# MIMOSA: Human-AI Co-Creation of Computational Spatial Audio Effects on Videos

Zheng Ning*
University of Notre Dame
Notre Dame, Indiana, USA
zning@nd.edu

Zheng Zhang*
University of Notre Dame
Notre Dame, Indiana, USA
zzhang37@nd.edu

Jerrick Ban
jban@nd.edu
University of Notre Dame
Notre Dame, Indiana, USA

Kaiwen Jiang
k1jiang@ucsd.edu
University of California San Diego
La Jolla, California, USA

Ruohong Gan
ruohongg@andrew.cmu.edu
Carnegie Mellon University
Pittsburgh, Pennsylvania, USA

Yapeng Tian
yapeng.tian@utdallas.edu
University of Texas at Dallas
Richardson, Texas, USA

Toby Jia-Jun Li
toby.j.li@nd.edu
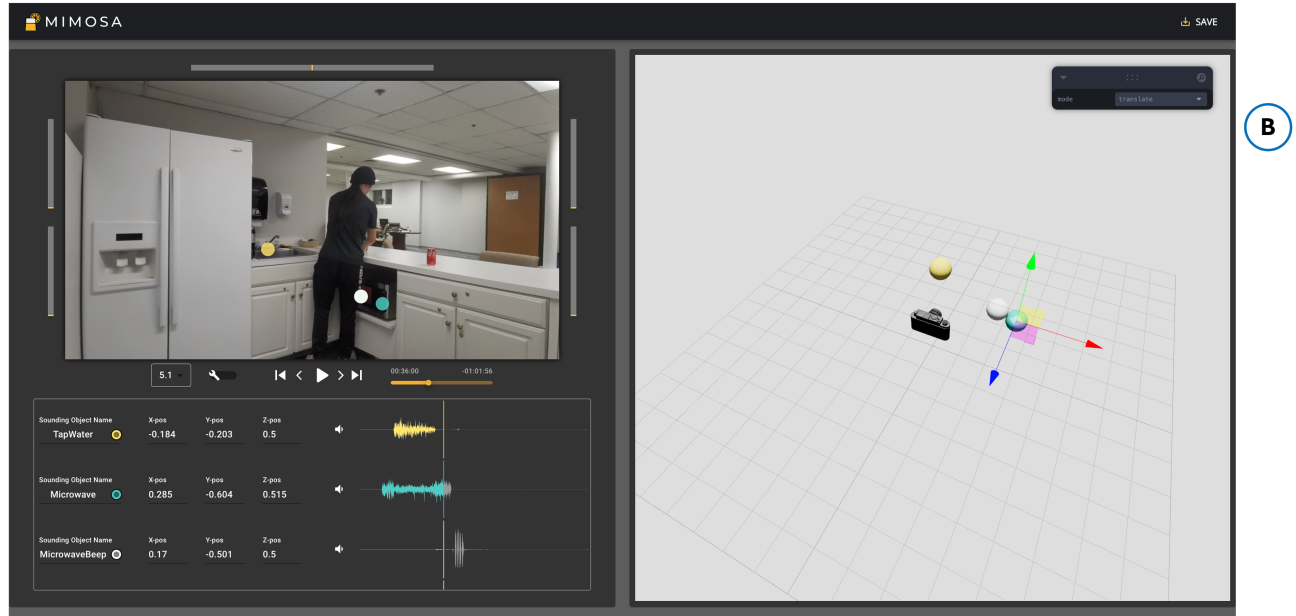University of Notre Dame
Notre Dame, Indiana, USA

Figure 1: The main interface of MIMOSA for spatial audio generation and manipulation

*Both authors contributed equally to this work.

## ABSTRACT

Spatial audio offers more immersive video consumption experiences to viewers; however, creating and editing spatial audio often expensive and requires specialized hardware equipment and skills, posing a high barrier for amateur video creators. We present MIMOSA, a human-AI co-creation tool that enables amateur users to computationally generate and manipulate spatial audio effects. For a video with only monaural or stereo audio, MIMOSA automatically grounds each sound source to the corresponding sounding object in the visual scene and enables users to further validate and fix errors in the location of the sounding objects. Users can also augment the spatial audio effect by flexibly manipulating the sounding source positions and creatively customizing the audio effect. The

design of Mimosa exemplifies a human-AI collaboration approach that, instead of utilizing state-of-art end-to-end "black-box" ML models, uses a multistep pipeline that aligns its interpretable intermediate results with the user's workflow. A lab user study with 15 participants demonstrates Mimosa's usability, usefulness, expressiveness, and capability in creating immersive spatial audio effects in collaboration with users.

## CCS CONCEPTS

• **Human-centered computing → User interface programming**.

## KEYWORDS

video, sound effects, multimodal, creator tools

## 1 INTRODUCTION

Audio effects play a critical role in enhancing the viewing experience for users. In particular, spatial audio, a type of audio effect that allows listeners to perceive the position of sound sources in a three-dimensional space through surround sound [52, 54], has been shown to significantly improve the immersion of video content for viewers. Enhanced immersion promotes user engagement, comprehension, and retention of video content [1, 10, 56].

However, the adoption of spatial audio among video creators, especially amateurs, remains limited. There are several reasons: first, recording spatial audio requires specialized hardware equipment, such as in-ear microphones (for binaural recording) and 360°microphones (for ambisonic recording) [55]. Although the equipment has become increasingly affordable, with entry-level ones costing hundreds of dollars, the cost and expertise required for these specialized microphones continue to hinder their widespread adoption. Second, millions of existing videos have been recorded with only monaural or stereo audio [35], lacking spatial audio information. As a result, content creators cannot directly utilize them as video resources for spatial audio. Lastly, even when the original audio of a video is recorded in a spatial format, post-processing remains challenging and requires specialized tools and expertise [8, 30].

Although end-to-end machine learning (ML) models have made significant progress in generating spatial audio from videos with monaural or stereo audio [12, 35, 76], they are not yet sufficiently useful for meeting the *practical* needs of content creators due to the following reasons:

(i) The creation of audio effects for videos is a personalized process where the quality of the effect is evaluated by the creator's subjective perception [26, 59]; however, the results generated from the ML models are usually evaluated on "ground-truth" quantitative metrics [12, 35, 76]. The "black-box" nature of the state-of-the-art end-to-end ML models prevents users from easily validating and revising their output,

as they lack the capability to produce meaningful intermediate results that end users can understand, validate, and edit.

(ii) Existing ML models are primarily designed to reconstruct the "ground truth" audio; however, video creators often seek to tailor effects to their personal preferences. For instance, the workflow often includes adjusting volume for each sound source separately and changing the spatial perception by redistributing the sound in a 3D space [23, 31]. They may also experiment with various audio settings [9]. In contrast, ML models are limited to predicting spatial audio that closely approximates how the original audio would sound if it had been recorded with spatial audio equipment.

Therefore, using the models alone may not always be the best option in various scenarios. An appropriate tool should support users in continuously exploring the creative space for immersive audio experiences, while providing sufficient flexibility and expressiveness for content creators to realize their ideal audio experiences.

To this end, we introduce Mimosa[1] (Fig. 1), a human-AI collaborative tool that helps amateur content creators create immersive spatial audio effects for videos with conventional monaural or stereo audio. Rather than relying on an end-to-end "black-box" machine learning approach, Mimosa employs a carefully designed audiovisual pipeline (Fig. 2) compromising object detection, depth estimation, soundtrack separation, audio tagging, and spatial audio rendering modules to produce useful intermediate results. Those results, such as the type and position of independent soundtracks of different sounding objects, the estimated 3D position of the sounding object that changes with time, are organized and presented to users through an interactive direct manipulation interface, which allows them to easily manipulate the spatial location and audio attributes of each sound source. To further facilitate real-world deployment of Mimosa, we developed an extension for Adobe Premiere Pro, a popular video editing application among our target users. This extension seamlessly integrates Mimosa with Premiere Pro, allowing users to directly invoke Mimosa while editing videos.

We assessed the quality of spatial audio effects generated by Mimosa through a subjective evaluation conducted by eight independent external evaluators. The evaluators reviewed and rated videos featuring various types of audio effects (including the generated spatial audio by using Mimosa and monaural or raw audio as comparisons) without prior knowledge of the type of audio effect being presented. The result demonstrates that videos featuring spatial audio effects produced by Mimosa were more immersive than the original video sound while maintaining a high degree of realism.

The usability and usefulness of Mimosa is evaluated through a user study with 15 participants who are either amateur content creators or users without prior experience in video or audio editing. Qualitative insights are discussed regarding how users utilize Mimosa's audiovisual interface to assist their audio editing process.

In summary, this paper makes the following contributions:

• Mimosa, an interactive system that enables amateur content creators to create spatial audio effects for videos with monaural or stereo audio only.

---

[1] Mimosa is an acronym for **M**agnifying **I**mmersion by **M**anipulating **O**bjects in **S**patial **A**udio

- A subjective evaluation (N=8) that evaluates the efficacy of the step-by-step spatial sound generation pipeline against the end-to-end models.
- Insights from a user study with 15 participants to create spatial audio effects using Mimosa.

## 2 RELATED WORK

Mimosa builds upon prior research in three key areas. First, we review computational techniques for modeling spatial audio, including approaches for simulating sound propagation and localizing sound sources in videos (Section 2.1). Next, we discuss the diverse applications of spatial audio, which aim to either accurately position sounds in 3D space or foster an immersive audio experience (Section 2.2). Finally, we situate Mimosa within the broader context of AI-enabled co-creation tools for multimedia, and highlight key design strategies that informed the development of our tool (Section 2.3).

