable Analysis and Signal Separation - 14th International Conference, LVA/ICA 2018, Guildford, UK, July 2-5, 2018, Proceedings*, volume 10891 of *Lecture Notes in Computer Science*, pages 329–339. Springer, 2018.

[Ser11] Xavier Serra. A multicultural approach in music information research. In *Klapuri A, Leider C, editors. ISMIR 2011: Proceedings of the 12th International Society for Music Information Retrieval Conference; 2011 October 24-28; Miami, Florida (USA). Miami: University of Miami; 2011.* International Society for Music Information Retrieval (ISMIR), 2011.

[Ser13] Xavier Serra. Exploiting domain knowledge in music information research. In *SMAC & SMC*, Stockholm (Sweden), 2013.

[SF20] Elona Shatri and György Fazekas. Optical music recognition: State of the art and major challenges. In Rama Gottfried, Georg Hajdu, Jacob Sello, Alessandro Anatrini, and John MacCallum, editors, *Proceedings of the International Conference on Technologies for Music Notation and Representation – TENOR'20/21*, pages 175–184, Hamburg, Germany, 2020. Hamburg University for Music and Theater.

[SF21] Elona Shatri and György Fazekas. Doremi: First glance at a universal omr dataset. *arXiv preprint arXiv:2107.07786*, 2021.

[SG21] Axel Sauer and Andreas Geiger. Counterfactual generative networks. In *International Conference on Learning Representations*, 2021.

[SGER14] Justin Salamon, Emilia Gómez, Daniel PW Ellis, and Gaël Richard. Melody extraction from polyphonic music signals: Approaches, applications, and challenges. *IEEE Signal Processing Magazine*, 31(2):118–134, 2014.

[SGH10] Joan Serra, Emilia Gómez, and Perfecto Herrera. Audio cover song identification and similarity: background, approaches, evaluation, and beyond. *Advances in music information retrieval*, pages 307–332, 2010.

[SGS+18] David Snyder, Daniel Garcia-Romero, Gregory Sell, Daniel Povey, and Sanjeev Khudanpur. X-vectors: Robust DNN embeddings for speaker recognition. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2018, Calgary, AB, Canada, April 15-20, 2018*, pages 5329–5333. IEEE, 2018.

[SGU+14] Markus Schedl, Emilia Gómez, Julián Urbano, et al. Music information retrieval: Recent developments and applications. *Foundations and Trends® in Information Retrieval*, 8(2-3):127–261, 2014.

[SGZ21] Aaqib Saeed, David Grangier, and Neil Zeghidour. Contrastive learning of general-purpose audio representations. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3875–3879. IEEE, 2021.

[SH22] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In *International Conference on Learning Representations*, 2022.

[SHB16] Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with subword units. In *ACL (1)*. The Association for Computer Linguistics, 2016.

[SHC+22] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(4):4713–4726, 2022.

[SHD21] Morgan Klaus Scheuerman, Alex Hanna, and Emily Denton. Do datasets have politics? disciplinary values in computer vision dataset development. *Proceedings of the ACM on Human-Computer Interaction*, 5(CSCW2):1–37, 2021.

[She64] Roger N Shepard. Circularity in judgments of relative pitch. *The Journal of the Acoustical Society of America*, 36(12):2346–2353, 1964.

[She82] Roger N Shepard. Geometrical approximations to the structure of musical pitch. *Psychological Review*, 89(4):305–333, 1982.

[SHLM22] Bowen Shi, Wei-Ning Hsu, Kushal Lakhotia, and Abdelrahman Mohamed. Learning Audio-Visual Speech Representation by Masked Multimodal Cluster Prediction. In *International Conference on Learning Representations*, Online, 2022.

[SIBT+19a] Bob L. T. Sturm, Maria Iglesias, Oded Ben-Tal, Marius Miron, and Emilia Gómez. Artificial Intelligence and Music: Open Questions of Copyright Law and Engineering Praxis. *Arts*, 8(3):115, September 2019.

[SIBT+19b] Bob LT Sturm, Maria Iglesias, Oded Ben-Tal, Marius Miron, and Emilia Gómez. Artificial intelligence and music: open questions of copyright law and engineering praxis. In *Arts*, volume 8(3), page 115. MDPI, 2019.

[SJS23a] Flavio Schneider, Zhijing Jin, and Bernhard Schölkopf. Moûsai: Text-to-music generation with long-context latent diffusion. *arXiv preprint arXiv:2301.11757*, 2023.

[SJS23b] Flavio Schneider, Zhijing Jin, and Bernhard Schölkopf. Moûsai: Text-to-music generation with long-context latent diffusion. *arXiv e-prints*, pages arXiv–2301, 2023.

[SKB+24] Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, et al. Dolma: An open corpus of three trillion tokens for language model pretraining research. *arXiv preprint arXiv:2402.00159*, 2024.

[SKM+21] Daniel Schneider, Nikolaus Korfhage, Markus Mühling, Peter Lüttig, and Bernd Freisleben. Automatic transcription of organ tablature music notation with deep neural networks. *Trans. Int. Soc. Music. Inf. Retr.*, 4(1):14–28, 2021.

[SKP+11] L. Schiebinger, I. Klinge, H. Paik, I. Sanchez, M. Schraudner, and M. Stefanick. Gendered innovations in science, health & medicine, engineering, and environment, 2011. Available at http://genderedinnovations.stanford.edu (Retrieved 14/06/2024).

[SL24] Tianyu Li Shijia Liao. Fish speech v1. https://github.com/fishaudio/fish-speech, 2024.

[SLH+23] Kun Su, Judith Yue Li, Qingqing Huang, Dima Kuzmin, Joonseok Lee, Chris Donahue, Fei Sha, Aren Jansen, Yu Wang, Mauro Verzetti, and Timo I. Denk. V2meow: Meowing to the visual beat via music generation. *CoRR*, abs/2305.06594, 2023.

[SLS20a] Kun Su, Xiulong Liu, and Eli Shlizerman. Audeo: Audio generation for a silent performance video. In Hugo Larochelle, Marc'Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Infor-*

[SLS20b] Kun Su, Xiulong Liu, and Eli Shlizerman. Multi-instrumentalist net: Unsupervised generation of music from body movements. *CoRR*, abs/2012.03478, 2020.

[SLS21] Kun Su, Xiulong Liu, and Eli Shlizerman. How does it sound? generation of rhythmic soundtracks for human movement videos. In *Conf. Neural Inf. Process. Syst*, volume 35, pages 0–10, 2021.

[SMB+13] Xavier Serra, Michela Magas, Emmanouil Benetos, Magdalena Chudy, Simon Dixon, Arthur Flexer, Emilia Gómez Gutiérrez, Fabien Gouyon, Herrera Boyer, Sergi Jordà Puig, et al. *Roadmap for music information research*. The MIReS Consortium, 2013.

[SME21] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021.

[Soh16] Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. *Advances in neural information processing systems*, 29, 2016.

[SPD+20] Igor André Pegoraro Santana, Fabio Pinhelli, Juliano Donini, Leonardo Catharin, Rafael Biazus Mangolin, Valéria Delisandra Feltrim, Marcos Aurélio Domingues, et al. Music4all: A new music database and its applications. In *2020 International Conference on Systems, Signals and Image Processing (IWSSIP)*, pages 399–404. IEEE, 2020.

[SRD14] Pierre Saurel, Francis Rousseaux, and Marc Danger. On the changing regulations of privacy and personal information in music information retrieval. In *15th International Society for Music Information Retrieval (ISMIR)*, 2014.

