

WAVCRAFT: AUDIO EDITING AND GENERATION WITH NATURAL LANGUAGE PROMPTS

Jinhua Liang¹, Huan Zhang¹, Haohe Liu², Yin Cao³, Qiuqiang Kong⁴, Xubo Liu²,
Wenwu Wang², Mark D. Plumbley², Huy Phan^{5,*}, Emmanouil Benetos^{1,6}

¹ Centre for Digital Music (C4DM), Queen Mary University of London, UK

² Centre for Vision, Speech and Signal Processing (CVSSP), University of Surrey, UK

³ Xi'an Jiaotong Liverpool University ⁴ The Chinese University of Hong Kong

⁵ Amazon, Cambridge, MA, USA ⁶ The Alan Turing Institute, UK

jinhua.liang@qmul.ac.uk

ABSTRACT

We introduce WavCraft, a collective system that leverages large language models (LLMs) to connect diverse task-specific models for audio content creation and editing. Specifically, WavCraft describes the content of raw sound materials in natural language and prompts the LLM conditioned on audio descriptions and users' requests. WavCraft leverages the in-context learning ability of the LLM to decompose users' instructions into several tasks and tackle each task collaboratively with audio expert modules. Through task decomposition along with a set of task-specific models, WavCraft follows the input instruction to create or edit audio content with more details and rationales, facilitating users' control. In addition, WavCraft is able to cooperate with users via dialogue interaction and even produce the audio content without explicit user commands. Experiments demonstrate that WavCraft yields a better performance than existing methods, especially when adjusting the local regions of audio clips. Moreover, WavCraft can follow complex instructions to edit and even create audio content on the top of input recordings, facilitating audio producers in a broader range of applications. Our implementation and demos are available at <https://github.com/JinhuaLiang/WavCraft>.

1 INTRODUCTION

Large language models (LLMs), such as ChatGPT (OpenAI, 2023) have remarkably promoted the development of artificial intelligence-generated content (AIGC). Driven by large-scale pre-training on massive high-quality textual tokens and reinforcement learning from human feedback (RLHF), LLMs demonstrate advanced capacity in language analysis, rationale, and interaction. While LLMs have attracted increasing amount of attention on topics such as chain-of-thought (Wei et al., 2023a), interpretability (Zhao et al., 2024), and in-context learning (Wei et al., 2023b), they are limited to textual data and fail to engage with a broader range of AIGC tasks.

AI-empowered agents have been devised to tackle more practical applications by equipping LLMs with task-specific modules (Qian et al., 2023). These agents (Shen et al., 2023; Huang et al., 2023) use the LLM to interpret a user query to some basic tasks and call task-specific modules (namely expert models) with an appropriate order. By using a modular approach, AI-driven agents are capable of solving intricate tasks without the requirement of additional training. In the audio domain, WavJourney (Liu et al., 2023e) proposed an AI-driven agent that synthesises an audio clip by connecting speech, audio, and music generative models. An audio script is created based on user instructions and compiled into an executable computer program. The computer program then instructs various audio generative models to synthesize a recording. Despite its success, the current audio agents cannot use audio clips as input, hindering themselves from a broader range of audio generation applications. Considering the collaborative ability of LLMs and real-world need for multimodal interactive creation, a natural question arises: *can we improve audio agents with the ability of audio analysis and transformation?*