### 2.1 Computational Techniques for Spatial Audio Generation

Spatial audio generation is a broad topic that contains multiple individual research questions. One of them is to model changes in head-related transfer function (HRTF) as sound waves propagate. The HRTF describes how sounds are modified due to bouncing, scattering, and diffraction as they approach and enter the ears of a listener [13, 34]. The modeling of HRTF and sound propagation is based on a large amount of sound source and environmental data. In this scenario, deep learning networks were often adopted to model the environment [68], human head and ear geometries [15, 51, 73] and the change in sound features [3, 45]

Another complex challenge arises when dealing with monaural audio recordings from videos that have a mixture of sounding objects. To generate spatial audio effects for such videos, current state-of-the-art approaches often leveraged deep neural networks to associate each independent sound source to their corresponding visual counterparts, and spatialize the video soundtrack with the additional information from the visual objects at different video frames [12, 35, 67, 76].

In our work, we focus on assisting amateur content creators produce realistic spatial audio effects from videos by allowing users to iteratively validate and edit the intermediate outputs. Furthermore, our work targets at allowing the users to augment the spatial effects for more immersive experience or express their creative goals of the audio effects in the video context.

### 2.2 Applications of Spatial Audio

The application of spatial audio can be broadly categorized into two groups based on its aim: (i) accurate 3D positioning and (ii) immersive illusion perception [54]. In the first category, spatial audio provides essential cues that help convey the locations and movements of sound sources. This has been applied in applications such as understanding the positions of multiple speakers [1, 37], augmenting spatial awareness for vision-impaired individuals [19, 39], aiding in navigation for the blind [17], and enriching object perception in augmented reality environments [69]. The second category focuses on enhancing viewer immersion by creating a spatial audio sound field,

which does not necessarily require precise localization of sound sources. Examples include designing the soundtrack of movies to foster an immersive cinema or home theater experience [7, 41], and enhancing audio perception in AR/VR environment [71].

Despite spatial audio has demonstrated utility across various application domains, its widespread adoption remains limited due to barriers in creator expertise and equipment availability. A key goal of Mimosa is to mitigate these barriers, promoting more extensive use of spatial audio in diverse application domains.

### 2.3 AI-Enabled Co-Creation Tools for Multimedia Experiences

Mimosa is part of an expanding group of AI-enabled co-creation tools designed for multimedia experiences. These tools generally aim to achieve two objectives in assisting content creators: (i) improving efficiency, lowering barriers, or reducing effort required in the content creation process; and (ii) expanding the flexibility and expressiveness of the creations [5, 11, 49]. While Mimosa's primary design goal aligns with the first category, it also contributes to the second category by enabling users to "go beyond the ground truth" through manipulation of inferred sounding object positions, ultimately creating customized spatial audio effects.

The design of Mimosa is informed by key design strategies and implications from previous human-AI co-creation efforts in various application scenarios, such as sketching [28], music creation [18, 58], and video creation [26]. Specifically, it ensures that users maintain control and consistently play a "leading role" in the co-creation process [42], fosters accurate mental models of the AI system by employing a pipeline that mirrors the cognitive and reasoning processes of human users [14, 16, 74, 75], and allows users to efficiently handle errors [40, 43, 44] during the process.

## 3 THE MIMOSA SYSTEM

We designed and implemented Mimosa, a human-AI collaborative tool for generating and manipulating spatial audio effects on videos. In this section, we start by outlining the key design challenges and goals of Mimosa. Then, we introduce the architecture of the system (Section 3.2), an example application scenario (Section 3.3), the key interactive features (Section 3.4), and the back-end spatial audio generation pipeline (Section 3.5). This section ends with a discussion of the system's implementation details (Section 3.6).

### 3.1 Technical Challenges and Design Goals

Informed by previous studies on creativity support tools [5, 6, 11] and computational approaches in audiovisual context [64], we identify the following technical challenges (TC) for designing and building Mimosa.

TC1 *Limitations in the performance and capabilities of existing fully-automated ML models.* As discussed in Section 2.1, current state-of-the-art ML models have limited performance in generating realistic and immersive spatial audio from videos containing only monaural or stereo audio. The quality of the generated spatial audio is significantly affected by the training data, video quality, and video scenarios.

TC2 *Lack of user control in the generation process of spatial audio.* Due to the black-box nature of current end-to-end models,
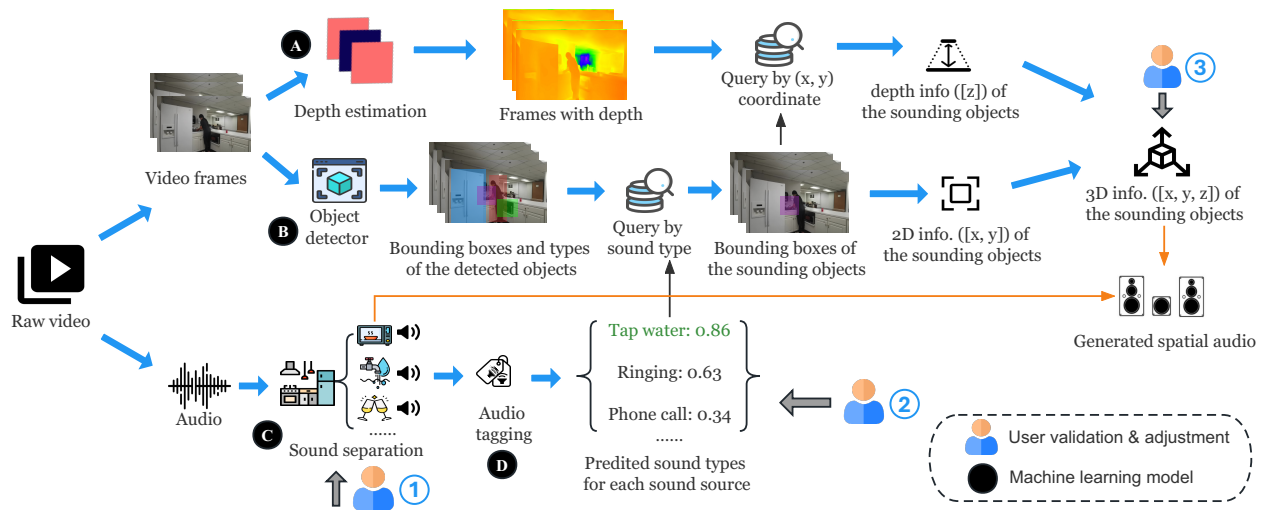
**Figure 2: Mimosa's human-AI collaborative audio spatialization pipeline. Users can validate and adjust the intermediate results in three ways. From left to right, 1: users can adjust the audio properties of each separated soundtrack; 2: users can manually fix the error in aligning the separated soundtrack to the visual object in the video; 3: users can customize the spatial effect for each sounding object by manipulating its corresponding visual position.**

users cannot easily repair the errors in the generated effect, or control the generation process to customize the final output.

**TC3** *Increased information overload and cognitive burden for users working with multimodal data.* As users have to handle both audio and visual data while editing the audio effect for a video, the discrepancies divert users' attention between modalities, making it more difficult for them to manage the task at hand [53]. This issue becomes more pronounced when visual and audio information are inconsistent in terms of their respective positions.

In light of these technical challenges, we identify three design goals for Mimosa.

**DG1** *Allow users to effectively repair the errors in the model-generated spatial audio effects.* To address the limitations of end-to-end models outlined in TC1, the system should provide user-interpretable intermediate results of ML models, allowing users to verify the results generated by the model and rectify errors.

**DG2** *Provide flexibility for users to augment the audio effect for a more immersive experience and creatively modify the spatial audio effects.* Instead of strictly follow the visual relative position of the sounding object to spatialize the audio, the system should also allow users to augment the spatial audio effect. Additionally, it should offer users ample expressiveness, such as flexibly moving the individual sound sources and adjusting the reference points to enable them experiencing different audio settings.

**DG3** *Coordinate user perceptions, cognition, and reasoning while users are handling audiovisual data.* The system should provide effective and easy-to-learn interaction strategies to help reduce the cognitive load when users edit the audio effects in the video context and using the interactive strategies and UI components of the system.

## 3.2 System Architecture

The architecture of Mimosa consists of two layers. When a video is loaded, the audiovisual spatializing pipeline (illustrated in Fig. 2) runs in the backend to process the video and generate a default spatial audio effect. The pipeline also produces intermediate data, enabling users to identify and address any issues in the generation process.

On the second layer, users can modify the intermediate results generated by the machine learning pipeline in the system's main interface (as shown in Fig. 1). It consists of three main components: (i) The video playback and 2D sounding source overlay panel (Fig. 1-A). The panel allows users to manipulate the inferred position of each individual sound source in a 2D space that refers to the video frame by moving the corresponding overlay colored dots. (ii) The simulated 3D manipulation panel (Fig. 1-B). This panel is designed to facilitate user interaction with sound source positioning within a simulated 3D space, it also allows for modifications to the reference point by manipulating the camera object through movements or rotations. (iii) The audio properties display and control panel (Fig. 1-C). This panel allows users to verify and manually specify the corresponding pairings between soundtracks and their associated visual objects within the frame, provides users with a comprehensive overview of the audio properties of each individual soundtrack, and enables users to interact with these properties.