[SRK+23] Kilian Schulze-Forster, Gaël Richard, Liam Kelley, Clement S. J. Doire, and Roland Badeau. Unsupervised music source separation using differentiable parametric source models. *IEEE ACM Trans. Audio Speech Lang. Process.*, 31:1276–1289, 2023.

[SS15] Swathi Swaminathan and E. Glenn Schellenberg. Current Emotion Research in Music Psychology. *Emotion Review*, 7(2):189–197, April 2015.

[SS18] Jiyoung Son and Yongtae Shin. Music lyrics summarization method using textrank algorithm. *Journal of Korea Multimedia Society*, 21(1):45–50, 2018.

[SSBTK16] Bob L Sturm, Joao Felipe Santos, Oded Ben-Tal, and Iryna Korshunova. Music transcription modelling and composition using deep learning. *arXiv preprint arXiv:1604.08723*, 2016.

[SSDK+21] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021.

[SSM+19] Kai Siedenburg, Charalampos Saitis, Stephen McAdams, Arthur N Popper, and Richard R Fay. *Timbre: Acoustics, Perception, and Cognition*. Springer, Cham, 2019.

[SST+23a] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving ai tasks with chatgpt and its friends in huggingface. In *NeurIPS 2023*, 2023.

[SST+23b] Yongliang Shen, Kaitao Song, Xu Tan, Wenqi Zhang, Kan Ren, Siyu Yuan, Weiming Lu, Dongsheng Li, and Yueting Zhuang. Taskbench: Benchmarking large language models for task automation. *arXiv preprint arXiv:2311.18760*, 2023.

[SSZ+22] Yingjie Song, Wei Song, Wei Zhang, Zhengchen Zhang, Dan Zeng, Zhi Liu, and Yang Yu. Singing voice synthesis with vibrato modeling and latent energy representation. In *2022 IEEE 24th International Workshop on Multimedia Signal Processing (MMSP)*, pages 1–6. IEEE, 2022.

[Sta94] Russel A. Stamets. Ain't Nothin' Like the Real Thing, Baby : The Right of Publicity and the Singing Voice. *Federal Communications Law Journal*, 46(2):347–372, 1994.

[SUV18] Peter Shaw, Jakob Uszkoreit, and Ashish Vaswani. Self-attention with relative position representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*. Association for Computational Linguistics, 2018.

[SV24] Catherine Stinson and Sofie Vlaad. A feeling for the algorithm: Diversity, expertise, and artificial intelligence. *Big Data & Society*, 11(1):20539517231224247, 2024.

[SVN37] Stanley Smith Stevens, John Volkmann, and Edwin Broomell Newman. A scale for the measurement of the psychological magnitude pitch. *The journal of the acoustical society of america*, 8(3):185–190, 1937.

[SVRS22] Dídac Surís, Carl Vondrick, Bryan Russell, and Justin Salamon. It's time for artistic correspondence in music and video. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 10554–10564. IEEE, 2022.

[SWR+21] Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Arun Raja, Manan Dey, et al. Multitask prompted training enables zero-shot task generalization. In *International Conference on Learning Representations*, 2021.

[SYKL21] Berrak Sisman, Junichi Yamagishi, Simon King, and Haizhou Li. An overview of voice conversion and its challenges: From statistical modeling to deep learning. *IEEE ACM Trans. Audio Speech Lang. Process.*, 29:132–157, 2021.

[SYZ+24] Jingjing Sun, Jingyi Yang, Guyue Zhou, Yucheng Jin, and Jiangtao Gong. Understanding Human-AI Collaboration in Music Therapy Through Co-Design with Therapists, February 2024.

[TDFH+22] Romal Thoppilan, Daniel De Freitas, Jamie Hall, Noam Shazeer, Apoorv Kulshreshtha, Heng-Tze Cheng, Alicia Jin, Taylor Bos, Leslie Baker, Yu Du, et al. Lamda: Language models for dialog applications. *arXiv preprint arXiv:2201.08239*, 2022.

[Teaa] AgentGPT Team. Agentgpt: Assemble, configure, and deploy autonomous ai agents in your browser. https://github.com/reworkd/AgentGPT. Accessed: 2024-01-21.

[Teab] AutoGPT Team. Autogpt: build & use ai agents. https://github.com/Significant-Gravitas/AutoGPT. Accessed: 2024-01-21.

[Tea23] MosaicML NLP Team. Introducing mpt-7b: A new standard for open-source, commercially usable llms, 2023. Accessed: 2023-05-05.

[Tea24] SVC Develop Team. so-vits-svc. https://github.com/svc-develop-team/so-vits-svc, 2024. Accessed: 2024-07-09.

[TEE+19] Cara Tannenbaum, Robert P Ellis, Friederike Eyssel, James Zou, and Londa Schiebinger. Sex and gender analysis improves science and engineering. *Nature*, 575(7781):137–146, 2019.

[Ter20] Petros Terzis. Onward for the freedom of others: marching beyond the ai ethics. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, pages 220–229, 2020.

[TH20] Hao Hao Tan and Dorien Herremans. Music FaderNets: Controllable Music Generation Based On High-Level Features via Low-Level Feature Modelling. In *Proceedings of the International Society for Music Information Retrieval (ISMIR) Conference*, Online, 2020.

[THDL23] John Thickstun, David Hall, Chris Donahue, and Percy Liang. Anticipatory music transformer. *arXiv preprint arXiv:2306.08620*, 2023.

[THED+22] Chad M Topaz, Jude Higdon, Avriel Epps-Darling, Ethan Siau, Harper Kerkhoff, Shivani Mendiratta, and Eric Young. Race-and gender-based under-representation of creative contributors: art, fashion, film, and music. *Humanities and Social Sciences Communications*, 9(1):1–11, 2022.

[THZY20] Qishou Tang, Zhaohui Huang, Huan Zhou, and Peijie Ye. Effects of music therapy on depression: A meta-analysis of randomized controlled trials. *PLOS ONE*, 15(11):e0240862, November 2020.

[TKC+22] Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. Galactica: A large language model for science. *arXiv preprint arXiv:2211.09085*, 2022.

[TLI+23a] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

[TLI+23b] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al.

Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

[TLR23] Bernardo Torres, Stefan Lattner, and Gaël Richard. Singer identity representation learning using self-supervised techniques. In *Proceedings of the 24th International Society for Music Information Retrieval Conference, ISMIR 2023, Milan, Italy, November 5-9, 2023*, pages 448–456, 2023.

[TLY+24] Zeyue Tian, Zhaoyang Liu, Ruibin Yuan, Jiahao Pan, Xiaoqiang Huang, Qifeng Liu, Xu Tan, Qifeng Chen, Wei Xue, and Yike Guo. Vidmuse: A simple video-to-music generation framework with long-short-term modeling. *arXiv preprint arXiv:2406.04321*, 2024.

[TMS+23] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton-Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurélien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models. *CoRR*, abs/2307.09288, 2023.

[TRH23] Ha Thi Phuong Thao, Gemma Roig, and Dorien Herremans. Emomv: Affective music-video correspondence learning datasets for classification and retrieval. *Information Fusion*, 91:64–79, 2023.

[TRHH22] Francesca Trevisan, Patrice Rusconi, Paul Hanna, and Peter Hegarty. Psychologising meritocracy: A historical account of its many guises. *Theory & Psychology*, 32(2):221–242, 2022.