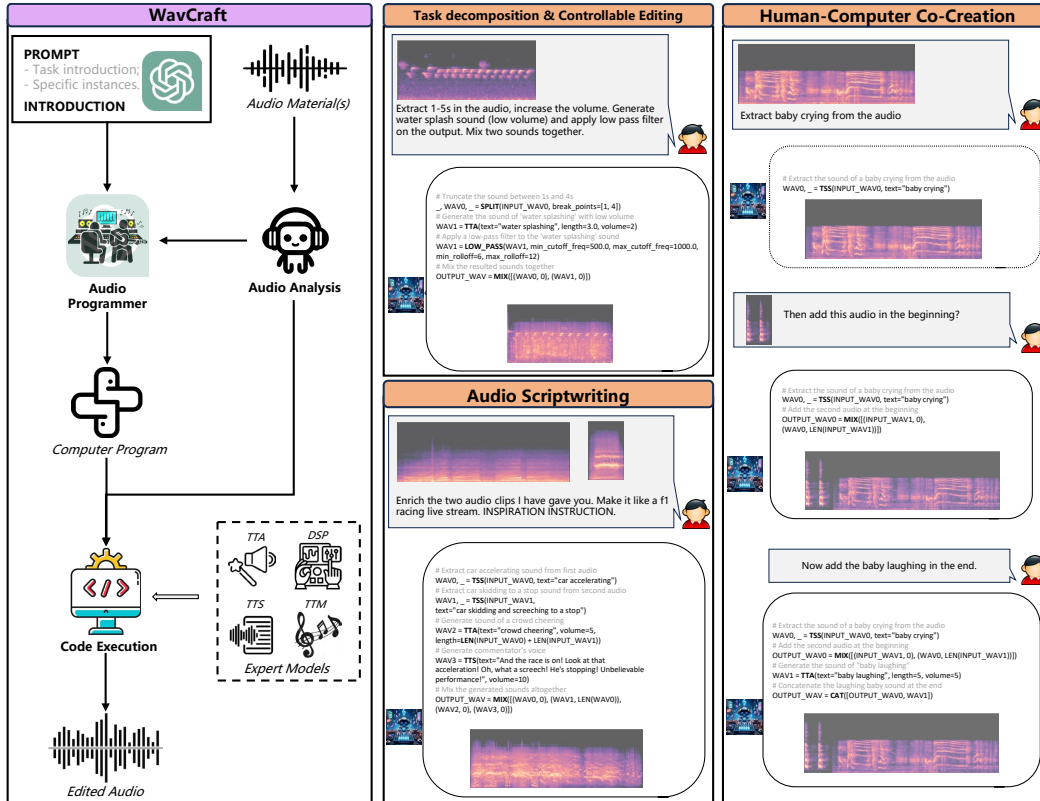


Figure 1: WavCraft overview. WavCraft processes the content of input audio clips and prompts the LLM to generate the code conditioned on user query and audio content. The generated code is then implemented as a computer program empowered by a set of expert models. WavCraft is capable to tackle cases involving: 1) task decomposition; 2) controllable editing; 3) human-computer co-creation; and 4) audio scripting.

In this work, we introduce WavCraft, an AI-empowered audio agent that leverages LLMs, together with a variety of expert models, to create audio content based on human instructions and available audio materials. Specifically, WavCraft describes the content of input audio clips with an audio analysis module. The audio description and the user query are wrapped up with a pre-defined instruction template and directed to an audio programmer module. The audio programmer applies an LLM to break down a complex audio content creation task into several basic ones, and generate an executable program to invoke modules like audio expert models, DSP functions, or logic operations. Following this modular approach, WavCraft is able to assemble a variety of audio creation tools with great flexibility. The overall framework of WavCraft and its featured use cases are shown in Figure 1.

Features of WavCraft includes: 1) Adjustable. WavCraft can take available audio clips as raw materials and create audio content based on both user instruction and input audio. Compared with existing audio agents (Liu et al., 2023e), WavCraft is capable of a broader range of audio content creation, such as audio editing; 2) Modular. WavCraft can break down a comprehensive instruction into several basic audio tasks and thus handle a wide range of audio content generation tasks. In addition, the decompositional framework of WavCraft presents an explicit pathway of content creation, enhancing the explainability in the eyes of users; 3) Interactive. Exploiting the language analysis ability of the LLM, WavCraft interacts with users in multiple dialogues. During multi-round co-creation, the generated audio clips stay consistent with each other; 4) Creative. Based on the analysis of audio content and user’s blueprint, WavCraft leverages the LLM to narrate a story, infers instructions for expert models, and creates the audio content that fulfill the storyline. We refer to such ability of WavCraft to generate audio content without explicit user instruction as *audio scripting*.

The contributions of this paper are summarised as follows:

- An LLM-driven audio agent named WavCraft is proposed to create or edit audio content based on user instructions and available audio clips.
- By coordinating various generative models, WavCraft produces audio content in a controllable manner. Our experiments demonstrate that WavCraft achieves better performance on audio generation and editing compared to current models.
- Additional experiments are conducted to evaluate WavCraft’s ability of audio scriptwriting where models should manipulate audio content without explicit user commands. We hope this will facilitate the process of audio production.

2 RELATED WORK

Language models for multi-modal tasks. Language models such as ChatGPT (OpenAI, 2023) and LLaMA (Touvron et al., 2023) have achieved considerable progress in natural language processing. These models, featuring billions of parameters, are trained with massive high-quality training data to handle a variety of text-related tasks with a single model. To extend the open-world knowledge to more domains, subsequent works (Alayrac et al., 2022; Li et al., 2023) aligned audio, video, and/or image to text and proposed multi-modal foundational models in their respective domains. While multi-modal language models achieved state-of-the-art performance in downstream tasks, most of them are restricted to common text-guided generation tasks, such as text-to-audio synthesis (Liu et al., 2023b), where the model is not required to analyse complex user instructions.

