

# JOINT MUSIC AND LANGUAGE ATTENTION MODELS FOR ZERO-SHOT MUSIC TAGGING

Xingjian Du<sup>1</sup>, Zhesong Yu<sup>1</sup>, Jiaju Lin<sup>1</sup>, Bilei Zhu<sup>1</sup>, Qiuqiang Kong<sup>2</sup>

<sup>1</sup>Bytedance <sup>2</sup>The Chinese University of Hongkong

## ABSTRACT

Music tagging is a task to predict the tags of music recordings. However, previous music tagging research primarily focuses on close-set music tagging tasks which can not be generalized to new tags. In this work, we propose a zero-shot music tagging system modeled by a joint music and language attention (JMLA) model to address the open-set music tagging problem. The JMLA model consists of an audio encoder modeled by a pretrained masked autoencoder and a decoder modeled by a Falcon7B. We introduce preceiver resampler to convert arbitrary length audio into fixed length embeddings. We introduce dense attention connections between encoder and decoder layers to improve the information flow between the encoder and decoder layers. We collect a large-scale music and description dataset from the internet. We propose to use ChatGPT to convert the raw descriptions into formalized and diverse descriptions to train the JMLA models. Our proposed JMLA system achieves a zero-shot audio tagging accuracy of 64.82% on the GTZAN dataset, outperforming previous zero-shot systems and achieves comparable results to previous systems on the FMA and the MagnaTagATune datasets.

**Index Terms**— Music tagging, joint music and language attention models, Music Foundation Model.

## 1. INTRODUCTION

Music tagging [1, 2] is a task to design a system that can automatically predict the tags of music recordings. Music tagging is an essential task in music information retrieval (MIR) and has attracted research interests in both academia and industry. Many existing music tagging systems focus on *closed-set* tasks [2, 3, 4, 5, 6], where the tags are predefined for each dataset. The closed-set music tagging tasks can be identifying the genres, moods, instruments, singers, and eras of music. However, those systems can not generalize to out-of-domain tags. In this work, we address *open-set* music tagging tasks [7] where the evaluation datasets have different tags from the training dataset.

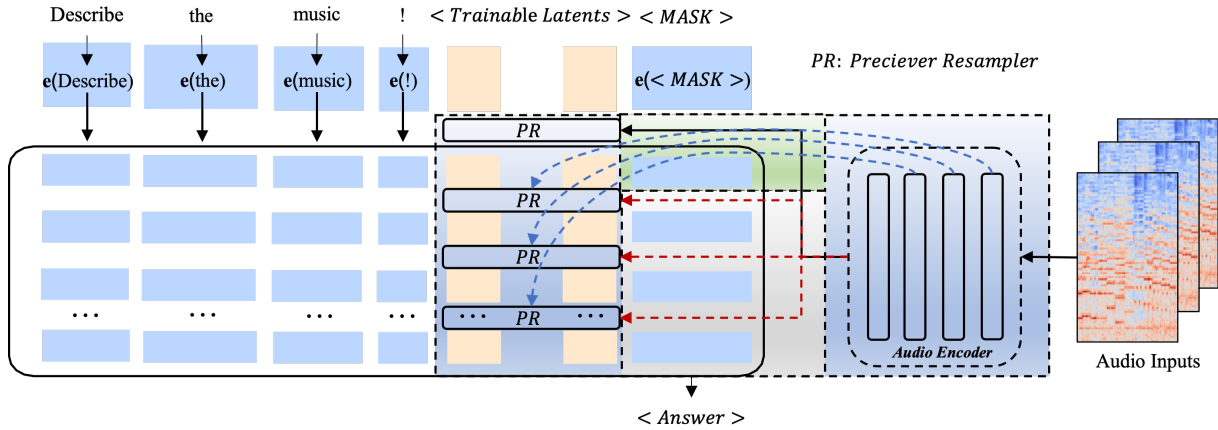
Multi-modal large language models (LLMs) have become a recent research hotspot [8, 9, 10], where connecting audio, image, and video with language models helps utilize their understanding and reasoning capabilities to better accomplish