[TSK+22] Joseph Turian, Jordie Shier, Humair Raj Khan, Bhiksha Raj, Björn W Schuller, Christian J Steinmetz, Colin Malloy, George Tzanetakis, Gissel Velarde, Kirk McNally, et al. Hear: Holistic evaluation of audio representations. In *NeurIPS 2021 Competitions and Demonstrations Track*, pages 125–145. PMLR, 2022.

[TTG+24] Andreas Triantafyllopoulos, Iosif Tsangko, Alexander Gebhard, Annamaria Mesaros, Tuomas Virtanen, and Björn Schuller. Computer audition: From task-specific machine learning to foundation models. *arXiv preprint arXiv:2407.15672*, 2024.

[TYS+23] Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun Ma, and Chao Zhang. Salmonn: Towards generic hearing abilities for large language models. *arXiv preprint arXiv:2310.13289*, 2023.

[TYS+24a] Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun MA, and Chao Zhang. SALMONN: Towards generic hearing abilities for large language models. In *The Twelfth International Conference on Learning Representations*, 2024.

[TYS+24b] Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, MA Zejun, and Chao Zhang. Salmonn: Towards generic hearing abilities for large language models. In *The Twelfth International Conference on Learning Representations*, 2024.

[VB+22] Marcel A Vélez Vásquez, John Ashley Burgoyne, et al. Tailed u-net: Multi-scale music representation learning. In *ISMIR*, pages 67–75, 2022.

[VBV22] Sergey Verbitskiy, Vladimir Berikov, and Viacheslav Vyshegorodtsev. Eranns: Efficient residual audio neural networks for audio pattern recognition. *Pattern Recognition Letters*, 161:38–44, 2022.

[VDG19] Gaurav Verma, Eeshan Gunesh Dhekane, and Tanaya Guha. Learning affective correspondence between music and image. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3975–3979. IEEE, 2019.

[VDODZ+16] Aaron Van Den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, Koray Kavukcuoglu, et al. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 12, 2016.

[vdOLV18] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.

[VDOV+17a] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.

[VDOV+17b] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in neural information processing systems*, 30, 2017.

[VHM+20] Andrea Vaglio, Romain Hennequin, Manuel Moussallam, Gaël Richard, and Florence d'Alché-Buc. Multilingual lyrics-to-audio alignment. In *Proceedings of the 21th International Society for Music Information Retrieval Conference, ISMIR 2020, Montreal, Canada, October 11-16, 2020*, pages 512–519, 2020.

[vRBKH23] Dimitri von Rütte, Luca Biggio, Yannic Kilcher, and Thomas Hofmann. FIGARO: Controllable music generation using learned and expert features. In *Proc. ICLR*, 2023.

[VSL+23] Apoorv Vyas, Bowen Shi, Matthew Le, Andros Tjandra, Yi-Chiao Wu, Baishan Guo, Jiemin Zhang, Xinyue Zhang, Robert Adkins, William Ngan, et al. Audiobox: Unified audio generation with natural language prompts. *arXiv preprint arXiv:2312.15821*, 2023.

[VSP+17a] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

[VSP+17b] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

[WBZ+21] Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In *International Conference on Learning Representations*, 2021.

[WCL+24] Haibin Wu, Xuanjun Chen, Yi-Cheng Lin, Kai-wei Chang, Ho-Lam Chung, Alexander H Liu, and Hung-yi Lee. Towards audio language modeling-an overview. *arXiv preprint arXiv:2402.13236*, 2024.

[WCW+23] Chengyi Wang, Sanyuan Chen, Yu Wu, Ziqiang Zhang, Long Zhou, Shujie Liu, Zhuo Chen, Yanqing Liu, Huaming Wang, Jinyu Li, et al. Neural codec language models are zero-shot text to speech synthesizers. *arXiv preprint arXiv:2301.02111*, 2023.

[WCZ+23a] Yusong Wu, Ke Chen, Tianyu Zhang, Yuchen Hui, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. Large-scale contrastive language-audio pretraining with feature fusion and keyword-to-caption augmentation. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.

[WCZ+23b] Yusong Wu, Ke Chen, Tianyu Zhang, Yuchen Hui, Taylor Berg-Kirkpatrick, and Shlomo Dubnov. Large-scale contrastive language-audio pretraining with feature fusion and keyword-to-caption augmentation. In *Proc. ICASSP*, 2023.

[WDWB24] Shih-Lun Wu, Chris Donahue, Shinji Watanabe, and Nicholas J Bryan. Music controlnet: Multiple time-varying controls for music generation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 2024.

[WERA16] Elliot Waite, Douglas Eck, Adam Roberts, and Daniel Abolafia. Project magenta: Generating long-term structure in songs and stories, 2016. https://magenta.tensorflow.org/2016/07/15/lookback-rnn-attention-rnn.

[WFQ+23] Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. Next-gpt: Any-to-any multimodal llm. *arXiv preprint arXiv:2309.05519*, 2023.

[WGI+21] Zhepei Wang, Ritwik Giri, Umut Isik, Jean-Marc Valin, and Arvindh Krishnaswamy. Semi-supervised singing voice separation with noisy self-training. In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2021, Toronto, ON, Canada, June 6-11, 2021*, pages 31–35. IEEE, 2021.

[WH17]   Derick T Wade and Peter W Halligan. The biopsychosocial model of illness: A model whose time has come. *Clinical Rehabilitation*, 31(8):995–1004, August 2017.

[WH19]   Yixue Wang and Emőke-Ágnes Horvát. Gender differences in the global music industry: Evidence from musicbrainz and the echo nest. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 13, pages 517–526, 2019.

[Whi23]   Meredith Whittaker. Origin stories: Plantations, computers, and industrial control. *Logic(s) Magazine*, 19, 2023.

[WHL24a]   Minz Won, Yun-Ning Hung, and Duc Le. A foundation model for music informatics. In *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1226–1230, 2024.

[WHL24b]   Minz Won, Yun-Ning Hung, and Duc Le. A Foundation Model for Music Informatics. In *ICASSP 2024 - 2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1226–1230, Seoul, Korea, April 2024. IEEE.

[WHY+22]   Da-Yi Wu, Wen-Yi Hsiao, Fu-Rong Yang, Oscar Friedman, Warren Jackson, Scott Bruzenak, Yi-Wen Liu, and Yi-Hsuan Yang. Ddsp-based singing vocoders: A new subtractive-based synthesizer and a comprehensive evaluation. *arXiv preprint arXiv:2208.04756*, 2022.

[Wid98]   Gerhard Widmer. Applications of machine learning to music research: Empirical investigations into the phenomenon of musical expression. In Ryszad Michalski, Ivan Bratko, and Miroslav Kubat, editors, *Machine Learning, Data Mining and Knowledge Discovery: Methods and Applications*. Wiley & Sons, Chichester (UK), 1998.

[Wig04]   Tony Wigram. *Improvisation: Methods and techniques for music therapy clinicians, educators, and students*. Jessica Kingsley Publishers, 2004.

[Wil92]   Frank Wilcoxon. Individual comparisons by ranking methods. In *Breakthroughs in Statistics: Methodology and Distribution*, pages 196–202. Springer, 1992.

[WJ23]   Jun-You Wang and Jyh-Shing Roger Jang. Training a singing transcription model using connectionist temporal classification loss and cross-entropy loss. *IEEE ACM Trans. Audio Speech Lang. Process.*, 31:383–396, 2023.