We illustrate the design rationales and detailed descriptions of the interface components in Section 3.4.

## 3.3 Example Usage Scenario

To demonstrate the use of Mimosa, we present a scenario in which an amateur video creator, Lucy (she/her)[2], wanted to enhance the soundtrack of a video featuring two music players playing the flute and violin in a room with immersive spatial audio effects. An example screenshot of a video frame was similar to Fig. 3.

Lucy started by launching Adobe Premiere Pro (Pr) and loading the target video. She selected the desired clip for editing and launched Mimosa from the Extension menu in Premiere Pro. The system automatically started processing the video, when it finished, the main interface of Mimosa (Fig. 1) appeared, allowing her to proceed with editing the default generated spatial effects.

Next, she played the video with the default generated spatial audio effects to *detect* potential issues. As the video played, the positions of dots (Fig.3-E), spheres (Fig.4-A), and numeric numbers (Fig.5-C) changed automatically over time, illustrating the detected sounding objects' positions. She could control the play/pause, play progress, and move backward/forward 1 second using the control buttons (Fig.3-D), similar to using a typical video player. The real-time volume indicators (Fig.3-A) visualized the volume from different channels to improve her understanding of the spatial effect at specific times. Additionally, she could check if the dots (Fig.3-E) representing the inferred spatial positions of sounding objects aligned with the actual visual objects in the frame.

During playback, Lucy noticed that the spatial sound of the flute is incorrect and wants to *repair* it. There were three types of errors in general. First, if the spatial positions of the flute and violin were reversed, she could re-specify their correspondence by editing the sounding object name in Fig.5-A. Second, if the volume of either instrument were too loud or too soft, she can adjust the volume of each instrument separately using the volume controller (Fig.5-D). Finally, if the correspondence was correct, but the prediction was inaccurate, she could: (i) dragged the dot representing the flute towards the desired direction in the video display panel; (ii) manipulated the sphere corresponding to each sound source in the 3D panel, which also allowed her to adjust the distance from the viewer to the sounding object; or (iii) manually changed the numeric 3D position of the flute object in the audio properties display and control panel.

As a video creator, Lucy wanted to enhance the spatial effect by not strictly adhering to the estimated visual positions of the sounding objects, but instead following her creative ideas. Specifically, in this case, she hoped to increase the distance between the two sound sources, thereby enabling the audience to more distinctly perceive sounds appearing from two separate directions. To achieve this, she could simply move the red dot further to the left and the yellow dot further to the right.

She also wanted to test how different positions of the two players affect spatial effects and how the audience perceived the scene from various angles. To achieve this, she could create desired spatial effects by manipulating the position of each sounding object using the three previously mentioned approaches. Additionally, she could change the viewing point in the 3D panel by adjusting the camera object's position and angle. The spatial audio effects were re-rendered in real-time.

Once satisfied with spatial audio creation, she could press the SAVE button on the interface. The edited spatial audio track for each of the sounding objects will be loaded to Pr automatically, allowing her to continue editing other aspects of the video with the functionalities from Pr.

## 3.4 Key Interface Components

Specifically, we designed the following UI components to enable amateur content creators to repair the errors introduced at each stage of the ML pipeline, expand the design space for customized spatial audio effects, and promote user creativity.
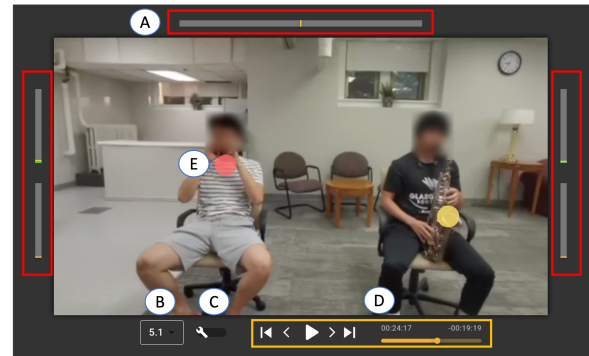


**Figure 3: (A) Users can view the volume of each channel using volume indicators. (B) Users can select the audio output format among monaural, stereo, quadraphonic, and 5.1 channels. (C) Users can toggle whether to create spatial effects based on the model-predicted spatial position or their own. (D) Video control buttons, from left to right: previous video, previous second, play/pause, next second, next video. (E) 2D sound source manipulation, where users can adjust the spatial position of each sounding object by moving the corresponding dot.**

*3.4.1 2D & 3D sound source manipulation panels.* The 2D and 3D direct manipulation panels visually represent sound sources as dots in 2D and spheres in 3D within distinct reference systems. These panels enable users to adjust the spatial positions of sound sources to fix errors in the ML model's predictions or create customized spatial effects.

The 2D panel allows users to manipulate sound sources within the original video frame using dot-like visual representations, as illustrated in Fig. 3-E. Each colored dot represents a specific type of sounding object, with its default [x, y] position predicted by the ML model and its size determined by the predicted depth ([z]) from the imaginary camera position. Users can reposition the objects by dragging the corresponding dots to adjust the [x, y] coordinates as needed.

In contrast, the 3D manipulation panel (Fig.4) operates within a simulated 3D space, independent of the original video frame. Predicted 3D positions of sounding objects are represented as colored spheres. The panel gives users a greater degree of creative freedom when manipulating the positions of sounding objects in several ways. Users can manipulate objects within the open 3D

---

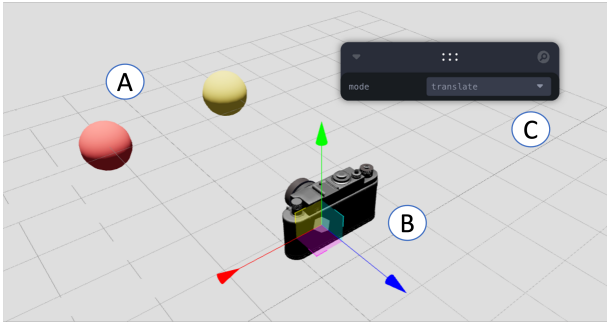[2]A pseudonym for demonstrating the usage scenario.

**Figure 4: (A) 3D sound source manipulation, where users can adjust the spatial position of each sounding object by moving the corresponding sphere. (B) Users can change the viewing point position using the camera object. (C) Users can choose to move or rotate objects, especially the camera object.**

space, modify the point of recording (Fig.4-B) through rotation or translation (Fig. 4-C), and rotate the entire simulated space to view object positions and the camera from various angles.

The goal of the 2D and 3D direct manipulation panels is to allow users selecting preferred reference systems for balancing the trade-off between error correction and creative support. Specifically, the 2D panel allows users to easily correct errors in predicted spatial positions by aligning the dot with the actual sounding object in the video frame. However, it limits depth correction and experimentation with spatial audio effects from different perspectives, thus restricting creativity. The 3D panel, on the other hand, replaces the original video frame reference with a more flexible 3D coordinate system and viewing point, allowing users to manipulate objects within an open space and incorporating diverse interactive strategies to support creativity.

*3.4.2    Video panel.* The video panel displays the video and includes several interactive elements for user control.

First, users can select the audio channel mode from a drop-down list, as seen in Fig. 3-B. Mimosa supports four modes: monaural (1-channel), stereo (2-channel), quadraphonic (4-channel), and 5.1-channel audio.

The number of channels corresponds to the number of volume indicators in Fig. 3-A. Each bar represents the sound intensity of its respective channel. For example, in a 5.1-channel setup, the bars at the upper right, center, upper left, bottom left, and bottom right indicate the sound intensity for the front left, center, front right, rear left, and rear right channels, respectively. The subwoofer's sound intensity is not displayed, as it does not convey directional information. These channel volume indicators aid users in the error discovery process by visually representing sound intensity across each channel.

Additionally, the toggle in Fig. 3-C allows users to enable or disable the use of model-predicted spatial coordinates. When enabled, the system automatically interpolates the position of a sounding object between two frames where the user modifies the positions, reducing the effort required for manual editing of individual frames. When disabled, the spatial coordinates of the sound source remain

unchanged until the user moves the object, providing full user control.

Lastly, users can control video playback using the control buttons in this panel, as displayed in Fig. 3-D.
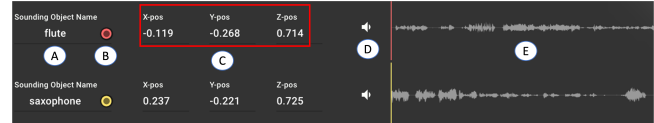


**Figure 5: (A) Users can specify the audio-visual correspondence by changing the name of the sounding object. (B) Object color indicator allows users to change the color of dots (in 2D manipulation) or spheres (in 3D manipulation) representing sounding objects. (C) Numeric spatial position enables users to view and modify the numeric coordinates of sounding objects. (D) By clicking the icon, users can control the volume through a volume slider. (E) Sound waveform display that visualize the current audio.**

*3.4.3    Audio properties display and control panel.* This panel allows users to correct misalignments between separated soundtracks and corresponding visual objects. By default, the identified sounding object name appears in Fig.5-A. Users can select the correct visual object name from a drop-down menu by clicking on the name field if the default name is incorrect. Additionally, users can click the ◎ button (Fig.5-B) to change the color of each sounding object's visual representation in the 2D and 3D manipulation panels.