tasks such as classification and captioning. In this paper, we explore how LLMs can aid in improving the open-set music tagging task. One of the most important questions in multi-modal LLM research is how to allow language models to obtain multi-modal information, or in other words, how to connect multi-modal modules with LLMs. Current mainstream connection methods usually connect the last layer of the multi-modal encoder with the LLM decoder [11]. It is generally believed that the embedding of the last layer of the encoder contains more high-level semantic information, while the middle and low-level semantic information is ignored. However, music tagging tasks have a very wide distribution of label categories, where some labels require high-level semantics, such as genre and emotion, while others require middle and low-level semantics, such as instruments. Therefore, in the music tagging task, it may also be important to output the information of the middle layer of the audio encoder to the language model.

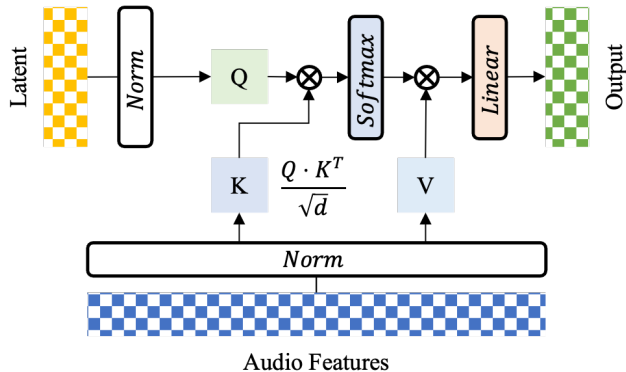
In this paper, we propose a new model called Joint Music and Language Attention Models (JMLAs). This model differs from previous multi-modal LLMs in that we designed a mechanism to perform cross-attention between multiple layers of the audio encoder and multiple layers of the LLM decoder, thereby enhancing the interaction of audio and text information at different semantic levels. This mechanism uses a module called the Perceiver Resampler, which introduces learnable parameters to the cross-attention, thereby bridging the gap between audio and language embeddings. Additionally, it can also reduce the dimensionality of the audio embedding to improve the computational efficiency of the JMLAs.

One challenge of training JMLAs is the lack of publicly available training data [12]. To address this problem, we collect a large-scale music and description dataset from the internet. However, the raw descriptions may contain noisy and irrelevant information about the music. To address this problem, we propose to use ChatGPT [13] to process the raw descriptions into formalized and diverse captions. We show that JMLAs trained with the ChatGPT processed captions achieve better results than the JMLAs trained with raw descriptions.

This paper is organized as follows. Section 2 introduces the joint music and language attention model. Section 3 introduces the text generation with ChatGPT. Section 4 shows experiments. Section 5 concludes this work.



**Fig. 1.** The framework of the JMLA model, consists of an audio encoder, a language decoder, and perceiver resampler-based modules. The inputs to the model including audio and text. The output of the model are the decoded texts. PR: perceiver resampler, uses the cross-attention mechanism to obtain a length-compressed audio representation for better efficiency and to decrease memory consumption while the LLM is processing the injected audio embedding. Red arrows: last-layer injection. Blue arrows: multiple layers injection.



**Fig. 2.** Perceiver resampler.

## 2. JOINT MUSIC AND LANGUAGE ATTENTION MODELS

Our proposed JMLA consists of an audio encoder, a language decoder, and an perceiver resampler-based multiple-layer attention module. Fig. 1 shows the framework of the JMLA.

### 2.1. Audio Encoder

We first transform a time-domain waveform into log mel spectrogram  $X$  with a shape of  $T \times F$ , where  $T$  is the frames number and  $F$  is the frequency bins number. We train a masked autoencoder (MAE) [14] on the log mel spectrogram of a music recording and take the audio encoder of the MAE as the audio encoder. The log mel spectrogram are split into patches  $16 \times 16$  patches along the time and frequency axes. We denote the number of patches as  $P$ . All patches are forwarded into fully connected layers.