LLM-based agents. More recently, LLM-based agents have attracted an increasing amount of attention to tackle challenging, intricate applications by integrating language models with a set of task-specific models. Toolformer (Schick et al., 2023) was trained to decide when and where an API should be called and how to assemble outputs from different APIs. HuggingGPT (Shen et al., 2023) applied ChatGPT (OpenAI, 2023) as a controller to allocate existing neural networks in Huggingface (Shen et al., 2023). HuggingGPT is thus capable of solving diverse AI tasks across natural language, visual, and audio domains. Meanwhile, ViperGPT (Surís et al., 2023), VisProg (Gupta & Kembhavi, 2023), and RVP (Ge et al., 2023) have demonstrated the promise of visual agents on image/video analysis tasks. LLaVA-Plus (Liu et al., 2023c) extended the input query to visual domain by replacing LLM with visual language model (VLM). In the audio domain, AudioGPT (Huang et al., 2023) connected multiple audio neural networks and used ChatGPT to classify the user query into a predefined task. WavJourney (Liu et al., 2023e) used LLM to screenwrite the audio scripts and then generate audio clips by calling diverse audio generative models. Although considerable progress has been made in previous works to extend open-world knowledge in language models to multiple modalities, few of them can be prompted by non-text inputs, restricting their ability to many practical applications, especially audio editing.

Audio Creation and Editing. Audio creation and editing are challenging parts of generative AI since they require models to not only understand the audio content but also modify the audio conditioned on input instructions. With the development of deep learning, generative models have demonstrated remarkable capacities to synthesise speech (Wang et al., 2023a), audio (Liu et al., 2023b; Kreuk et al., 2023; Borsos et al., 2022; Team et al.), and music (Copet et al., 2023; Agostinelli et al., 2023). Existing audio generation methods are mainly dedicated to synthesising audio conditioned on different types of prompts, such as text description, voice style, and music melody. However, these methods are trained to generate audio from scratch and thus are ill-suited to editing tasks on existing audio. Recently, AUDIT (Wang et al., 2023b) was proposed to learn an end-to-end diffusion model to modify the audio content based on both text instructions and input audio. While AUDIT is capable of various basic editing tasks, including adding, removal, replacement, super-resolution, and infilling, it suffers from two drawbacks: 1) it cannot address complex editing tasks that combines these basic tasks; and 2) it cannot perform local changes on designated audio regions, limiting its application in real-world scenario.

Table 1: List of the audio APIs and their implementation used by WavCraft.

Task	Input	Output	API Name	Model Name
<i>Task-specific models for audio manipulation</i>				
Text-to-Audio	Text	Audio	TTA	AudioGen (Kreuk et al., 2023)
Text-to-Speech	Text	Audio	TTS	Bark
Text-guided Source Separation	Text, Audio	Audio	TSS	AudioSep (Liu et al., 2023d)
Extract	Audio	Audio	EXTRACT	AudioSep (Liu et al., 2023d)
Text-to-Music	Text	Audio	TTM	MusicGen (Copet et al., 2023)
Super resolution	Audio	Audio	SR	AudioSR
Drop	Text	Audio	DROP	AudioSep (Liu et al., 2023d)
Inpaint	Audio	Audio	INPANT	AudioLDM (Liu et al., 2023b)
<i>Basic audio processing functions</i>				
Mix	Audio	Audio	MIX	<code>numpy.add</code>
Length	Audio	Text	LEN	<code>len</code>
Concatenate	Audio	Audio	CAT	<code>numpy.concatenate</code>
Clip	Audio	Audio	CLIP	<code>numpy.ndarray</code>
Adjust Volume	Audio	Audio	ADJUST_VOL	<code>torchaudio.Vol</code>
Low Pass	Audio	Audio	LOW_PASS	<code>audiomentations.LowPassFilter</code>
High Pass	Audio	Audio	HIGH_PASS	<code>audiomentations.HighPassFilter</code>
Room Simulate	Audio	Audio	ROOM_SIM	<code>audiomentations.RoomSimulaor</code>
Impulse Response	Audio	Audio	ADD_RIR	<code>audiomentations.ApplyImpulseResponse</code>

3 WAVCRAFT

3.1 OVERALL FRAMEWORK

WavCraft is an LLM-driven system equipped with a set of task-specific audio networks, capable of audio editing and creation. The overall framework of WavCraft can be found in Figure 1, highlighting three core designs: 1) *audio analysis*: WavCraft initially describes the content of audio clips using natural language; 2) *task decomposition*: given a user query and input audio descriptions, WavCraft formulates a set of instructions from a predefined template by prompting ChatGPT (OpenAI, 2023) directly; 3) *code execution*: WavCraft calls the APIs of expert audio models to execute the generated computer program. We will detail these core designs in the following:

Audio analysis. Complex audio editing requires models to modify audio clips based on user queries and the content of input recordings. Therefore, WavCraft applies the audio analysis module to describe input audio clips in natural language. We apply an audio question and answering model to describe sounds using a template question “write an audio caption to describe the sound”. Please note that WavCraft can integrate any models for audio question and answering (Liang et al., 2023; Deshmukh et al., 2023), or even an audio captioning model (Mei et al., 2023), as the audio analysis module.