[WJT+23]   Yuancheng Wang, Zeqian Ju, Xu Tan, Lei He, Zhizheng Wu, Jiang Bian, and Sheng Zhao. Audit: Audio editing by following instructions with latent diffusion models. *arXiv preprint arXiv:2304.00830*, 2023.

[WKT+21]   Ho-Hsiang Wu, Chieh-Chi Kao, Qingming Tang, Ming Sun, Brian McFee, Juan Pablo Bello, and Chao Wang. Multi-task self-supervised pre-training for music classification. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 556–560. IEEE, 2021.

[WLL+23]   Jun-You Wang, Chon-In Leong, Yu-Chen Lin, Li Su, and Jyh-Shing Roger Jang. Adapting pretrained speech model for mandarin lyrics transcription and alignment. In *IEEE Automatic Speech Recognition and Understanding Workshop, ASRU 2023, Taipei, Taiwan, December 16-20, 2023*, pages 1–8. IEEE, 2023.

[WLQ+22]   Qingyang Wu, Zhenzhong Lan, Kun Qian, Jing Gu, Alborz Geramifard, and Zhou Yu. Memformer: A memory-augmented transformer for sequence modeling. In *Findings of the Association for Computational Linguistics: AACL-IJCNLP 2022*, pages 308–318, 2022.

[WLS23]   Shangda Wu, Xiaobing Li, and Maosong Sun. 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, 2023.

[WLT+22]   Chao Wang, Zhonghao Li, Benlai Tang, Xiang Yin, Yuan Wan, Yibiao Yu, and Zejun Ma. Towards high-fidelity singing voice conversion with acoustic reference and contrastive predictive coding. In *23rd Annual Conference of the International Speech Communication Association, Interspeech 2022, Incheon, Korea, September 18-22, 2022*, pages 4287–4291. ISCA, 2022.

[WLW+22]   Luyu Wang, Pauline Luc, Yan Wu, Adria Recasens, Lucas Smaira, Andrew Brock, Andrew Jaegle, Jean-Baptiste Alayrac, Sander Dieleman, Joao Carreira, et al. Towards learning universal audio representations. In *ICASSP 2022-*

*2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4593–4597. IEEE, 2022.

[WLYS23]   Shangda Wu, Xiaobing Li, Feng Yu, and Maosong Sun. Tunesformer: Forming irish tunes with control codes by bar patching. *arXiv preprint arXiv:2301.02884*, 3528, 2023.

[WMA+22]   Yizhong Wang, Swaroop Mishra, Pegah Alipoormolabashi, Yeganeh Kordi, Amirreza Mirzaei, Atharva Naik, Arjun Ashok, Arut Selvan Dhanasekaran, Anjana Arunkumar, David Stap, et al. Super-naturalinstructions: Generalization via declarative instructions on 1600+ nlp tasks. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 5085–5109, 2022.

[WMB+24]   Benno Weck, Ilaria Manco, Emmanouil Benetos, Elio Quinton, György Fazekas, and Dmitry Bogdanov. Muchomusic: Evaluating music understanding in multimodal audio-language models. In *Proceedings of the 25th International Society for Music Information Retrieval Conference (ISMIR)*, 2024.

[WMF+24]   Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6):186345, 2024.

[WMG22]   Adrienne Williams, Milagros Miceli, and Timnit Gebru. The exploited labor behind artificial intelligence. *Noema Magazine*, 22, 2022.

[WMX24]   Ziyu Wang, Lejun Min, and Gus Xia. Whole-song hierarchical generation of symbolic music using cascaded diffusion models. *arXiv preprint arXiv:2405.09901*, 2024.

[WMZ+23]   Zihao Wang, Le Ma, Chen Zhang, Bo Han, Yikai Wang, Xinyi Chen, HaoRong Hong, Wenbo Liu, Xinda Wu, and Kejun Zhang. SongDriver2: Real-time Emotion-based Music Arrangement with Soft Transition, May 2023.

[WRK+19]   B. N. Walker, J. M. Rehg, A. Kalra, R. M. Winters, P. Drews, J. Dascalu, E. O. David, and A. Dascalu. Dermoscopy diagnosis of cancerous lesions utilizing dual deep learning algorithms via visual and audio (sonification) outputs: Laboratory and prospective observational studies. *eBioMedicine*, 40:176–183, February 2019.

[WS22]   Shangda Wu and Maosong Sun. Exploring the efficacy of pre-trained checkpoints in text-to-music generation task. *arXiv preprint arXiv:2211.11216*, 2022.

[WSBB20]   Yu Wang, Justin Salamon, Nicholas J Bryan, and Juan Pablo Bello. Few-shot sound event detection. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 81–85. IEEE, 2020.

[WSBB22]   Yu Wang, Daniel Stoller, Rachel M Bittner, and Juan Pablo Bello. Few-shot musical source separation. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 121–125. IEEE, 2022.

[WSC+20]   Yu Wang, Justin Salamon, Mark Brozier Cartwright, Nicholas J. Bryan, and Juan Pablo Bello. Few-shot drum transcription in polyphonic music. In *International Society for Music Information Retrieval Conference*, 2020.

[WSKB22]   Ho-Hsiang Wu, Prem Seetharaman, Kundan Kumar, and Juan Pablo Bello. Wav2clip: Learning robust audio representations from clip. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4563–4567. IEEE, 2022.

[WSRS+17]   Yuxuan Wang, RJ Skerry-Ryan, Daisy Stanton, Yonghui Wu, Ron J Weiss, Navdeep Jaitly, Zongheng Yang, Ying Xiao, Zhifeng Chen, Samy Bengio, et al. Tacotron: Towards end-to-end speech synthesis. *arXiv preprint arXiv:1703.10135*, 2017.

[WST+24]   Yuning Wu, Jiatong Shi, Yuxun Tang, Shan Yang, Qin Jin, et al. Toksing: Singing voice synthesis based on discrete tokens. *arXiv preprint arXiv:2406.08416*, 2024.

[WTY19]   Xin Wang, Shinji Takaki, and Junichi Yamagishi. Neural source-filter waveform models for statistical parametric speech synthesis. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 28:402–415, 2019.

[WTY+23]   Xianke Wang, Bowen Tian, Weiming Yang, Wei Xu, and Wenqing Cheng. Musicyolo: A vision-based framework for automatic singing transcription. *IEEE ACM Trans. Audio Speech Lang. Process.*, 31:229–241, 2023.

[WWS+22]   Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022.

[WWW+22]   Chengyi Wang, Yiming Wang, Yu Wu, Sanyuan Chen, Jinyu Li, Shujie Liu, and Furu Wei. Supervision-guided codebooks for masked prediction in speech pre-training. In *Proc. Interspeech*, 2022.

[WX22]   Shiqi Wei and Gus Xia. Learning long-term music representations via hierarchical contextual constraints. *arXiv preprint arXiv:2202.06180*, 2022.

[WY20]   Shih-Lun Wu and Yi-Hsuan Yang. The jazz transformer on the front line: Exploring the shortcomings of ai-composed music through quantitative measures. In *Proc. ISMIR*, 2020.

[WY23a]   Shih-Lun Wu and Yi-Hsuan Yang. Compose & Embellish: Well-structured piano performance generation via a two-stage approach. In *Proc. ICASSP*, 2023.

[WY23b]   Shih-Lun Wu and Yi-Hsuan Yang. MuseMorphose: Full-Song and Fine-Grained Piano Music Style Transfer With One Transformer VAE. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 31:1953–1967, 2023. Conference Name: IEEE/ACM Transactions on Audio, Speech, and Language Processing.