The panel also displays the spatial location of individual sounding objects in a numeric format (Fig. 5-C). Users can modify the position of each object by directly editing the corresponding numeric value.

Moreover, users can refer to the waveform of each separated sound source (Fig.5-E) to facilitate navigation to specific timestamps in the video. The volume of each sound source can be adjusted by clicking the 🔊 button (Fig.5-D), which activates a volume slider.

## 3.5    Pipeline for Audio-Visual Spatialization

In this section, we introduce our audiovisual pipeline (Fig. 2) for generating spatial audio effects for each sounding object.

*3.5.1    Inferring the 3D positions of the visual objects in 2D videos.* To achieve this, the video is first sampled into individual frames. A pre-trained object detector is then used to extract the 2D positions and respective categories of objects within each video frame (Fig.2-B). We adopt a pre-trained image depth estimation model[20] to generate a pixel-level depth map for each video frame (Fig. 2-A). Using the object's 2D position from the previous step, we query the depth map to obtain the depth information for each object. This approach allows us to accurately determine and track the 3D position of each object over time.

The object detector used in Mimosa was built using the Faster R-CNN neural network architecture [50] and implemented with the Detectron2 framework[3]. The model was pre-trained on the MS-COCO dataset [27], a widely used benchmark dataset for object

---

[3]https://github.com/facebookresearch/detectron2

detection in computer vision. The object detector does not require fine-tuning when integrated into our system.

*3.5.2 Separating individual soundtracks and aligning them with visual counterparts.* At the audio level, a sound separation model [60, 65] is used to extract individual soundtracks from the original sound mixture. Each soundtrack contains sound from a single sounding object (Fig. 2-C).

To align auditory information with visual information obtained in Section 3.5.1, an audio tagging model [21] (Fig. 2-D) is used to predict the sound types for each sound source (e.g., speech, telephone, music, etc.). The separated soundtracks are then mapped to the corresponding category of visual objects, using their names as bridges.

*3.5.3 Real-time spatial audio rendering.* This module takes the separated soundtrack for each sound source, mixes them together based on the corresponding 3D positions of each object, and distributes the mixed sound signal into different channels depending on the speaker setup.

The module is built using the `PannerNode` object in the WebAudio API[4], which simulates sound effects in a simulated 3D space given the 3D coordinates of the sounding object and the listening position. Since the original `PannerNode` only supports two-channel audio rendering, we combined two `PannerNodes` to compute spatial effects separately for the sounding sources from the front and the back.

## 3.6 Implementation

The front-end of Mimosa is hosted on a virtual machine through Google Cloud. The Adobe Premiere Pro plugin for Mimosa is built with the Adobe Common Extensibility Platform (CEP)[5]. Mimosa's front-end UI was implemented in React[6] rather than using the native CEP APIs in Premiere Pro to enable flexible interactions. The backend spatializing pipeline runs on a workstation with an AMD Ryzen Threadripper 3960X CPU and an NVIDIA RTX A6000 GPU.

## 4 TECHNICAL EVALUATION

In this independent evaluation, we want to answer two questions:

(i) How effective is the audiovisual sound spatialization pipeline against the offline end-to-end model?

(ii) How is the quality of the spatial audio effect generated by users using Mimosa?

## 4.1 Dataset

We recorded six sample videos across various scenarios. The videos were recorded with original spatial audio effects, serving as one of the contrast conditions in subsequent analyses. Each video is approximately one minute long to maintain a balance between the overall duration of the user study and the number of videos participants we can test. Detailed information on each video is shown in Table 1.

| Video ID | Video Scenario | Duration | Number of Sound Sources |
|---|---|---|---|
| V1 | Vehicle honking | 0:52 | 2 |
| V2 | Man speaking | 0:54 | 2 |
| V3 | Music duet | 0:43 | 2 |
| V4 | Man playing basketball | 1:26 | 2 |
| V5 | Dog barking | 1:06 | 2 |
| V6 | Cooking in the kitchen | 1:37 | 4 |

**Table 1: Summary of the video details for the user study.**

Each video clip with the ground-truth spatial audio was recorded with an iPhone 13 camera and a ZOOM H3-VR microphone[7], a popular economic 360° audio recorder with a four-capsule ambisonic microphone array.

## 4.2 Evaluation Method

We recruited eight independent evaluators to subjectively evaluate different types of audio. An evaluator viewed a video with an associated audio type at a time. The evaluators were 20 to 30 years old and were recruited from the local community. They had no previous experience with Mimosa and did not participate in the user study, ensuring that their ratings remained unbiased. Each evaluator reported having experience consuming videos with and without spatial audio effects.

Each video has five different audio types:

(1) Raw audio (RA): the pre-recorded original spatial audio.
(2) Monaural audio (MA): The pre-recorded original audio played in a single-channel setting (no spatial effect).
(3) Spatial audio generated by an offline end-to-end model [67] (OA).
(4) Default-generated spatial audio by Mimosa's pipeline without any user intervention (DA).
(5) Spatial audio created by users using the Mimosa system (UA).

We shuffled the order of video clips created by participants (UA) along with the other four versions (RA, MA, OA and DA) and presented them to the evaluators in random order. Each evaluator rated about 12 videos, each with five different audio types.

Evaluators were asked to rate each video on `Immersion` and `Realism` separately on a 7-point semantic differential scale. `Immersion` evaluates the extent to which the spatial audio effect was perceptible, and how well the spatial locations of objects could be inferred from the audio effect. `Realism` focuses on identifying any distortion or dissonance that could be perceived as unreal.

## 4.3 Results

The quantitative rating results are shown in Table 2. The value of each cell represents the average rating for the corresponding audio type in the video. We found **videos with spatial audio effect after user editing with Mimosa (UA) were more immersive than other types of audio effects.** We adopted the Friedman test to measure the difference among the five audio types. The result showed that there was a significant difference among the mean

---

[4]https://developer.mozilla.org/en-US/docs/Web/API/Web_Audio_API
[5]https://github.com/Adobe-CEP
[6]https://reactjs.org

[7]https://zoomcorp.com/en/us/handheld-recorders/handheld-recorders/h3-vr-360-audio-recorder/

| Metric | Video ID | Audio type | | | | |
|--------|----------|------------|---|---|---|---|
| | | Monaural (MA) | Raw (RA) | Offline model (OA) | Mɪᴍᴏsᴀ default (DA) | User generated (UA) |
| **Immersion** | V1 | 2.56 | 4.5 | 3.25 | 4.81 | 6.12 |
| | V2 | 1.69 | 5.69 | 2.36 | 5.12 | 6.31 |
| | V3 | 1.5 | 4.88 | 4.69 | 4.25 | 6.00 |
| | V4 | 2.06 | 3.81 | 2.25 | 5.00 | 6.25 |
| | V5 | 1.63 | 4.69 | 2.06 | 4.38 | 5.50 |
| | V6 | 2.25 | 3.5 | 1.94 | 3.12 | 5.94 |
| | Average | 1.95 | 4.51 | 2.76 | 4.47 | 6.03 |
| **Realism** | V1 | 5.94 | 5.69 | 3.81 | 3.06 | 5.63 |
| | V2 | 5.69 | 6.31 | 3.44 | 3.13 | 5.63 |
| | V3 | 5.44 | 5.81 | 4.81 | 3.63 | 5.38 |
| | V4 | 5.75 | 6.06 | 4.38 | 4.19 | 5.44 |
| | V5 | 5.38 | 5.88 | 3.25 | 3.75 | 5.88 |
| | V6 | 5.88 | 6.44 | 4.06 | 4.25 | 5.56 |
| | Average | 5.68 | 6.03 | 3.96 | 3.67 | 5.58 |

**Table 2: Average rating scores of immersion and realism for videos with varied audio types by external evaluators.**

immersion scores of various audio types ($p < 0.001$). Furthermore, we compared the difference between every two audio types with Wilcoxon signed-rank test. The result shows that the differences between the mean of UA and all other audio types are also statistically significant ($p < 0.001$).

Additionally, we found that immersion score of DA ($\mu = 4.47$; $\sigma = 0.94$) was close to RA ($\mu = 4.51$; $\sigma = 1.19$). The difference of ratings between DA and RA is not significant ($p = 0.54$). The results suggest that spatial audio effects generated using Mɪᴍᴏsᴀ's pipeline, even in a fully automated way, achieve an immersive experience comparable to that of raw audio recorded with a 360-degree audio recorder. The user's validation and revision in Mɪᴍᴏsᴀ further improved the immersion of the generated spatial audio effects. Noticeably, the default-generated effects are better in *realism* and *immersion* than OA, which is generated by an end-to-end ML model.

Last, we found **spatial effects after user editing with Mɪᴍᴏsᴀ (UA) compromised realism compared to raw audio (RA) but show improvement over the fully-automated outcome (DA).** We identified a significant difference among various audio types through a Friedman test in Realism ($p < 0.001$). Specifically, RA received the highest mean score ($\mu = 6.03$; $\sigma = 0.59$), followed by MA ($\mu = 5.68$; $\sigma = 0.33$). This observation suggests that the computationally generated spatial audio inevitably introduces distortion, leading to a decrease in user-perceived realism.

## 5 USER STUDY

Informed by prior research in measuring creativity support and heuristic evaluation of digital tools [4, 38], we conducted an in-person lab study[8] to evaluate the usability, usefulness, and user experience of Mɪᴍᴏsᴀ. The study seeks to answer the following research questions.