Task decomposition. The audio programmer module in WavCraft then drives the LLM to generate an executable script conditioned on both the user query and the content of input audio clips. Specifically, WavCraft fills a pre-defined template with the input query and audio description and then directs the instructions to the LLM. Compared with other AI programmers (Gupta & Kembhavi, 2023), WavCraft generates not only the code but also the comment for each line and the audio script. We found that these comments and the audio script facilitate the audio programmer module to generate code step by step, leading to high-quality output and explainable operations.

Code execution. WavCraft executes the generated scripts by calling a set of audio expert models. Table 1 lists the APIs constituting WavCraft. WavCraft consists of various publicly-available expert models: AudioGen (Kreuk et al., 2023) was used for text-to-audio generation; MusicGen (Copet et al., 2023) was adopted to music synthesis due to its high-fidelity performance based on text and/or melody. For text-to-speech generation, we use Bark¹, a state-of-the-art model that generates matched speech conditioned on the tone, pitch, emotion, and prosody of a given voice preset. For text-guided source separation, AudioSep (Liu et al., 2023d) is used to separate targeted sound tracks conditioned on language queries. AudioSR (Liu et al., 2023a) and AudioLDM (Liu et al., 2023b) are used for super-resolution and audio infilling, respectively. In addition, a series of DSP modules are introduced as well. We implement the DSP modules by using torchaudio (Yang et al., 2022) and

¹<https://huggingface.co/spaces/suno/bark>

audiomentations². It is noteworthy that these task-specific modules can be easily replaced with the alternative architectures.

3.2 FEATURES

Empowered by LLMs, WavCraft is capable of intricate audio editing and creation tasks. WavCraft highlighted four advanced features:

Modular operations. WavCraft decomposes a user instruction into several basic tasks and thus is capable of more complex editing applications in an explainable manner.

Controllable editing. WavCraft translates user requests into executable lines such that it can edit the targeted attributes while keeping the rest unchanged.

Human-AI co-creation. WavCraft leverages large language models to edit audio in a interactive manner, facilitating human producers to create audio content through multi-round refinement. Moreover, WavCraft generates the audio script and comment lines to explain the process of audio content creation. This chain-of-thought method improves the interpretability and transparency of WavCraft.

Audio scriptwriting. Beyond audio content generation under explicit user guidance, WavCraft can produce the sounds in a creative approach, following a high-level outline. We name this ability to devise a plot itself and then manipulate the audio content as *audio scriptwriting*. To make an audio drama, WavCraft creates a script conditioned on input audio, together with the outline, and then sonifies the script with a variety of expert models.

4 TASKS

WavCraft provides a flexible framework that can address a diverse range of audio generative and editing tasks. We evaluate WavCraft on 5 basic tasks, involving adding, removal, replacement, super-resolution, and infilling. We also assess the advanced features of WavCraft on complex tasks.

Adding. Given two audio clips A and B , the model is required to output a mixture M by combining A and B . Suppose C_A is the caption (i.e., text description) of A , an example of the instruction can be “Add C_A in the background of C_B ”.

Removal. Given a mixture M and one of its sound track A , the model is required to output a new audio clip B by removing A from M . Suppose C_A and C_M are the caption of A , respectively, an example of the instruction can be “Remove C_A from C_M ”.

Replacement. Given an audio clip A , a mixture M and one of its sound track B , the model is required to output a new audio clip C by replacing B with A in the same time slot. An example of the instruction can be “Replace C_B with C_A ”.

Super-resolution. Given a low-resolution audio clip A , the model is required to output a new audio clip A' with a higher sampling rate. An example of the instruction can be “Increase resolution of A ”.

Audio infilling. Given an audio clip A where some parts are randomly masked, the model is required to complete the audio by filling the masked areas. An example of the instruction can be “Inpaint A ”.

In addition to the basic tasks, we also evaluate the advanced features of WavCraft through a case study.

5 EXPERIMENTS

5.1 EXPERIMENTS SETUP

To build up WavCraft, we used the GPT-4 model (OpenAI, 2023) as audio programming module and LTU (Gong et al., 2023) for audio analysis. We applied the publicly available models as audio expert models (shown in Table 1). We use 16 kHz sampling rate throughout the pipeline of audio editing and generation in line with the sampling rate of many integrated generative models (Kreuk

²<https://github.com/iver56/audiomentations>

Table 2: Objective evaluation results on five different editing tasks.

Task	AUDIT (Wang et al., 2023b)				WavCraft			
	FAD ↓	IS ↑	KL ↓	LSD ↓	FAD ↓	IS ↑	KL ↓	LSD ↓
Add	9.27	3.87	3.00	1.95	0.63	6.05	1.45	1.59
Removal	17.57	3.27	4.40	3.46	3.48	6.38	1.72	2.07
Replacement	10.24	2.86	3.10	2.55	0.72	6.09	2.16	1.77
Infilling	12.61	3.88	2.86	3.40	3.31	6.37	1.00	2.10
Super-resolution	13.68	2.62	4.25	2.50	5.98	5.96	1.26	1.93

Table 3: Objective evaluation results on the AudioCaps evaluation set.