[WY23c]   Shih-Lun Wu and Yi-Hsuan Yang. MuseMorphose: Full-song and fine-grained piano music style transfer with one transformer vae. *IEEE/ACM T-ASLP*, 2023.

[WYQ+23]   Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv preprint arXiv:2303.04671*, 2023.

[WYTS23a]   Shangda Wu, Dingyao Yu, Xu Tan, and Maosong Sun. Clamp: Contrastive language-music pre-training for cross-modal symbolic music information retrieval. *24th International Society for Music Information Retrieval Conference.*, 2023.

[WYTS23b]   Shangda Wu, Dingyao Yu, Xu Tan, and Maosong Sun. Clamp: Contrastive language-music pre-training for cross-modal symbolic music information retrieval. In Augusto Sarti, Fabio Antonacci, Mark Sandler, Paolo Bestagini, Simon Dixon, Beici Liang, Gaël Richard, and Johan Pauwels, editors, *Proceedings of the 24th International Society for Music Information Retrieval Conference, ISMIR 2023, Milan, Italy, November 5-9, 2023*, pages 157–165, 2023.

[WYW+24]   Shangda Wu, Yue Yang, Zhaowen Wang, Xiaobing Li, and Maosong Sun. Generating chord progression from melody with flexible harmonic rhythm and controllable harmonic density. *EURASIP Journal on Audio, Speech, and Music Processing*, 2024(1):4, 2024.

[WZH22]   Chunhui Wang, Chang Zeng, and Xing He. Xiaoicesing 2: A high-fidelity singing voice synthesizer based on generative adversarial network. *arXiv preprint arXiv:2210.14666*, 2022.

[WZZ+20a]   Ziyu Wang, Yiyi Zhang, Yixiao Zhang, Junyan Jiang, Ruihan Yang, Junbo Zhao, and Gus Xia. Pianotree VAE: Structured representation learning for polyphonic music. In *Proceedings of the 21st International Society for Music Information Retrieval Conference, ISMIR*, 2020.

[WZZ+20b]   Ziyu Wang, Yiyi Zhang, Yixiao Zhang, Junyan Jiang, Ruihan Yang, Junbo Zhao, and Gus Xia. PIANOTREE VAE: Structured Representation Learning for Polyphonic Music. In *Proceedings of the International Society for Music Information Retrieval (ISMIR) Conference*, Online, August 2020.

[XCG+23]   Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, et al. The rise and potential of large language model based agents: A survey. *arXiv preprint arXiv:2309.07864*, 2023.

[XD17]   Guangyu Xia and Roger B Dannenberg. Improvised duet interaction: learning improvisation techniques for automatic accompaniment. In *NIME*, pages 110–114, 2017.

[XGL22]   Gao Xiaoxue, Chitralekha Gupta, and Haizhou Li. Music-robust automatic lyrics transcription of polyphonic music. In *Proc. SMC 2022*, 04 2022.

[XHT+24]   Yazhou Xing, Yingqing He, Zeyue Tian, Xintao Wang, and Qifeng Chen. Seeing and hearing: Open-domain visual-audio generation with diffusion latent aligners. *CoRR*, abs/2402.17723, 2024.

[XLB+22]   Hu Xu, Juncheng Li, Alexei Baevski, Michael Auli, Wojciech Galuba, Florian Metze, Christoph Feichtenhofer, et al. Masked autoencoders that listen. *arXiv preprint arXiv:2207.06405*, 2022.

[XLHL20]   Qizhe Xie, Minh-Thang Luong, Eduard H. Hovy, and Quoc V. Le. Self-training with noisy student improves imagenet classification. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020*, pages 10684–10695. Computer Vision Foundation / IEEE, 2020.

[XOG+22]   Wenhan Xiong, Barlas Oguz, Anchit Gupta, Xilun Chen, Diana Liskovich, Omer Levy, Scott Yih, and Yashar Mehdad. Simple local attentions remain competitive for long-context tasks. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1975–1986, 2022.

[XPD+24]   Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy S Liang, Quoc V Le, Tengyu Ma, and Adams Wei Yu. Doremi: Optimizing data mixtures speeds up language model pretraining. *Advances in Neural Information Processing Systems*, 36, 2024.

[XRV21]   Huang Xie, Okko Räsänen, and Tuomas Virtanen. Zero-shot audio classification with factored linear and nonlinear acoustic-semantic projections. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 326–330. IEEE, 2021.

[XSA17]   Yongqin Xian, Bernt Schiele, and Zeynep Akata. Zero-shot learning-the good, the bad and the ugly. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4582–4591, 2017.

[YCY17]   Li-Chia Yang, Szu-Yu Chou, and Yi-Hsuan Yang. Midinet: A convolutional generative adversarial network for symbolic-domain music generation. *arXiv preprint arXiv:1703.10847*, 2017.

[YGG+19]   Zihao Ye, Qipeng Guo, Quan Gan, Xipeng Qiu, and Zheng Zhang. Bp-transformer: Modelling long-range context via binary partitioning, 2019.

[YHF+21]   Yin-Cheng Yeh, Wen-Yi Hsiao, Satoru Fukayama, Tetsuro Kitahara, Benjamin Genchel, Hao-Min Liu, Hao-Wen Dong, Yian Chen, Terence Leong, and Yi-Hsuan Yang. Automatic melody harmonization with triad chords: A comparative study. *Journal of New Music Research*, 50(1):37–51, 2021.

[YHH+24]   Weijia Yang, Chih-Fang Huang, Hsun-Yi Huang, Zixue Zhang, Wenjun Li, and Chunmei Wang. Research on the Improvement of Children's Attention Through Binaural Beats Music Therapy in the Context of AI Music Generation. In Xiaobing Li, Xiaohong Guan, Yun Tie, Xinran Zhang, and Qingwen Zhou, editors, *Music Intelligence*, pages 19–31, Singapore, 2024. Springer Nature.

[YL20]   Li-Chia Yang and Alexander Lerch. On the evaluation of generative models in music. *Neural Computing and Applications*, 2020.

[YLG+23]   Lijun Yu, José Lezama, Nitesh B Gundavarapu, Luca Versari, Kihyuk Sohn, David Minnen, Yong Cheng, Agrim Gupta, Xiuye Gu, Alexander G Hauptmann, et al. Language model beats diffusion–tokenizer is key to visual generation. *arXiv preprint arXiv:2310.05737*, 2023.

[YLH+23]   Dongchao Yang, Songxiang Liu, Rongjie Huang, Jinchuan Tian, Chao Weng, and Yuexian Zou. Hifi-codec: Group-residual vector quantization for high fidelity audio codec. *arXiv preprint arXiv:2305.02765*, 2023.

[YLK+22]   Jiahui Yu, Xin Li, Jing Yu Koh, Han Zhang, Ruoming Pang, James Qin, Alexander Ku, Yuanzhong Xu, Jason Baldridge, and Yonghui Wu. Vector-quantized image modeling with improved VQGAN. In *International Conference on Learning Representations*, 2022.

[YLL+24]   Zonghan Yang, An Liu, Zijun Liu, Kaiming Liu, Fangzhou Xiong, Yile Wang, Zeyuan Yang, Qingyuan Hu, Xinrui Chen, Zhenhe Zhang, Fuwen Luo, Zhicheng Guo, Peng Li, and Yang Liu. Towards unified alignment between agents, humans, and environment. In *Proceedings of the 41th International Conference on Machine Learning*, 2024.