We denote the MAE encoder output as  $\mathbf{e} = f_{\text{enc}}(X)$  and the MAE decoder output  $\hat{X} = f_{\text{dec}}(\mathbf{e})$  where  $\hat{X}$  is the es-

timated spectrogram. During training, the MAE in applies a large position of masks that randomly remove 75% of the patches in the input  $X$  [14]. In inference, we remain the MAE encoder layers and remove the MAE decoder layers to train JMLAs. The MAE encoder weights are frozen during the training of JMLAs.

### 2.2. Perceiver Resampler

The disadvantage of previous works [11] is that the computation cost increase quadratically with the length of audio. To address this problem, we propose to use a Perceiver resampler [15, 10] to convert arbitrary length audio embedding into a fixed length embedding. Fig. 2 shows that the Perceiver uses a cross-attention module to project the audio embedding  $\mathbf{e}$  with arbitrary length into a fixed length latent bottleneck  $\mathbf{h} \in \mathbb{R}^{L \times D}$ , where  $L$  is the sequence length of  $\mathbf{h}$ . More specifically, the key  $K$  and value  $V$  are projections of the audio embedding and query  $Q$  is a projection of a learned latent array with a length  $L \ll P$ . The learned latent array contains learnable parameters during training.

### 2.3. Multiple-layer Cross Attention

In previous encoder-decoder architectures [16, 17], the encoder output is directly input to the decoder. There is a lack of information flow between the encoder and decoder layers. To address this problem, we propose to introduce multiple layer connections between the intermediate layers of the encoder and decoder.

We introduce two types of multiple-layer cross attention in JMLA. The first type is to inject the *last-layer output* of the audio embedding into individual perceiver resamplers and multiple decoder layers as shown in the red arrows in Fig. 1.

The second type is to inject *multiple-layer output* of the audio embedding into individual perceiver resamplers and multiple decoder layers as shown in the blue arrow in Fig. 1. By this means, there are more information flow between the audio encoder and language decoder layers.

To preserve the knowledge of the language models, we apply a prefix tuning [18] strategy that freeze the parameters of the language models part while only train the parameters of the perceiver resamplers part. The yellow blocks of Fig. 1 shows the trainable parameters. By this means, the JMLA model can maximize the utilization of the LLM.

## 2.4. Language Decoder

We input the audio representation  $\mathbf{h}$  into the language decoder [19] to autoregressively decode descriptions.

The input to the language decoder is the concatenation of audio representation  $\mathbf{h}$  and the texts of question  $\mathbf{q} = \{q_1, \dots, q_M\}$ , where  $M$  is the number of words of the question. The target of the decoder is the texts of answer  $\mathbf{d} = \{d_1, \dots, d_T\}$ , where  $N$  is the number of words of the description.

The language encoder has a multiple layer causal Transformer architecture. We apply the pretrained Falcon7B [19] as the language decoder. During the training of JMLA, the loss function can be written as:

$$\mathcal{L} = - \sum_{t=1}^T \log(d_t | \mathbf{e}, \mathbf{q}, d_{<t}). \quad (1)$$

Equation (1) shows that the JMLA autoregressively decode the texts of answer in a casual way.

## 3. DATASET

### 3.1. Music-description Dataset Collection

Different from previous music tagging datasets [20, 21, 22] that only contain the close-set tags of music and due to the lack of publicly available music description datasets [12], we collect a large-scale dataset that contains natural language descriptions of music. The descriptions contains plentiful information of genre, speed, era, instrument, emotion, key of the music. We crawl music data from the inhouse data store. We collected 1.5 million audios and description pairs.