Model	FAD ↓	KL ↓	IS ↑
AudioLDM (Liu et al., 2023b)	4.65	1.89	7.91
WavJourney(Liu et al., 2023e)	3.38	1.53	7.94
WavCraft	2.95	1.68	8.07

et al., 2023; Copet et al., 2023). For the volume control of the generated audio content, we adopt the Loudness Unit Full Scale (LUFS) standard (International Telecommunication Union, 2020).

We evaluated WavCraft on audio editing and generation tasks separately. We compared WavCraft with AUDIT (Wang et al., 2023b), an state-of-the-art audio editing model, on diverse downstream tasks. For text-to-audio generation, we used AudioLDM (Liu et al., 2023b) and WavJourney (Liu et al., 2023e) for comparison. It is noteworthy that while WavJourney is also an LLM-based agent for audio content generation, it cannot take waveforms as inputs.

5.2 EVALUATION METRICS

For objective evaluation, we follows the evaluation protocols of existing audio generative models (Liu et al., 2023b,e; Wang et al., 2023b) to calculate several measurements: Frechet Audio Distance (FAD), Kullback-Leibler Divergence (KL), Inception Score (IS) and Log Spectral Distance (LSD) for evaluation. FAD measures the Frechet distance between reference and generated audio distributions of the embeddings extracted by a pre-trained VGGish model (Gemmeke et al., 2017). KL computes the similarity between logit distributions of two audio groups by using an audio tagging model, namely Patch-out Transformer (Koutini et al., 2022). IS reflects the variety and diversity of the generated audio group. Log Spectral Distance (LSD) calculates the distance between frequency spectrograms of output samples and target samples. A lower score of FAD, KL, or LSD indicates a better audio fidelity while a higher IS indicates a more diverse audio group (and thus more desirable for generated audio).

5.3 OBJECTIVE EVALUATION

Here we evaluate the performance of WavCraft on audio editing and text-to-audio generation tasks separately.

Evaluation on audio editing tasks. Table 2 shows the objective results of AUDIT and the proposed WavCraft. The WavCraft achieves a better performance than AUDIT in all objective evaluation across different tasks. Compared with AUDIT, WavCraft is 8.97, 14.09, 9.52 lower in FAD on the add, removal, and replacement tasks, respectively. On the audio infilling task, WavCraft achieves the FAD score of 3.31 while having the LSD score of 1.93.

Evaluation on text-to-audio generation. Table 3 shows the objective results of our WavCraft and the compared methods. WavCraft achieves the best FAD and IS score among the three evaluated models on the AudioCaps evaluation dataset. WavCraft also yields the KL score of 1.68, close to the performance of WavJourney (Liu et al., 2023e).

5.4 CASE STUDY

Leveraging LLM’s ability of natural language processing, WavCraft decomposes user’s requests into basic applications and thus is capable of the tasks beyond the common tasks described in Sect. 5.3. As shown in Figure 2, we hereby discuss the advanced features of the proposed WavCraft by studying two cases:

Case study 1: Audio scriptwriting. WavCraft took two raw audio materials (i.e., the beginning and end of a fan heading to a soccer match field) and a user instruction as inputs. It first analysed the content of input audio and wrote an audio script conditioned on both the user instruction and audio descriptions. The script was written in the format of python coding and applied to allocate diverse modules, such as target source separation, text-to-audio generation models, and the room simulator, for audio editing. The output of activated modules was mixed together in line with the generated audio script. To the best of our knowledge, WavCraft is the only audio agent capable of such complex editing task without an explicit user command.

Case study 2: Human-computer co-creation. We illustrate how WavCraft interacts with a user during the process of audio production. WavCraft started with a basic replacement editing task: replacing the female speech in the audio with another female saying. After the user went through the generated recording, namely OUTPUT1.WAV, the user further instructed the system to drop the audio content between 6-10s. WavCraft kept track of the current chat and generated OUTPUT2.WAV based on the new user instruction and the previous execution lines. Likewise, WavCraft continued to generate OUTPUT3.WAV in response to user’s request about “Add more cheers sound in the end” while taking into consideration the previous instructions. Throughout the process of audio production, in addition to the generated audio, WavCraft also provides users with the executable code together with their comments. We hope this will improve the explainability of WavCraft operations.