[YLW+22]   Botao Yu, Peiling Lu, Rui Wang, Wei Hu, Xu Tan, Wei Ye, Shikun Zhang, Tao Qin, and Tie-Yan Liu. Museformer: Transformer with fine-and coarse-grained attention for music

generation. *Advances in Neural Information Processing Systems*, 35:1376–1388, 2022.

[YLW+24a] Ruibin Yuan, Hanfeng Lin, Yi Wang, Zeyue Tian, Shangda Wu, Tianhao Shen, Ge Zhang, Yuhang Wu, Cong Liu, Ziya Zhou, et al. Chatmusician: Understanding and generating music intrinsically with llm. *arXiv preprint arXiv:2402.16153*, 2024.

[YLW+24b] Ruibin Yuan, Hanfeng Lin, Yi Wang, Zeyue Tian, Shangda Wu, 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 B. Dannenberg, Wei Xue, Shiyin Kang, and Yike Guo. Chatmusician: Understanding and generating music intrinsically with LLM. *CoRR*, abs/2402.16153, 2024.

[YLY+24] Zonghan Yang, Peng Li, Ming Yan, Ji Zhang, Fei Huang, and Yang Liu. React meets actre: Autonomous annotations of agent trajectories for contrastive self-training. In *Conference on Language Modeling*, 2024.

[YML+23] Ruibin Yuan, Yinghao Ma, Yizhi Li, Ge Zhang, Xingran Chen, Hanzhi Yin, Le Zhuo, Yiqi Liu, Jiawen Huang, Zeyue Tian, Binyue Deng, Ningzhi Wang, Chenghua Lin, Emmanouil Benetos, Anton Ragni, Norbert Gyenge, Roger B. Dannenberg, Wenhu Chen, Gus Xia, Wei Xue, Si Liu, Shi Wang, Ruibo Liu, Yike Guo, and Jie Fu. MARBLE: music audio representation benchmark for universal evaluation. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023.

[YML+24] Ruibin Yuan, Yinghao Ma, Yizhi Li, Ge Zhang, Xingran Chen, Hanzhi Yin, Yiqi Liu, Jiawen Huang, Zeyue Tian, Binyue Deng, et al. Marble: Music audio representation benchmark for universal evaluation. *Advances in Neural Information Processing Systems*, 36, 2024.

[YMZ+23] Guanrou Yang, Ziyang Ma, Zhisheng Zheng, Yakun Song, Zhikang Niu, and Xie Chen. Fast-hubert: an efficient training framework for self-supervised speech representation learning. In *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, pages 1–7. IEEE, 2023.

[YNZ+23] Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, et al. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. *arXiv preprint arXiv:2311.16502*, 2023.

[YPRF23] Chin-Yun Yu, Emilian Postolache, Emanuele Rodolà, and György Fazekas. Zero-shot duet singing voices separation with diffusion models. *arXiv preprint arXiv:2311.07345*, 2023.

[YQMB19] Chi-Hua Yu, Zhao Qin, Francisco J. Martin-Martinez, and Markus J. Buehler. A Self-Consistent Sonification Method to Translate Amino Acid Sequences into Musical Compositions and Application in Protein Design Using Artificial Intelligence. *ACS Nano*, 13(7):7471–7482, July 2019.

[YRC+22] Zhuoyuan Yao, Shuo Ren, Sanyuan Chen, Ziyang Ma, Pengcheng Guo, and Lei Xie. Tessp: text-enhanced self-supervised speech pre-training. *arXiv preprint arXiv:2211.13443*, 2022.

[YSC+24] Siyu Yuan, Kaitao Song, Jiangjie Chen, Xu Tan, Yongliang Shen, Ren Kan, Dongsheng Li, and Deqing Yang. Easytool: Enhancing llm-based agents with concise tool instruction. *arXiv preprint arXiv:2401.06201*, 2024.

[YSL+23] Dingyao Yu, Kaitao Song, Peiling Lu, Tianyu He, Xu Tan, Wei Ye, Shikun Zhang, and Jiang Bian. Musicagent: An ai agent for music understanding and generation with large language models. *arXiv preprint arXiv:2310.11954*, 2023.

[YSN23] Sangeon Yong, Li Su, and Juhan Nam. A phoneme-informed neural network model for note-level singing transcription. In *IEEE International Conference on Acoustics, Speech and Signal Processing ICASSP 2023, Rhodes Island, Greece, June 4-10, 2023*, pages 1–5. IEEE, 2023.

[YTHT24] Sarthak Yadav, Sergios Theodoridis, Lars Kai Hansen, and Zheng-Hua Tan. Masked Autoencoders with Multi-Window Local-Global Attention Are Better Audio Learners. In *International Conference on Learning Representations*, Vienna, Austria, 2024. arXiv.

[YTT+23] Dongchao Yang, Jinchuan Tian, Xu Tan, Rongjie Huang, Songxiang Liu, Xuankai Chang, Jiatong Shi, Sheng Zhao, Jiang Bian, Xixin Wu, et al. Uniaudio: An audio foundation model toward universal audio generation. *arXiv preprint arXiv:2310.00704*, 2023.

[YV16] Haizi Yu and Lav R Varshney. Towards deep interpretability (mus-rover ii): Learning hierarchical representations of tonal music. In *International Conference on Learning Representations*, 2016.

[YWT23] Reo Yoneyama, Yi-Chiao Wu, and Tomoki Toda. Source-filter hifi-gan: Fast and pitch controllable high-fidelity neural vocoder. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.

[YWW+19] Ruihan Yang, Dingsu Wang, Ziyu Wang, Tianyao Chen, Junyan Jiang, and Gus Xia. Deep Music Analogy Via Latent Representation Disentanglement. In *Proceedings of the International Society for Music Information Retrieval (ISMIR) Conference*, Delft, The Netherlands, 2019.

[YWW+23] Weitao Yuan, Shengbei Wang, Jianming Wang, Masashi Unoki, and Wenwu Wang. Unsupervised deep unfolded representation learning for singing voice separation. *IEEE ACM Trans. Audio Speech Lang. Process.*, 31:3206–3220, 2023.

[YWZ+22] Fan Yang, Kai Wu, Shuyi Zhang, Guannan Jiang, Yong Liu, Feng Zheng, Wei Zhang, Chengjie Wang, and Long Zeng. Class-aware contrastive semi-supervised learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14421–14430, 2022.

[YXK+22] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin Li, Han Zhang, Jason Baldridge, and Yonghui Wu. Scaling autoregressive models for content-rich text-to-image generation. *Transactions on Machine Learning Research*, 2022. Featured Certification.

[YXT+23] Zhen Ye, Wei Xue, Xu Tan, Jie Chen, Qifeng Liu, and Yike Guo. Comospeech: One-step speech and singing voice synthesis via consistency model. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 1831–1839, 2023.

[YYW+23] Dongchao Yang, Jianwei Yu, Helin Wang, Wen Wang, Chao Weng, Yuexian Zou, and Dong Yu. Diffsound: Discrete diffusion model for text-to-sound generation. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 31:1720–1733, 2023.

[YZT+24] Kaixing Yang, Xukun Zhou, Xulong Tang, Ran Diao, Hongyan Liu, Jun He, and Zhaoxin Fan. Beatdance: A beat-based model-agnostic contrastive learning framework for music-dance retrieval. In *Proceedings of the 2024 International Conference on Multimedia Retrieval*, ICMR '24, page 11–19, New York, NY, USA, 2024. Association for Computing Machinery.