### 3.2. Dataset processing with ChatGPT

There are two main problems of the raw descriptions of the downloaded dataset. First, although the descriptions are abundant in quantity, the descriptions can be noisy and irrelevant to the music recordings. For example, some descriptions may focus on the historical background of the music while others concentrate on the instrument analysis. Such dissimilarity in expression will result in difficulty in training.

**Table 1.** Zero-shot tagging results with different systems.

	GTZAN	FMA	MagnaTagATune	
	Acc	Acc	PR	AUC
AudioFlamingo [10] †	38.62	39.62	17.43	61.30
Pengi [11]	32.25	-	-	-
CLAP HTS-AT [23] †	57.24	45.38	31.73	79.60
CLAP MAE (MuLan) [12] †	60.04	33.00	10.63	64.41
JMLA	58.28	41.88	18.64	70.66
JMLA-DenseDec	60.34	42.00	21.06	74.42
JMLA-DenseEncDec	64.82	41.00	20.78	73.51

The † symbol indicates our reproduced system.

To address this problem, we format the descriptions into question-answer formats by using ChatGPT [13]. We design prompts and ask GPT-3.5-turbo to generate question-answer pairs based on given descriptions to increase the diversity.

## 4. EXPERIMENTS

### 4.1. Datasets

We evaluate our proposed music language model on diverse audio tagging datasets to compare our music language model with other systems, including the GTZAN [20], Free Music Archive (FMA) [21], and the MagnaTagATune datasets [22]. The GTZAN dataset consists of 1,000 30-second audio clips with 10 genres. The FMA dataset FMA dataset small set contains 8,000 audio clips with 8 genres. The MagnaTagATune dataset contains 25,863 30-second clips with 188 genres. Due to the large number of labels in this dataset, the top 50 most commonly used labels are usually employed. We adopt the classification accuracy as the evaluation metric.

### 4.2. Zero-shot Music Tagging Results

Table 1 shows the music tagging results of different systems. We compare our system with the Audio Flamingo [10], Pengi [11], CLAP HTS-AT [23], and the CLAP MAE [12] systems using the same data as our system. All systems are evaluated in a zero-shot way. That is, only the evaluation subset of the datasets are used, without trained or finetuned using the training subset of the datasets.

(1) **AudioFlamingo** [10] is a multimodal architecture, with a cross-attention module and resampling blocks. Flamingo is trained on large-scale multimodal web corpora containing arbitrarily interleaved text and images.

(2) **Pengi** [11] is a novel audio language model that leverages transfer learning by framing all audio tasks as text-generation task. In Pengi, audio embedding sequences are combined as a prefix to prompt a pre-trained frozen language model. We use the officially released model weights for evaluation.

(3) **CLAP-HTS-AT**[23] is a language-audio model that aligns language and audio representation in the same space by

**Table 2.** Zero-shot tagging results with different prompts.

Input	Input Example	Output	Postprocess	GTZAN	FMA
				Acc	Acc
Prompt	What is the genre of the music?	Sentence	Similarity	35.80	22.25
Prompt + TagsList	What is the genre of the music? The genres include pop, blues, ...	Sentence	Similarity	40.80	35.00
Prompt + TagsList	What is the genre of the music? Answer one word from pop, blues, ...	OneHot	N/A	56.55	39.75
Prompt + All Candidates	What is the genre of the music? The answer is {pop   blues   ...} .	N/A	Log-likelihood	64.82	40.50

**Table 3.** Zero-shot tagging results with different training data.

	GTZAN	FMA	MagnaTagATune	
	Acc	Acc	PR	AUC
RawCaption	50.68	28.87	25.58	77.37
GPT-QA	60.34	42.00	21.06	74.42
RawCaption + GPT-QA	53.79	41.87	23.90	77.12
GPT-QA Finetune	62.40	43.60	23.63	74.30

contrastive learning. The model is trained on 633,526 audio-text pairs from different data sources. We use the HTS-AT version, where the audio encoder is realized by HTS-AT.