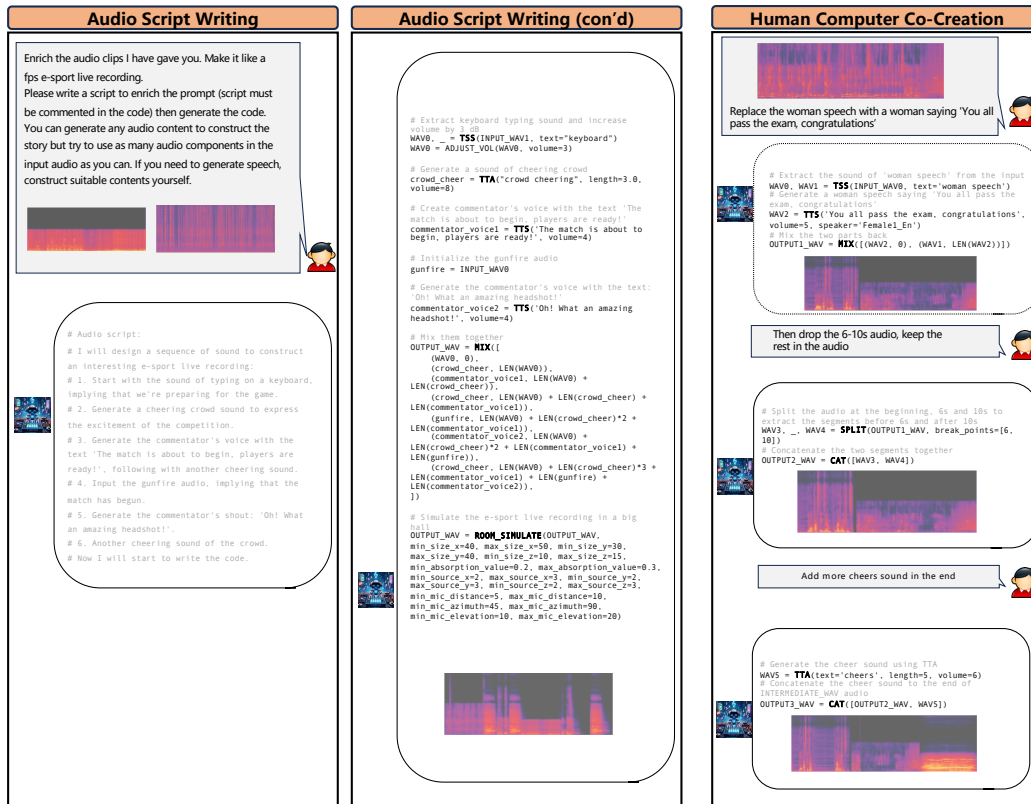


Figure 2: Case studies on audio scriptwriting and human-AI co-creation.

6 LIMITATIONS

Despite WavCraft demonstrating a desirable ability of audio editing and generation, there still exists some limitations: 1) *Audio analysis*: While WavCraft involves an audio analysis module to describe the raw materials, the performance of existing audio analysis models is limited, hindering WavCraft from precisely describing input audio with temporal relationship. 2) *Inference speed*: WavCraft needs to call diverse APIs to solve a complex task, which introduces time costs during inference. Reducing the inference time will facilitate a more seamless human-AI co-creation, meeting the requirement of more practical applications.

7 CONCLUSION

This work presented WavCraft, an agent system that integrates the LLM with diverse audio expert models, to create audio content conditioned on input audio clips and user queries. WavCraft decomposes an intricate editing work into individual basic audio tasks, after comprehending users' queries and the content of given recordings. The output of basic tasks is then assembled under the instructions formulated by audio programming module, contributing to the final output. Case studies conducted on several real-world scenarios have demonstrated the potential of WavCraft in audio production applications. WavCraft shows the feasibility of AIGC in the audio domain in a transparent, interpretable, and interactive manner.

ACKNOWLEDGEMENTS

This work is supported by the Engineering and Physical Sciences Research Council [grant number EP/T518086/1]. E. Benetos is supported by a RAEng/Leverhulme Trust Research Fellowship [grant number LTRF2223-19-106]. H. Zhang is supported by the UKRI Centre for Doctoral Training in Artificial Intelligence and Music, funded by UK Research and Innovation [grant number EP/S022694/1]. This work is also partly supported by Engineering and Physical Sciences Research Council (EPSRC) Grant EP/T019751/1 "AI for Sound (AI4S)", and a PhD scholarship from the Centre for Vision, Speech and Signal Processing (CVSSP) at the University of Surrey and BBC R&D. For the purpose of open access, the authors have applied a Creative Commons Attribution (CC BY) licence to any Author Accepted Manuscript version arising.