RQ1: How usable is Mɪᴍᴏsᴀ in assisting users to edit spatial effects and support their expressiveness?

RQ2: What is the user perception towards the interactive strategies of Mɪᴍᴏsᴀ?

RQ3: How would users utilize Mɪᴍᴏsᴀ to augment the spatial audio effects while editing?

### 5.1 Participants

15 participants (9 men, 6 women) aged 20–30 were recruited from the local community to use the Mɪᴍᴏsᴀ system to generate spatial audio effects for videos. 8 participants were amateur video creators who had posted their videos on websites such as YouTube, while the other 7 did not consider themselves video content creators. 6 participants were experienced with video editing tools such as Adobe Premier Pro; 7 were novice users of these tools, and the rest 2 had no video editing experience. Each participant was compensated with $15 USD for their time.

### 5.2 Study Procedure

The study was conducted in a usability lab (Fig. 6). The lab is equipped with a set of Logitech Z606 5.1 surround speakers[9] that provided audio playback during the study.

At the beginning of each session, the participant signed the consent form and completed a demographic questionnaire. After the researcher gave a brief introduction to the study, the participant watched a five-minute tutorial video on how to use Mɪᴍᴏsᴀ and a one-minute 5.1 channel spatial audio test video to check the surround speakers.

In the study session, each participant completed two tasks for each of the six video clips. The order of the six video clips was randomized, but the tasks for each video clip followed the same order.

The first task seeks to evaluate Mɪᴍᴏsᴀ's capability in facilitating effective human-AI collaboration for users to validate model-generated results and repair any errors they find. In the first task, the user was asked to create *realistic* spatial sound effects for the

---

[8]The study protocol has been reviewed and approved by the IRB at our institution.

[9]https://www.logitech.com/en-us/products/speakers/z606-surround-sound-system.980-001328.html

**Figure 6: The setup of the user study.**

video by repairing the inaccuracy and errors in the video with default generated spatial audio (DA) using Mimosa.

The second task evaluates Mimosa's capability to support flexible customization of spatial audio effects. To familiarize the participant with creation process, we demonstrated the function by guiding the users to create two simple scenarios. The two scenarios were: *"Please move the Saxophone player so that it gradually moves away from you.", "Please turn your face back to the basketball player."* The researcher would answer any questions from the participants during this warm-up process. After creating the effects according to the provided instructions, the user was asked to use Mimosa to create a customized spatial audio effect to their own liking for the current video clip.

The study session ended with a post-study questionnaire and a 10-minute semi-structured interview. In the questionnaire, participants rate statements on the usability and user experience of Mimosa and its key interaction features on a seven-point Likert scale from *"strongly disagree"* to *"strongly agree"* (The questions and results are shown in Fig. 7). During the interview, we asked follow-up questions about their responses to the post-study questionnaire, especially when they gave negative ratings to any aspects of Mimosa. Following established open coding methods [2, 24], an author conducted a thematic analysis of interview transcripts to identify insights and findings particularly related to user experiences, challenges, system usability, and suggestions for new features using an inductive approach.

All 15 participants successfully completed both tasks for all six video clips. All user study sessions were video recorded with the consent of the participants.

## 5.3 Results and Findings

Overall, participants were satisfied with their experience with Mimosa. They found it easy to use, useful, expressive, and capable of generating immersive spatial audio effects. Specifically, Mimosa scored 6.47 ($\sigma = 0.52$) in *"Mimosa is useful"*, 6.20 ($\sigma = 0.68$) in *"the spatial audio effects created through Mimosa is immersive"*, 6.27 ($\sigma = 0.96$) in *"Mimosa allows me to freely create the spatial audio effects that I like"*, and 5.87 ($\sigma = 1.06$) in *"Mimosa is easy to use"*. Regarding the effectiveness of Mimosa's interaction features, the

3D direct manipulation panel scores 6.07 ($\sigma = 1.09$), and the 2D dots visualization overlay scores 6.53 ($\sigma = 0.83$).

**Generate spatial audio effects with ease.** Participants stated that they could *"quickly get familiar with the functions of the system."* (P1) and *"speed up the editing process after finishing editing the first several videos."* (P13). Additionally, most participants mentioned that the spatial effects after editing were more immersive. When asked about the factors contributing to the immersion, participants provided examples such as *"I can clearly feel the car is moving from left to right."* (P8); and *"When I moved the Saxophone to my back, the sound was actually coming from that position."* (P3).

However, several participants noted issues, such as the back-to-front sound transition not being as natural as expected. Users reported hearing a *"crunch sound"* (P9) during playback. Additionally, some users mentioned that the visual position of the sounding objects in the interface did not always align with the auditory position they perceived in the playback. For example, P13 mentioned that *"In one of the videos, I wish I could move the sound source further away, but even though I did it in the interface, the sound I heard was still louder than I expected."*

**Support creativity and expressiveness by intuitive manipulation and real-time playback.** Informed by the existing evaluation dimensions for creativity support tools [4], we asked users about how enjoy, expressive, and immersive when they were using the tool, and how helpful the tool was in allowing them tracking different ideas. Participants expressed that they could easily *"test a variety of different audio settings"* (P14) due to the flexible manipulation strategies offered by Mimosa. P11 also mentioned that *"Aligning the dots and the sounding objects in the video frame cost nearly no labor to me, so I felt really excited playing around with different settings."* Participants also appreciated the flexibility of augmenting existing spatial effect by intentionally *"Increase the distance between the two sounding objects or bring them closer to the reference point."* (P3). For the 3D manipulation panel especially, participants applauded it as it allows for flexibly changing the viewing point, which offered them more creativity control. For example, P13 stated that *"The 3D panel allows me to listen to different spatial effects from varied perspectives in a simulated 3D space by moving the viewing point."* More specifically, P15 stated that *"It (3D/manipulation panel) allows me to imagine I can walk into the video scene and pretend to be at a*
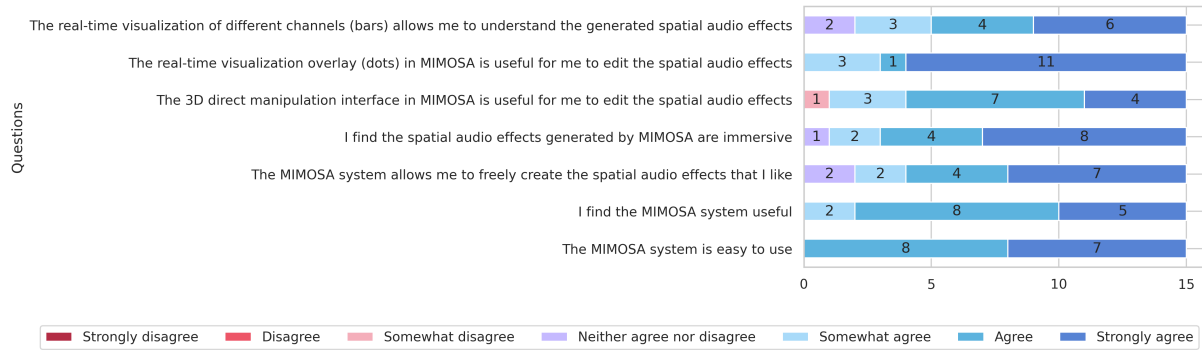
**Figure 7: Results of the post-study questionnaire**

*particular position to hear the music instruments playing."* Additionally, a number of participants emphasized that being able to hear the effects after they made changes further encourages them to test more scenarios due to the minimized cognitive load. For instance, P8 expressed that: *"getting instant feedback is super important in the editing process...this allows me to do the editing and evaluating tasks simultaneously."*

**Handle audio error by using visual clues.** We found participants who were amateur content creators relied on using the visual clues to discover and repair errors in the generated spatial audio effects. We observed that most users first discovered the inaccuracy in the generated spatial effects by noticing the difference between the position of the colored dots and the corresponding visual objects. For example, P13 noted that comparing with listening to the generated spatial audio and trying to identify the errors, by looking the misalignment of the representative dots overlay, he could *"easily found out unnoticed errors by just 'hearing'."* P15 also mentioned that: *"(the strategy) makes fixing the error easy and intuitive, as I just need to drag the dots to make it overlap with the objects."* The user feedback suggests that MIMOSA is able to augment the audio perception and understanding of the amateur creators when perceiving the spatial audio effect, achieved by providing the associated visual information.

**Selectively choose the manipulation method depending on personalized user needs.** MIMOSA provides three distinct approaches for participants to manipulate the location of sounding objects in space (Figure 3-E, Figure 4, and Fig. 5-C). Most participants found manipulating sounding objects using 2D or 3D manipulation panels by *dragging* was more user-friendly and effective than editing the numerical values of their coordinates. For instance, participant P2 noted that: *"I only used the numerical values as a reference, because the spatial effect changes were more immediate when directly moving the dot or spheres."* However, some participants preferred numerical input as the primary editing method. Participant P1, for example, stated that: *"I felt more confident editing spatial audio effects when using the numerical input, especially when I tried to create a new spatial audio effect and wanted to position different objects at various places."* The result suggests that users have various preferences in choosing the way of manipulation. The decision depends on whether they prefer a more deterministic way

or a free-form one. Noticeably, there are some participants (P6, P10, P11) mentioned that they do not have significant preferences but to mix-use all the functions.