(4) **CLAP-MAE** is our re-implementation of CLAP. Regarding the limited representation ability of HTS-AT, we substitute Audio MAE for HTS-AT as the audio encoder. For fair comparison, we re-train the system with our music-description dataset.

In order to evaluate the correctness of the downstream classification task, we employ a vocabulary ranking method following previous work [11]. We concatenate  $N$  possible candidates after the question and constructing  $N$  complete sentences. The highest log-likelihood score of each candidates is select as the final predicted result.

Table 1 shows that the AudioFlamingo achieves an accuracy of 38.62% on the GTZAN dataset. The Pengi system achieves a lower accuracy of 32.25%. The CLAP HTS-AT or the CLAP-MAE system improve the accuracy to 57.24% and 60.04%, respectively. Our proposed JMLA system without multiple-layer attention achieves an accuracy of 58.28%. Our proposed JMLA-DenseDec refers to the system that the output of the encoder is input to multiple decoder layers improves the accuracy to 60.34%. Our proposed JMLA-DenseEncDec refers to the system that the mutiple layer output of the encoder are input to multiple decoder layers improves tehe accuracy to 64.82%. On the FMA dataset, the JMLA system achieves an accuracy of 41.00%, outperforming the AudioFlamingo and CLAP MAE system, while slightly underperforms the CLAP HTS-AT system. The MagnaTagATune dataset is a more challenging task than the GTZAN and the FMA datasets due to there are 50 genres. The JMLA system achieves a PR of 20.78%, outperforming the AudioFlamingo and the CLAP MAE system.

### 4.3. Different Promots Results

Table 2 shows the JMLA results with different prompts. The prompt consists of a question and the token of the audio. Table 2 shows that different types of prompts will lead to different results. By limiting the predicted genres to be one of the genres in the dataset, the accuracy improves from 35.80% to 40.80% and 56.55%, respectively. In addition, we show that by taking the predicting with log-likelihood the accuracy improves to 64.82%. This results indicate designing prompts are important to the zero-shot music tagging.

### 4.4. Different Training Data Results

Table 3 shows the results of JMLA trained with four types of data: 1) **Raw caption** data crawled from the internet; 2) **GPT-QA** data created by using ChatGPT to process and clean the raw caption data; 3) **RawCaption + GPT-QA** data that concatenates the RawCaption and the GPT-QA data; 4) **GPT-QA Finetune** data that use GPT-QA data to finetune the system pretrained with RawCaption + GPT-QA data.

Table 3 shows that the JMLA system trained with the RawCaption data achieves the lowest accuracy of 50.68% on the GTZAN dataset. When training on the GPT-QA dataset, the accuracy significantly improves the accuracy to 60.34%. This result shows the effectiveness of the GPT-QA data to train the JMLA system. Table 3 shows that when training on the concatenation of RawCaption + GPT-QA data, the accuracy decreases to 53.79%. This result shows that the RawCaption data has negative effect to train the JMLA system. Furthermore, we show that after finetuning the RawCaption + GPT-QA system with the GPT-QA data, the accuracy further improves to 62.40%.

## 5. CONCLUSION

In this work, we propose a joint music language attention model (JMLA) model consists of an audio encoder modeled by audio masked auto-encoder (MAE), a perceiver and mutiple-layer attention module, and a language decoder. Our proposed JMLA can address the open vocabulary audio tagging tasks. We collect a large-scale music language dataset from the website and process the data with a ChatGPT to process the dataset. We show that the JMLA outperforms and achieves comparable results to previous zero-shot audio tagging system on the GTZAN, FMA, and MagnaTagATune

datasets, respectively. In future, we will investigate more music question answer problems with JMLA.