REFERENCES

- Andrea Agostinelli, Timo I. Denk, Zalán Borsos, Jesse Engel, Mauro Verzetti, Antoine Caillon, Qingqing Huang, Aren Jansen, Adam Roberts, Marco Tagliasacchi, Matt Sharifi, Neil Zeghidour, and Christian Frank. MusicLM: Generating Music From Text, January 2023. URL <http://arxiv.org/abs/2301.11325>. arXiv:2301.11325 [cs, eess].
- Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: A Visual Language Model for Few-Shot Learning. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.), *Advances in Neural Information Processing Systems*, 2022.
- Zalán Borsos, Raphaël Marinier, Damien Vincent, Eugene Kharitonov, Olivier Pietquin, Matt Sharifi, Olivier Teboul, David Grangier, Marco Tagliasacchi, and Neil Zeghidour. AudioLM: a Language Modeling Approach to Audio Generation, September 2022. URL <http://arxiv.org/abs/2209.03143>. arXiv:2209.03143 [cs, eess].
- Jade Copet, Felix Kreuk, Itai Gat, Tal Remez, David Kant, Gabriel Synnaeve, Yossi Adi, and Alexandre Défossez. Simple and Controllable Music Generation, November 2023. URL <http://arxiv.org/abs/2306.05284>. arXiv:2306.05284 [cs, eess].
- Soham Deshmukh, Benjamin Elizalde, Rita Singh, and Huaming Wang. Pengi: An Audio Language Model for Audio Tasks, May 2023. arXiv:2305.11834.

- Jiaxin Ge, Sanjay Subramanian, Baifeng Shi, Roei Herzig, and Trevor Darrell. Recursive Visual Programming, December 2023. URL <http://arxiv.org/abs/2312.02249>. arXiv:2312.02249 [cs].
- Jort F. Gemmeke, Daniel P. W. 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)*, pp. 776–780, New Orleans, LA, March 2017. IEEE. ISBN 978-1-5090-4117-6. doi: 10.1109/ICASSP.2017.7952261.
- Yuan Gong, Hongyin Luo, Alexander H. Liu, Leonid Karlinsky, and James Glass. Listen, Think, and Understand, May 2023. arXiv:2305.10790.
- Tanmay Gupta and Aniruddha Kembhavi. Visual Programming: Compositional Visual Reasoning Without Training. pp. 14953–14962, 2023. URL https://openaccess.thecvf.com/content/CVPR2023/html/Gupta_Visual_Programming_Compositional_Visual_Reasoning_Without_Training_CVPR_2023_paper.html.
- Rongjie Huang, Mingze Li, Dongchao Yang, Jiatong Shi, Xuankai Chang, Zhenhui Ye, Yuning Wu, Zhiqing Hong, Jiawei Huang, Jinglin Liu, Yi Ren, Zhou Zhao, and Shinji Watanabe. AudioGPT: Understanding and Generating Speech, Music, Sound, and Talking Head, April 2023. URL <http://arxiv.org/abs/2304.12995>. arXiv:2304.12995 [cs, eess].
- International Telecommunication Union. ITU-R BS.1770-4: Algorithms to measure audio programme loudness and true-peak audio level, 2020. URL <https://www.itu.int/rec/R-REC-BS.1770>.
- Khaled Koutini, Jan Schlüter, Hamid Eghbal-zadeh, and Gerhard Widmer. Efficient Training of Audio Transformers with Patchout, March 2022. URL <http://arxiv.org/abs/2110.05069>. arXiv:2110.05069 [cs, eess].
- 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. URL <https://openreview.net/forum?id=CYK7RfcOzQ4>.
- Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. BLIP-2: Bootstrapping Language-Image Pre-training with Frozen Image Encoders and Large Language Models, January 2023. arXiv:2301.12597.
- Jinhua Liang, Xubo Liu, Wenwu Wang, Mark D. Plumbley, Huy Phan, and Emmanouil Benetos. Acoustic Prompt Tuning: Empowering Large Language Models with Audition Capabilities, November 2023. URL <http://arxiv.org/abs/2312.00249>. arXiv:2312.00249 [eess].
- Haohe Liu, Ke Chen, Qiao Tian, Wenwu Wang, and Mark D. Plumbley. AudioSR: Versatile Audio Super-resolution at Scale, September 2023a. URL <https://arxiv.org/abs/2309.07314v1>.
- 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 Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pp. 21450–21474. PMLR, July 2023b. URL <https://proceedings.mlr.press/v202/liu23f.html>.
- Shilong Liu, Hao Cheng, Haotian Liu, Hao Zhang, Feng Li, Tianhe Ren, Xueyan Zou, Jianwei Yang, Hang Su, Jun Zhu, Lei Zhang, Jianfeng Gao, and Chunyuan Li. LLaVA-Plus: Learning to Use Tools for Creating Multimodal Agents, November 2023c. URL <http://arxiv.org/abs/2311.05437>. arXiv:2311.05437 [cs].
- Xubo Liu, Qiuqiang Kong, Yan Zhao, Haohe Liu, Yi Yuan, Yuzhuo Liu, Rui Xia, Yuxuan Wang, Mark D. Plumbley, and Wenwu Wang. Separate Anything You Describe, August 2023d. URL <http://arxiv.org/abs/2308.05037>. arXiv:2308.05037 [cs, eess].