**Suggestions for new features.** During the interviews, participants offered feedback and recommendations for potential new features in MIMOSA. P8 proposed the addition of a separate drop-down menu within the system, allowing users to select specific objects they wish to manipulate. This menu would facilitate distinguishing between various audible objects when their corresponding spheres (in the 3D panel) or dots (in the 2D panel) overlap. Another concern expressed by participants was the difficulty in accurately navigating back to previous actions. P7 suggested incorporating a new panel that records actions along with their corresponding timestamps. This addition would enable users to effortlessly review and modify their action history, consequently enhancing the editing efficiency.

## 6 DISCUSSION

### 6.1 Error Discovery and Repair as a Foundation for Human-AI Co-Creation

Effective human-AI collaboration requires a clear understanding of mutual goals [62] and trust between the human and the AI system [47]. However, most machine learning models will inevitably introduce errors into predictions that threaten mutual understanding and trust, making human-AI collaboration systems uniquely difficult to design [70]. The study results of MIMOSA show that an effective error discovery and repair mechanism is crucial in the context of human-AI co-creation where each party needs to work closely with the other party and build upon each other's intermediate results.

The design of MIMOSA exemplifies the use of error discovery and repair as a foundation to support effective human-AI co-creation. Instead of training an end-to-end model to predict spatial audio, MIMOSA uses a multi-step pipeline that produces useful and interpretable intermediate results at each of its steps for users to review, validate, and edit the model results. For example, the 2D panel overlays the results of model-inferred sounding objects on top of the original video, allowing users to easily validate whether the inferred positions align well with the underlying visual objects and repair them through direct manipulation. As discussed in Section 5.3, many study participants found that the interaction

strategies in Mimosa made it easy for them to discover and repair the errors. The effectiveness of such a method was also evidenced in the improvement in *realism* from DV to EV in our evaluation.

The success of Mimosa in error discovery and handling demonstrates the adoption of classic theories in multi-modal interactions [43, 44] in the new context of human-AI co-creation of multimedia contents. Specifically, the design of Mimosa utilizes the *mutual disambiguation* paradigm. Because the generated result of the model is auditory, it is more difficult and cognitively demanding for the user to validate whether it aligns well with the visual information in the original video. To address this problem, Mimosa simultaneously represents the result of generated spatial audio effects visually on the 2D/3D panels, allowing users to easily compare its alignment with the visual information in the video. To make edits, the user can either manipulate it visually (i.e., moving the visual indicator) or adjust its auditory properties directly through the soundtrack information panel. The availability of both interaction modalities at the same time enables the user to choose whichever modality is most natural for their task goals (as we learned from observations and interviews from the user study), while the synchronous representation of the result in both modalities allows the users to easily identify issues in either modality.

## 6.2 Going Beyond the "Ground Truth"

Compared to traditional media editing tools, creative tools often promise users full expressiveness that allows them to flexibly explore the creative space beyond "realism". However, this characteristic poses significant challenges when incorporating machine learning models in creative tools for human-AI co-creation: These models were built to predict the results of "ground truth" learned from the training data. Although some generative models such as [29, 48] and recent commercial systems such as DALL-E[10] can generate impressive "artistic" work on their own, they lack the support of user control, which is necessary for enabling true co-creation.

The design of Mimosa illustrates an approach to support user-initiated augmentation and creation based on model predictions. Throughout the pipeline, the model's strategies to generate and present its output are aligned with the user's existing mental model and workflow of the task. For example, the 3D panel visualizes the positions of each inferred sounding object and the reference of viewing point (camera) in a simulated 3D space. The user can either move the sounding objects, move the camera, or point the camera in a different direction (changing the viewing angle). All these operations provide close analogies to the manipulation of real-world objects and therefore pose a lower learning curve to the users. The user can also easily amplify the presence of a particular sounding object by directly manipulating the intensity of its sound in the sound track information panel. As a result, users can easily understand and leverage the model predictions in this specific context. In this way, users can offload laborious lower-level sub-tasks to ML models and focus on pursuing higher-level creation goals.

## 6.3 Human-AI Collaboration for Augmentation of Multimedia Content

Mimosa's approach to creating spatial audio for videos with only mono or stereo audio illustrates a approach that computationally *augments* existing multimedia contents to a new format with richer information through human-AI co-creation. The AI model can predict the outcome of the new format, while users can help validate and fix the alignment between the old format and the new format.

We plan to continue exploring this approach in other multimedia domains such as augmenting regular 2D video content into AR/VR content in 3D. Content creation for VR/AR faces similar challenges to what we encountered in our problem domain—it requires specialized recording equipment and significant creator expertise. While there are millions of existing 2D video footage available, they cannot be easily used in the 3D VR/AR space. Although ML models are available for synthesizing 3D contents from 2D ones [22, 32, 66], the limited accuracy and applicability hinder their practical use. Besides, another content creation domain that can benefit from this approach is authoring audio descriptions for video content. While AI models like [63] can generate scene descriptions for video content in an end-to-end fashion, they fall short on accuracy, coherence, naturalness, and context awareness [36], which could be fixed by human validation and repair through a human-AI collaboration workflow [72]. An easy-to-use and effective human-AI collaboration framework like Mimosa would also make AI assistance in these domains of content creation more accessible to amateur creators.

## 7 LIMITATIONS & FUTURE WORK

The current version of Mimosa presents several technical and user study limitations for future work.

**Support for more general video types.** Due to the limited support for different topic domains in its sound separation, audio tagging, and object detection models, Mimosa currently supports only a limited range of sounding objects. Additionally, if a specific type of sounding object is not present in the training stage of the model, the model would struggle to recognize the visual object, separate its soundtrack, and predict its sound type. A promising approach to address this limitation is to incorporate an active learning approach that allows ML models to incrementally learn from the user's manipulation in real-time, thereby improving the model's performance.

**Model the interaction between the sound and environment.** Currently, the spatialization strategy used by Mimosa only considers the spatial position of the sounding object, without accounting for other factors such as reflection, absorption, and diffraction of sound. While new algorithms have been proposed for this problem [45, 46], we expect it to remain a challenge in the near future due to the level of complexity involved. These algorithms require extensive information about the scene's geometry and material properties, which may not be readily available or easily inferred from videos.

**Handle out-of-sight objects.** The current version of Mimosa is limited to handling sounding objects that are *visible* in the video frame. When the object detector is out of sight, the spatialization pipeline fails to work. To address this issue, Mimosa currently

---

[10]https://openai.com/dall-e-3/

employs linear interpolation for transient out-of-sight situations within video clips. In other cases, the system relies on users to manually specify the locations of sounding objects, which can be time-consuming and labor-intensive. While it is challenging to infer the location of objects that are never in sight, in future versions of Mimosa, we will explore the use of state-of-the-art computer vision techniques in trajectory prediction [25, 33, 57, 61] to predict in-between positions for objects that are *temporarily* out-of-sight due to their movements.

**Deployment study** As part of our future work, we plan to conduct a field deployment study involving content creators using Mimosa in their own video projects. Due to the constrained duration of the lab study, we did not systematically assess how users could integrate the tool into their actual creative workflows, such as incorporating Premiere Pro (Pr), although Mimosa supports an integrated workflow as a Pr plug-in. A deployment study will enable us to examine the long-term in-situ usage of Mimosa within its intended context. This will also allow us to assess the real-world effectiveness of the system for content creators and verify its ecological validity.

## 8 CONCLUSION

In this paper, we presented Mimosa, a human-AI co-creation tool that enabled amateur users to computationally generate and interactively manipulate spatial audio effects in videos that only had monaural or stereo audio. Mimosa featured a human-AI collaborative spatializing pipeline that produces user-interpretable and controllable intermediate results to support effective error discovery and repair. A controlled user study of Mimosa demonstrated that Mimosa's approach could support amateur content creators to generate immersive and realistic spatial audio and enable the flexible creation of customized spatial effects. Our findings provided design implications for AI-assisted content creation and future human-AI collaboration tools for working with multi-modal data.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Jessica J. Baldis. 2001. Effects of spatial audio on memory, comprehension, and preference during desktop conferences. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Seattle, Washington, USA) *(CHI '01)*. Association for Computing Machinery, New York, NY, USA, 166–173. https://doi.org/10.1145/365024.365092

[2] Virginia Braun and Victoria Clarke. 2006. Using thematic analysis in psychology. *Qualitative research in psychology* 3, 2 (2006), 77–101.

[3] Chakravarty R Alla Chaitanya, Nikunj Raghuvanshi, Keith W Godin, Zechen Zhang, Derek Nowrouzezahrai, and John M Snyder. 2020. Directional sources and listeners in interactive sound propagation using reciprocal wave field coding. *ACM Transactions on Graphics (TOG)* 39, 4 (2020), 44–1.

[4] Erin Cherry and Celine Latulipe. 2014. Quantifying the Creativity Support of Digital Tools through the Creativity Support Index. *ACM Trans. Comput.-Hum. Interact.* 21, 4, Article 21 (jun 2014), 25 pages. https://doi.org/10.1145/2617588

[5] John Joon Young Chung, Shiqing He, and Eytan Adar. 2021. The intersection of users, roles, interactions, and technologies in creativity support tools. In *Designing Interactive Systems Conference 2021*. 1817–1833.