## 6. REFERENCES

- [1] Zhouyu Fu, Guojun Lu, Kai Ming Ting, and Dengsheng Zhang, “A survey of audio-based music classification and annotation,” *IEEE Transactions on Multimedia*, vol. 13, no. 2, pp. 303–319, 2010.
- [2] Keunwoo Choi, György Fazekas, Mark Sandler, and Kyunghyun Cho, “Convolutional recurrent neural networks for music classification,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2017, pp. 2392–2396.
- [3] Keunwoo Choi, George Fazekas, and Mark Sandler, “Automatic tagging using deep convolutional neural networks,” in *International Society of Music Information Retrieval (ISMIR)*, 2016.
- [4] Minz Won, Andres Ferraro, Dmitry Bogdanov, and Xavier Serra, “Evaluation of CNN-based automatic music tagging models,” in *Sound and Music Computing Conference (SMC)*, 2020.
- [5] Minz Won, Keunwoo Choi, and Xavier Serra, “Semi-supervised music tagging transformer,” in *International Society for Music Information Retrieval (ISMIR)*, 2021.
- [6] 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,” *arXiv preprint arXiv:2306.00107*, 2023.
- [7] Kevin Liu, Julien DeMori, and Kobi Abayomi, “Open set recognition for music genre classification,” *arXiv preprint arXiv:2209.07548*, 2022.
- [8] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi, “Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models,” *arXiv preprint arXiv:2301.12597*, 2023.
- [9] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi, “Instructblip: Towards general-purpose vision-language models with instruction tuning,” 2023.
- [10] 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,” in *Advances in Neural Information Processing Systems*, 2022.
- [11] Soham Deshmukh, Benjamin Elizalde, Rita Singh, and Huaming Wang, “Pengi: An audio language model for audio tasks,” *arXiv preprint arXiv:2305.11834*, 2023.
- [12] Qingqing Huang, Aren Jansen, Joonseok Lee, Ravi Ganti, Judith Yue Li, and Daniel PW Ellis, “MuLan: A joint embedding of music audio and natural language,” in *International Society for Music Information Retrieval Conference (ISMIR)*, 2022.
- [13] OpenAI, “GPT-4 Technical Report,” *arXiv preprint arXiv:2303.08774*, 2023.
- [14] 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*, 2022, vol. 166.
- [15] Andrew Jaegle, Felix Gimeno, Andy Brock, Oriol Vinyals, Andrew Zisserman, and Joao Carreira, “Perceiver: General perception with iterative attention,” in *International Conference on Machine Learning (ICML)*, 2021, pp. 4651–4664.
- [16] 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 (NeurIPS)*, vol. 30, 2017.
- [17] 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,” in *Association for Computational Linguistics*, 2022.
- [18] Xiang Lisa Li and Percy Liang, “Prefix-tuning: Optimizing continuous prompts for generation,” *arXiv preprint arXiv:2101.00190*, 2021.
- [19] Ebtesam Almazrouei, Hamza Alobeidli, Abdulaziz Alshamsi, Alessandro Cappelli, Ruxandra Cojocaru, Merouane Debbah, Etienne Goffinet, Daniel Heslow, Julien Launay, Quentin Malartic, Badreddine Noune, Baptiste Pannier, and Guilherme Penedo, “Falcon-40B: an open large language model with state-of-the-art performance,” 2023.
- [20] George Tzanetakis and Perry Cook, “Musical genre classification of audio signals,” *IEEE Transactions on speech and audio processing*, vol. 10, no. 5, pp. 293–302, 2002.
- [21] Michaël Defferrard, Kirell Benzi, Pierre Vandergheynst, and Xavier Bresson, “FMA: A dataset for music analysis,” *arXiv preprint arXiv:1612.01840*, 2016.
- [22] Edith Law, Kris West, Michael I Mandel, Mert Bay, and J Stephen Downie, “Evaluation of algorithms using games: The case of music tagging,” in *International Society for Music Information Retrieval Conference (ISMIR)*, 2009, pp. 387–392.
- [23] 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 *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2022, pp. 646–650.