[YZWY17] Lantao Yu, Weinan Zhang, Jun Wang, and Yong Yu. Seqgan: Sequence generative adversarial nets with policy gradient. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA*, pages 2852–2858. AAAI Press, 2017.

[ZAPX24] Dorothy Zhao, Jerone TA Andrews, Orestis Papakyriakopoulos, and Alice Xiang. Measuring diversity in datasets. In *Data-centric Machine Learning Research Workshop @ ICLR*, 2024.

[ZBRR23] Yueyue Zhu, Jared Baca, Banafsheh Rekabdar, and Reza Rawassizadeh. A survey of ai music generation tools and models. *arXiv preprint arXiv:2308.12982*, 2023.

[ZCCC+24] Huan Zhang, Shreyan Chowdhury, Carlos Eduardo Cancino-Chacón, Jinhua Liang, Simon Dixon, and Gerhard Widmer. Dexter: Learning and controlling performance expression with diffusion models. *Applied Sciences*, 14(15), 2024.

[ZCS+23] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging LLM-as-a-judge with MT-bench and chatbot arena. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023.

[ZCX+22]   Yongmao Zhang, Jian Cong, Heyang Xue, Lei Xie, Pengcheng Zhu, and Mengxiao Bi. Visinger: Variational inference with adversarial learning for end-to-end singing voice synthesis. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7237–7241. IEEE, 2022.

[ZCYC19]   Yichao Zhou, Wei Chu, Sam Young, and Xin Chen. Bandnet: A neural network-based, multi-instrument beatles-style midi music composition machine. In *Proceedings of the 20th International Society for Music Information Retrieval Conference, ISMIR*, 2019.

[ZDJ+22]   Ge Zhu, Jordan Darefsky, Fei Jiang, Anton Selitskiy, and Zhiyao Duan. Music source separation with generative flow. *IEEE Signal Process. Lett.*, 29:2288–2292, 2022.

[ZDY+24]   Jun Zhan, Junqi Dai, Jiasheng Ye, Yunhua Zhou, Dong Zhang, Zhigeng Liu, Xin Zhang, Ruibin Yuan, Ge Zhang, Linyang Li, et al. Anygpt: Unified multimodal llm with discrete sequence modeling. *arXiv preprint arXiv:2402.12226*, 2024.

[ZGD+20]   Manzil Zaheer, Guru Guruganesh, Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. Big bird: transformers for longer sequences. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NeurIPS 2020, 2020.

[ZGR+18]   Hang Zhao, Chuang Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio Torralba. The sound of pixels. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018.

[Zha03]   Tong Zhang. Automatic singer identification. In *Proceedings of the 2003 IEEE International Conference on Multimedia and Expo, ICME 2003, 6-9 July 2003, Baltimore, MD, USA*, pages 33–36. IEEE Computer Society, 2003.

[ZHCZ22]   Huangjie Zheng, Pengcheng He, Weizhu Chen, and Mingyuan Zhou. Truncated diffusion probabilistic models and diffusion-based adversarial auto-encoders. *arXiv preprint arXiv:2202.09671*, 2022.

[ZHJL24]   Jingyi Zhang, Jiaxing Huang, Sheng Jin, and Shijian Lu. Vision-language models for vision tasks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.

[ZHL+23a]   Renrui Zhang, Jiaming Han, Chris Liu, Peng Gao, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, and Yu Qiao. Llama-adapter: Efficient fine-tuning of language models with zero-init attention. *arXiv preprint arXiv:2303.16199*, 2023.

[ZHL+23b]   Renrui Zhang, Jiaming Han, Chris Liu, Peng Gao, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, and Yu Qiao. Llama-adapter: Efficient fine-tuning of language models with zero-init attention. *arXiv preprint arXiv:2303.16199*, 2023.

[ZHL+24]   Yu Zhang, Rongjie Huang, Ruiqi Li, Jinzheng He, Yan Xia, Feiyang Chen, Xinyu Duan, Baoxing Huai, and Zhou Zhao. Stylesinger: Style transfer for out-of-domain singing voice synthesis. In *Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada*, pages 19597–19605. AAAI Press, 2024.

[ZJJH21]   Huan Zhang, Yiliang Jiang, Tao Jiang, and Peng Hu. Learn by Referencing: Towards Deep Metric Learning for Singing Assessment. In *Proceedings of the 22nd International Society for Music Information Retrieval Conference (ISMIR)*, 2021.

[ZJM+21]   Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stéphane Deny. Barlow twins: Self-supervised learning via redundancy reduction. In *International conference on machine learning*, pages 12310–12320. PMLR, 2021.

[ZJXD22]   Yixiao Zhang, Junyan Jiang, Gus Xia, and Simon Dixon. Interpreting song lyrics with an audio-informed pre-trained language model. In Preeti Rao, Hema A. Murthy, Ajay Srinivasamurthy, Rachel M. Bittner, Rafael Caro Repetto, Masataka Goto, Xavier Serra, and Marius Miron, editors, *Proceedings of the 23rd International Society for Music Information Retrieval Conference, ISMIR 2022, Bengaluru, India, December 4-8, 2022*, pages 19–26, 2022.

[ZKD+23]   Huan Zhang, Emmanouil Karystinaios, Simon Dixon, Gerhard Widmer, and Carlos Eduardo Cancino-Chacón. Symbolic music representations for classification tasks: A systematic evaluation. In *Proceeding of the 24th International Society on Music Information Retrieval (ISMIR)*, Milan, Italy, 2023.

[ZKW+20]   Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. Bertscore: Evaluating text generation with bert. In *International Conference on Learning Representations*, 2020.

[ZLD24]   Huan Zhang, Jinhua Liang, and Simon Dixon. From audio encoders to piano judges: Benchmarking performance understanding for solo piano. In *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, San Francisco, USA, 2024.

[ZLL+23]   Ce Zhou, Qian Li, Chen Li, Jun Yu, Yixin Liu, Guangjing Wang, Kai Zhang, Cheng Ji, Qiben Yan, Lifang He, et al. A comprehensive survey on pretrained foundation models: A history from bert to chatgpt. *arXiv preprint arXiv:2302.09419*, 2023.

[ZLL+24]   Xiangyu Zhao, Bo Liu, Qijiong Liu, Guangyuan Shi, and Xiao-Ming Wu. Easygen: Easing multimodal generation with a bidirectional conditional diffusion model and llms, 2024.

[ZLO+21]   Neil Zeghidour, Alejandro Luebs, Ahmed Omran, Jan Skoglund, and Marco Tagliasacchi. SoundStream: An End-to-End Neural Audio Codec. *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, 30:495–507, November 2021.

[ZLW+24]   Dong Zhang, Zhaowei Li, Pengyu Wang, Xin Zhang, Yaqian Zhou, and Xipeng Qiu. Speechagents: Human-communication simulation with multi-modal multi-agent systems. *arXiv preprint arXiv:2401.03945*, 2024.

[ZLZ+23]   Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan, Pengyu Wang, Yaqian Zhou, and Xipeng Qiu. Speechgpt: Empowering large language models with intrinsic cross-modal conversational abilities. *arXiv preprint arXiv:2305.11000*, 2023.

[ZMX+23]   Yixiao Zhang, Akira Maezawa, Gus Xia, Kazuhiko Yamamoto, and Simon Dixon. Loop copilot: Conducting ai ensembles for music generation and iterative editing. *arXiv preprint arXiv:2310.12404*, 2023.