[6] John Joon Young Chung, Shiqing He, and Eytan Adar. 2022. Artist Support Networks: Implications for Future Creativity Support Tools. In *Proceedings of the 2022 ACM Designing Interactive Systems Conference* (Virtual Event, Australia) *(DIS '22)*. Association for Computing Machinery, New York, NY, USA, 232–246. https://doi.org/10.1145/3532106.3533505

[7] H. Clark, G. Dutton, and P. Vanderlyn. 1957. The "Stereosonic" recording and reproducing system. *IRE Transactions on Audio* AU-5, 4 (1957), 96–111. https://doi.org/10.1109/TAU.1957.1166013

[8] Philip Coleman, Andreas Franck, Philip JB Jackson, Richard J Hughes, Luca Remaggi, and Frank Melchior. 2017. Object-based reverberation for spatial audio. *Journal of the Audio Engineering Society* 65, 1/2 (2017), 66–77.

[9] Robert Dalton, Jimmy Tobin, and David Grunzweig. 2016. Rondo360: Dysonics' Spatial Audio Post-Production Toolkit for 360 Media. In *Audio Engineering Society Convention 141*. https://www.aes.org/e-lib/browse.cfm?elib=18387

[10] Chris Dede. 2009. Immersive Interfaces for Engagement and Learning. *Science* 323, 5910 (2009), 66–69. https://doi.org/10.1126/science.1167311 arXiv:https://www.science.org/doi/pdf/10.1126/science.1167311

[11] Jonas Frich, Lindsay MacDonald Vermeulen, Christian Remy, Michael Mose Biskjaer, and Peter Dalsgaard. 2019. Mapping the landscape of creativity support tools in HCI. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. 1–18.

[12] Ruohan Gao and Kristen Grauman. 2018. 2.5D Visual Sound. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2018), 324–333. https://api.semanticscholar.org/CorpusID:54628402

[13] William G Gardner and Keith D Martin. 1995. HRTF measurements of a KEMAR. *The Journal of the Acoustical Society of America* 97, 6 (1995), 3907–3908.

[14] Simret Araya Gebreegziabher, Zheng Zhang, Xiaohang Tang, Yihao Meng, Elena Glassman, and Toby Jia-Jun Li. 2023. PaTAT: Human-AI Collaborative Qualitative Coding with Explainable Interactive Rule Synthesis. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems (CHI '23)*. ACM.

[15] Israel D Gebru, Dejan Marković, Alexander Richard, Steven Krenn, Gladstone A Butler, Fernando De la Torre, and Yaser Sheikh. 2021. Implicit hrtf modeling using temporal convolutional networks. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 3385–3389.

[16] Matthew Guzdial, Nicholas Liao, Jonathan Chen, Shao-Yu Chen, Shukan Shah, Vishwa Shah, Joshua Reno, Gillian Smith, and Mark O. Riedl. 2019. Friend, Collaborator, Student, Manager: How Design of an AI-Driven Game Level Editor Affects Creators. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (Glasgow, Scotland Uk) *(CHI '19)*. Association for Computing Machinery, New York, NY, USA, 1–13. https://doi.org/10.1145/3290605.3300854

[17] Florian Heller and Johannes Schöning. 2018. NavigaTone: seamlessly embedding navigation cues in mobile music listening. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. 1–7.

[18] Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Noam Shazeer, Curtis Hawthorne, AM Dai, MD Hoffman, and D Eck. 2018. Music transformer: Generating music with long-term structure (2018). *arXiv preprint arXiv:1809.04281* (2018).

[19] Gaurav Jain, Basel Hindi, Connor Courtien, Xin Yi Therese Xu, Conrad Wyrick, Michael Malcolm, and Brian A. Smith. 2023. Front Row: Automatically Generating Immersive Audio Representations of Tennis Broadcasts for Blind Viewers. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology* (San Francisco, CA, USA) *(UIST '23)*. Association for Computing Machinery, New York, NY, USA, Article 39, 17 pages. https://doi.org/10.1145/3586183.3606830

[20] Doyeon Kim, Woonghyun Ga, Pyungwhan Ahn, Donggyu Joo, Sehwan Chun, and Junmo Kim. 2022. Global-Local Path Networks for Monocular Depth Estimation with Vertical CutDepth. *arXiv preprint arXiv:2201.07436* (2022).

[21] Qiuqiang Kong, Yin Cao, Turab Iqbal, Yuxuan Wang, Wenwu Wang, and Mark D. Plumbley. 2020. PANNs: Large-Scale Pretrained Audio Neural Networks for Audio Pattern Recognition. arXiv:1912.10211 [cs.SD]

[22] Janusz Konrad, Meng Wang, Prakash Ishwar, Chen Wu, and Debargha Mukherjee. 2013. Learning-based, automatic 2D-to-3D image and video conversion. *IEEE Transactions on Image Processing* 22, 9 (2013), 3485–3496.

[23] Simon Langford. 2013. *Digital audio editing: correcting and enhancing audio in Pro Tools, Logic Pro, Cubase, and Studio One*. Routledge. https://doi.org/10.4324/9780203512890

[24] Jonathan Lazar, Jinjuan Heidi Feng, and Harry Hochheiser. 2017. Preface. In *Research Methods in Human Computer Interaction (Second Edition)* (second edition ed.), Jonathan Lazar, Jinjuan Heidi Feng, and Harry Hochheiser (Eds.). Morgan Kaufmann, Boston. https://doi.org/10.1016/B978-0-12-805390-4.09987-8

[25] Junwei Liang, Lu Jiang, and Alexander Hauptmann. 2020. Simaug: Learning robust representations from simulation for trajectory prediction. In *European Conference on Computer Vision*. Springer, 275–292.

[26] David Chuan-En Lin, Anastasis Germanidis, Cristóbal Valenzuela, Yining Shi, and Nikolas Martelaro. 2023. Soundify: Matching Sound Effects to Video. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology* (San Francisco, CA, USA) *(UIST '23)*. Association for Computing Machinery, New York, NY, USA, Article 18, 13 pages. https://doi.org/10.1145/3586183.3606823

[27] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *European conference on computer vision*. Springer, 740–755.

[28] Yuyu Lin, Jiahao Guo, Yang Chen, Cheng Yao, and Fangtian Ying. 2020. It is your turn: collaborative ideation with a co-creative robot through sketch. In *Proceedings of the 2020 CHI conference on human factors in computing systems*. 1–14.

[29] Ryan Louie, Andy Coenen, Cheng Zhi Huang, Michael Terry, and Carrie J Cai. 2020. Novice-AI music co-creation via AI-steering tools for deep generative models. In *Proceedings of the 2020 CHI conference on human factors in computing systems*. 1–13.

[30] Leo McCormack and Archontis Politis. 2019. SPARTA & COMPASS: Real-time implementations of linear and parametric spatial audio reproduction and processing methods. In *Audio Engineering Society Conference: 2019 AES International Conference on Immersive and Interactive Audio*. Audio Engineering Society.

[31] Logan Middleton. 2022. Mix It Up, Mash It Up: Arrangement, Audio Editing, and the Importance of Sonic Context. (2022), 29. https://doi.org/10.37514/PRA-B.2022.1688.2.01

[32] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. 2021. Nerf: Representing scenes as neural radiance fields for view synthesis. *Commun. ACM* 65, 1 (2021), 99–106.

[33] Abduallah Mohamed, Deyao Zhu, Warren Vu, Mohamed Elhoseiny, and Christian Claudel. 2022. Social-Implicit: Rethinking Trajectory Prediction Evaluation and The Effectiveness of Implicit Maximum Likelihood Estimation. *arXiv preprint arXiv:2203.03057* (2022).

[34] Henrik Møller, Michael Friis Sørensen, Dorte Hammershøi, and Clemen Boje Jensen. 1995. Head-related transfer functions of human subjects. *Journal of the Audio Engineering Society* 43, 5 (1995), 300–321.

[35] Pedro Morgado, Nuno Nvasconcelos, Timothy Langlois, and Oliver Wang. 2018. Self-supervised generation of spatial audio for 360 video. *Advances in neural information processing systems* 31 (2018).

[36] Rosiana Natalie, Jolene Loh, Huei Suen Tan, Joshua Tseng, Hernisa Kacorri, and Kotaro Hara. 2021. Uncovering Patterns in Reviewers' Feedback to Scene Description Authors. In *Proceedings of the 23rd International ACM SIGACCESS Conference on Computers and Accessibility* (Virtual Event, USA) (*ASSETS '21*). Association for Computing Machinery, New York, NY, USA, Article 93, 4 pages. https://doi.org/10.1145/3441852.3476550

[37] Cuong Nguyen, Stephen DiVerdi, Aaron Hertzmann, and Feng Liu. 2017. CollaVR: collaborative in-headset review for VR video. In *Proceedings of the 30th Annual ACM Symposium on User Interface Software and Technology*. 267–277.

[38] Jakob Nielsen and Rolf Molich. 1990. Heuristic evaluation of user interfaces. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Seattle, Washington, USA) (*CHI '90*). Association for Computing Machinery, New York, NY, USA, 249–256. https://doi.org/10.1145/97243.97281

[39] Zheng Ning, Brianna L. Wimer, Kaiwen Jiang, Keyi Chen, Jerrick Ban, Yapeng Tian, Yuhang Zhao, and Toby Jia-Jun Li. 2024. SPICA: Interactive Video Content Exploration through Augmented Audio Descriptions for Blind or Low-Vision Viewers. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*.