- Xubo Liu, Zhongkai Zhu, Haohe Liu, Yi Yuan, Meng Cui, Qiushi Huang, Jinhua Liang, Yin Cao, Qiuqiang Kong, Mark D. Plumbley, and Wenwu Wang. WavJourney: Compositional Audio Creation with Large Language Models, July 2023e. URL <http://arxiv.org/abs/2307.14335>. arXiv:2307.14335 [cs, eess].
- Xinhao Mei, Chutong Meng, Haohe Liu, Qiuqiang Kong, Tom Ko, Chengqi Zhao, Mark D. Plumbley, Yuexian Zou, and Wenwu Wang. WavCaps: A ChatGPT-Assisted Weakly-Labelled Audio Captioning Dataset for Audio-Language Multimodal Research, March 2023. arXiv:2303.17395.
- OpenAI. GPT-4 Technical Report, March 2023.
- Chen Qian, Xin Cong, Wei Liu, Cheng Yang, Weize Chen, Yusheng Su, Yufan Dang, Jiahao Li, Juyuan Xu, Dahai Li, Zhiyuan Liu, and Maosong Sun. Communicative Agents for Software Development, December 2023. URL <http://arxiv.org/abs/2307.07924>. arXiv:2307.07924 [cs].
- 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, February 2023. URL <http://arxiv.org/abs/2302.04761>. arXiv:2302.04761 [cs].
- Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. HuggingGPT: Solving AI Tasks with ChatGPT and its Friends in HuggingFace, March 2023. URL <http://arxiv.org/abs/2303.17580>. arXiv:2303.17580 [cs].
- Dídac Surís, Sachit Menon, and Carl Vondrick. ViperGPT: Visual Inference via Python Execution for Reasoning, March 2023. URL <http://arxiv.org/abs/2303.08128>. arXiv:2303.08128 [cs].
- Audiobox Team, Apoorv Vyas, Bowen Shi, Matthew Le, Andros Tjandra, Yi-Chiao Wu, Baishan Guo, Jiemin Zhang, Xinyue Zhang, Robert Adkins, William Ngan, Jeff Wang, Ivan Cruz, Bapi Akula, Akinniyi Akinyemi, Brian Ellis, Rashel Moritz, Yael Yungster, Alice Rakotoarison, Liang Tan, Chris Summers, Carleigh Wood, Joshua Lane, Mary Williamson, and Wei-Ning Hsu. Audiobox: Unified Audio Generation with Natural Language Prompts.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu 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, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open Foundation and Fine-Tuned Chat Models, July 2023. arXiv:2307.09288.
- Chengyi Wang, Sanyuan Chen, Yu Wu, Ziqiang Zhang, Long Zhou, Shujie Liu, Zhuo Chen, Yanqing Liu, Huaming Wang, Jinyu Li, Lei He, Sheng Zhao, and Furu Wei. Neural Codec Language Models are Zero-Shot Text to Speech Synthesizers, January 2023a. URL <http://arxiv.org/abs/2301.02111>. arXiv:2301.02111 [cs, eess].
- 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, April 2023b. URL <http://arxiv.org/abs/2304.00830>. arXiv:2304.00830 [cs, eess].
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-Thought Prompting Elicits Reasoning in Large Language Models, January 2023a. URL <http://arxiv.org/abs/2201.11903>. arXiv:2201.11903 [cs].

Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu, Da Huang, Denny Zhou, and Tengyu Ma. Larger language models do in-context learning differently, March 2023b. URL <http://arxiv.org/abs/2303.03846>. arXiv:2303.03846 [cs].

Yao-Yuan Yang, Moto Hira, Zhaoheng Ni, Artyom Astafurov, Caroline Chen, Christian Puhersch, David Pollack, Dmitriy Genzel, Donny Greenberg, Edward Z. Yang, Jason Lian, Jeff Hwang, Ji Chen, Peter Goldsborough, Sean Narenthiran, Shinji Watanabe, Soumith Chintala, and Vincent Quenneville-Bélair. Torchaudio: Building Blocks for Audio and Speech Processing. In *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 6982–6986, May 2022. doi: 10.1109/ICASSP43922.2022.9747236. URL <https://ieeexplore.ieee.org/document/9747236?denied=>. ISSN: 2379-190X.

Haiyan Zhao, Hanjie Chen, Fan Yang, Ninghao Liu, Huiqi Deng, Hengyi Cai, Shuaiqiang Wang, Dawei Yin, and Mengnan Du. Explainability for Large Language Models: A Survey. *ACM Transactions on Intelligent Systems and Technology*, January 2024. ISSN 2157-6904. doi: 10.1145/3639372. URL <https://dl.acm.org/doi/10.1145/3639372>. Just Accepted.