[ZNAP95]   Edward N Zalta, Uri Nodelman, Colin Allen, and John Perry. Stanford encyclopedia of philosophy, 1995.

[ZOW+22]   Ye Zhu, Kyle Olszewski, Yu Wu, Panos Achlioptas, Menglei Chai, Yan Yan, and Sergey Tulyakov. Quantized GAN for complex music generation from dance videos. In Shai Avidan, Gabriel J. Brostow, Moustapha Cissé, Giovanni Maria Farinella, and Tal Hassner, editors, *Computer Vision - ECCV 2022 - 17th European Conference, Tel Aviv, Israel, October 23-27, 2022, Proceedings, Part XXXVII*, volume 13697 of *Lecture Notes in Computer Science*, pages 182–199. Springer, 2022.

[ZPW+23]   Pengfei Zhu, Chao Pang, Shuohuan Wang, Yekun Chai, Yu Sun, Hao Tian, and Hua Wu. Ernie-music: Text-to-waveform music generation with diffusion models. *CoRR*, abs/2302.04456, 2023.

[ZQL+24]   Ge Zhang, Scott Qu, Jiaheng Liu, Chenchen Zhang, Chenghua Lin, Chou Leuang Yu, Danny Pan, Esther Cheng, Jie Liu, Qunshu Lin, et al. Map-neo: Highly capable and transparent bilingual large language model series. *arXiv preprint arXiv:2405.19327*, 2024.

[ZQY+21]   Xulong Zhang, Jiale Qian, Yi Yu, Yifu Sun, and Wei Li. Singer identification using deep timbre feature learning with KNN-NET. In *IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2021, Toronto, ON, Canada, June 6-11, 2021*, pages 3380–3384. IEEE, 2021.

[ZRA23]   Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023.

[ZRZY23]   Chen Zhang, Yi Ren, Kejun Zhang, and Shuicheng Yan. SDMuse: Stochastic differential music editing and generation via hybrid representation. *IEEE T-MM*, 2023.

[ZTdCQT21]   Neil Zeghidour, Olivier Teboul, Félix de Chaumont Quitry, and Marco Tagliasacchi. LEAF: A learnable frontend for audio classification. In *Proc. ICLR*, 2021.

[ZTR+22]   Huan Zhang, Jingjing Tang, Syed Rafee, Simon Dixon, and George Fazekas. ATEPP: A Dataset of Automatically Transcribed Expressive Piano Performance. In *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, Bengaluru, India, 2022.

[ZTS21]   Biao Zhang, Ivan Titov, and Rico Sennrich. Sparse attention with linear units. *arXiv preprint arXiv:2104.07012*, 2021.

[ZTW+21]  Mingliang Zeng, Xu Tan, Rui Wang, Zeqian Ju, Tao Qin, and Tie-Yan Liu. MusicBERT: Symbolic Music Understanding with Large-Scale Pre-Training. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 791–800, Online, 2021. Association for Computational Linguistics.

[ZUR17]   Jiaming Zeng, Berk Ustun, and Cynthia Rudin. Interpretable classification models for recidivism prediction. *Journal of the Royal Statistical Society Series A: Statistics in Society*, 180(3):689–722, 2017.

[ZWCX22]  Xulong Zhang, Jianzong Wang, Ning Cheng, and Jing Xiao. Singer identification for metaverse with timbral and middle-level perceptual features. In *International Joint Conference on Neural Networks, IJCNN 2022, Padua, Italy, July 18-23, 2022*, pages 1–7. IEEE, 2022.

[ZWM17]   Frank Zalkow, Christof Weiß, and Meinard Müller. Exploring tonal-dramatic relationships in richard wagner's ring cycle. In *ISMIR*, pages 642–648, 2017.

[ZWO+23]  Ye Zhu, Yu Wu, Kyle Olszewski, Jian Ren, Sergey Tulyakov, and Yan Yan. Discrete contrastive diffusion for cross-modal music and image generation. In *The Eleventh International Conference on Learning Representations*, 2023.

[ZWW10]   Wei Zhao, Xinxi Wang, and Ye Wang. Automated sleep quality measurement using eeg signal: first step towards a domain specific music recommendation system. In *Proceedings of the 18th ACM International Conference on Multimedia*, MM '10, page 1079–1082, New York, NY, USA, 2010. Association for Computing Machinery.

[ZWW+23a] Le Zhuo, Zhaokai Wang, Baisen Wang, Yue Liao, Chenxi Bao, Stanley Peng, Songhao Han, Aixi Zhang, Fei Fang, and Si Liu. Video background music generation: Dataset, method and evaluation. In *IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023*, pages 15591–15601. IEEE, 2023.

[ZWW+23b] Le Zhuo, Zhaokai Wang, Baisen Wang, Yue Liao, Chenxi Bao, Stanley Peng, Songhao Han, Aixi Zhang, Fei Fang, and Si Liu. Video background music generation: Dataset, method and evaluation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15637–15647, 2023.

[ZX21]    Jingwei Zhao and Gus Xia. Accomontage: Accompaniment arrangement via phrase selection and style transfer. *arXiv preprint arXiv:2108.11213*, 2021.

[ZXW23]   Jingwei Zhao, Gus Xia, and Ye Wang. Accomontage-3: Full-band accompaniment arrangement via sequential style transfer and multi-track function prior. *arXiv preprint arXiv:2310.16334*, 2023.

[ZYP+23]  Le Zhuo, Ruibin Yuan, Jiahao Pan, Yinghao Ma, Yizhi Li, Ge Zhang, Si Liu, Roger B. Dannenberg, Jie Fu, Chenghua Lin, Emmanouil Benetos, Wenhu Chen, Wei Xue, and Yike Guo. Lyricwhiz: Robust multilingual zero-shot lyrics transcription by whispering to chatgpt. In *Proceedings of the 24th International Society for Music Information Retrieval Conference, ISMIR 2023, Milan, Italy, November 5-9, 2023*, pages 343–351, 2023.

[ZZL+23]  Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 2023.

[ZZQ+22]  Xueyao Zhang, Jinchao Zhang, Yao Qiu, Li Wang, and Jie Zhou. Structure-enhanced pop music generation via harmony-aware learning. In *Proc. ACM MM*, 2022.

[ZZS+22]  Runbang Zhang, Yixiao Zhang, Kai Shao, Ying Shan, and Gus Xia. Vis2mus: Exploring multimodal representation mapping for controllable music generation. *CoRR*, abs/2211.05543, 2022.

[ZZW+23]  Ziqiang Zhang, Long Zhou, Chengyi Wang, Sanyuan Chen, Yu Wu, Shujie Liu, Zhuo Chen, Yanqing Liu, Huaming Wang, Jinyu Li, et al. Speak foreign languages with your own voice: Cross-lingual neural codec language modeling. *arXiv preprint arXiv:2303.03926*, 2023.

[ZZZ+22]  Hang Zhao, Chen Zhang, Bilei Zhu, Zejun Ma, and Kejun Zhang. S3t: Self-supervised pre-training with swin transformer for music classification. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 606–610. IEEE, 2022.

[ZZZ+23]  Qiushi Zhu, Long Zhou, Ziqiang Zhang, Shujie Liu, Binxing Jiao, Jie Zhang, Lirong Dai, Daxin Jiang, Jinyu Li, and Furu Wei. VatLM: Visual-Audio-Text Pre-Training with Unified Masked Prediction for Speech Representation Learning. *IEEE Transactions on Multimedia*, pages 1–11, 2023. Conference Name: IEEE Transactions on Multimedia.