[40] Zheng Ning, Zheng Zhang, Tianyi Sun, Yuan Tian, Tianyi Zhang, and Toby Jia-Jun Li. 2023. An Empirical Study of Model Errors and User Error Discovery and Repair Strategies in Natural Language Database Queries. In *Proceedings of the 28th International Conference on Intelligent User Interfaces* (Sydney, NSW, Australia) (*IUI '23*). Association for Computing Machinery, New York, NY, USA, 633–649. https://doi.org/10.1145/3581641.3584067

[41] Adrian North and David Hargreaves. 2008. *The Social and Applied Psychology of Music*. Oxford University Press. https://doi.org/10.1093/acprof:oso/9780198567424.001.0001

[42] Changhoon Oh, Jungwoo Song, Jinhan Choi, Seonghyeon Kim, Sungwoo Lee, and Bongwon Suh. 2018. I Lead, You Help but Only with Enough Details: Understanding User Experience of Co-Creation with Artificial Intelligence. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems* (Montreal QC, Canada) (*CHI '18*). Association for Computing Machinery, New York, NY, USA, 1–13. https://doi.org/10.1145/3173574.3174223

[43] Sharon Oviatt. 1999. Mutual Disambiguation of Recognition Errors in a Multi-model Architecture. In *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (Pittsburgh, Pennsylvania, USA) (*CHI '99*). Association for Computing Machinery, New York, NY, USA, 576–583. https://doi.org/10.1145/302979.303163

[44] Sharon Oviatt. 1999. Ten Myths of Multimodal Interaction. *Commun. ACM* 42, 11 (nov 1999), 74–81. https://doi.org/10.1145/319382.319398

[45] Nikunj Raghuvanshi and John Snyder. 2014. Parametric wave field coding for precomputed sound propagation. *ACM Transactions on Graphics (TOG)* 33, 4 (2014), 1–11.

[46] Nikunj Raghuvanshi and John Snyder. 2018. Parametric directional coding for precomputed sound propagation. *ACM Transactions on Graphics (TOG)* 37, 4 (2018), 1–14.

[47] Summer Rebensky, Kendall Carmody, Cherrise Ficke, Daniel Nguyen, Meredith Carroll, Jessica Wildman, and Amanda Thayer. 2021. Whoops! Something went wrong: Errors, trust, and trust repair strategies in human agent teaming. In

[48] Scott Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. 2016. Generative adversarial text to image synthesis. In *International conference on machine learning*. PMLR, 1060–1069.

[49] Christian Remy, Lindsay MacDonald Vermeulen, Jonas Frich, Michael Mose Biskjaer, and Peter Dalsgaard. 2020. Evaluating creativity support tools in HCI research. In *Proceedings of the 2020 ACM designing interactive systems conference*. 457–476.

[50] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. 2015. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems* 28 (2015), 91–99.

[51] Alexander Richard, Dejan Markovic, Israel D. Gebru, Steven Krenn, Gladstone Alexander Butler, Fernando Torre, and Yaser Sheikh. 2021. Neural Synthesis of Binaural Speech From Mono Audio. In *International Conference on Learning Representations*. https://openreview.net/forum?id=uAX8q61EVRu

[52] Agnieszka Roginska and Paul Geluso. 2017. *Immersive sound: The art and science of binaural and multi-channel audio*. https://doi.org/10.4324/9781315707525

[53] Joshua S Rubinstein, David E Meyer, and Jeffrey E Evans. 2001. Executive control of cognitive processes in task switching. *Journal of experimental psychology: human perception and performance* 27, 4 (2001), 763.

[54] Francis Rumsey. 2012. *Spatial audio*. Routledge. https://doi.org/10.4324/9780080498195

[55] Francis Rumsey and Tim McCormick. 2021. *Sound and recording: applications and theory*. Routledge. https://doi.org/10.4324/9781003092919

[56] Jaime Sánchez and Héctor Flores. 2003. Memory enhancement through audio. In *Proceedings of the 6th international ACM SIGACCESS conference on Computers and accessibility*. 24–31.

[57] Adam Ścibior, Vasileios Lioutas, Daniele Reda, Peyman Bateni, and Frank Wood. 2021. Imagining the road ahead: Multi-agent trajectory prediction via differentiable simulation. In *2021 IEEE International Intelligent Transportation Systems Conference (ITSC)*. IEEE, 720–725.

[58] Ian Simon, Dan Morris, and Sumit Basu. 2008. MySong: automatic accompaniment generation for vocal melodies. In *Proceedings of the SIGCHI conference on human factors in computing systems*. 725–734.

[59] Paris Smaragdis. 2009. User guided audio selection from complex sound mixtures. In *Proceedings of the 22nd Annual ACM Symposium on User Interface Software and Technology* (Victoria, BC, Canada) (*UIST '09*). Association for Computing Machinery, New York, NY, USA, 89–92. https://doi.org/10.1145/1622176.1622193

[60] Yapeng Tian, Di Hu, and Chenliang Xu. 2021. Cyclic co-learning of sounding object visual grounding and sound separation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2745–2754.

[61] Chuhua Wang, Yuchen Wang, Mingze Xu, and David J Crandall. 2022. Stepwise goal-driven networks for trajectory prediction. *IEEE Robotics and Automation Letters* 7, 2 (2022), 2716–2723.

[62] Dakuo Wang, Elizabeth Churchill, Pattie Maes, Xiangmin Fan, Ben Shneiderman, Yuanchun Shi, and Qianying Wang. 2020. From human-human collaboration to Human-AI collaboration: Designing AI systems that can work together with people. In *Extended abstracts of the 2020 CHI conference on human factors in computing systems*. 1–6.

[63] Yujia Wang, Wei Liang, Haikun Huang, Yongqi Zhang, Dingzeyu Li, and Lap-Fai Yu. 2021. Toward Automatic Audio Description Generation for Accessible Videos. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems* (Yokohama, Japan) (*CHI '21*). Association for Computing Machinery, New York, NY, USA, Article 277, 12 pages. https://doi.org/10.1145/3411764.3445347

[64] Yake Wei, Di Hu, Yapeng Tian, and Xuelong Li. 2022. Learning in Audio-visual Context: A Review, Analysis, and New Perspective. *CoRR* abs/2208.09579 (2022). https://doi.org/10.48550/ARXIV.2208.09579 arXiv:2208.09579

[65] Scott Wisdom, Hakan Erdogan, Daniel PW Ellis, Romain Serizel, Nicolas Turpault, Eduardo Fonseca, Justin Salamon, Prem Seetharaman, and John R Hershey. 2021. What's all the fuss about free universal sound separation data?. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 186–190.

[66] Junyuan Xie, Ross Girshick, and Ali Farhadi. 2016. Deep3d: Fully automatic 2d-to-3d video conversion with deep convolutional neural networks. In *European conference on computer vision*. Springer, 842–857.

[67] Xudong Xu, Hang Zhou, Ziwei Liu, Bo Dai, Xiaogang Wang, and Dahua Lin. 2021. Visually informed binaural audio generation without binaural audios. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 15485–15494.

[68] Kazuhiko Yamamoto and Takeo Igarashi. 2017. Fully perceptual-based 3D spatial sound individualization with an adaptive variational autoencoder. *ACM Transactions on Graphics (TOG)* 36, 6 (2017), 1–13.

[69] Jing Yang and Friedemann Mattern. 2019. Audio Augmented Reality for Human-Object Interactions. In *Adjunct Proceedings of the 2019 ACM International Joint Conference on Pervasive and Ubiquitous Computing and Proceedings of the 2019 ACM International Symposium on Wearable Computers* (London, United Kingdom) (*UbiComp/ISWC '19 Adjunct*). Association for Computing Machinery, New York, NY, USA, 408–412. https://doi.org/10.1145/3341162.3349302

[70] Qian Yang, Aaron Steinfeld, Carolyn Rosé, and John Zimmerman. 2020. Re-Examining Whether, Why, and How Human-AI Interaction Is Uniquely Difficult to Design. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (Honolulu, HI, USA) *(CHI '20)*. Association for Computing Machinery, New York, NY, USA, 1–13. https://doi.org/10.1145/3313831.3376301

[71] Seraphina Yong and Hao-Chuan Wang. 2018. Using spatialized audio to improve human spatial knowledge acquisition in virtual reality. In *Proceedings of the 23rd International Conference on Intelligent User Interfaces Companion*. 1–2.

[72] Beste F. Yuksel, Pooyan Fazli, Umang Mathur, Vaishali Bisht, Soo Jung Kim, Joshua Junhee Lee, Seung Jung Jin, Yue-Ting Siu, Joshua A. Miele, and Ilmi Yoon. 2020. Human-in-the-Loop Machine Learning to Increase Video Accessibility for Visually Impaired and Blind Users. In *Proceedings of the 2020 ACM Designing Interactive Systems Conference* (Eindhoven, Netherlands) *(DIS '20)*. Association for Computing Machinery, New York, NY, USA, 47–60. https://doi.org/10.1145/3357236.3395433

[73] Mengfan Zhang, Jui-Hsien Wang, and Doug L James. 2021. Personalized HRTF Modeling Using DNN-Augmented BEM. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 451–455.

[74] Zheng Zhang, Jie Gao, Ranjodh Singh Dhaliwal, and Toby Jia-Jun Li. 2023. VISAR: A Human-AI Argumentative Writing Assistant with Visual Programming and Rapid Draft Prototyping. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology* (San Francisco, CA, USA) *(UIST '23)*. Association for Computing Machinery, New York, NY, USA, Article 5, 30 pages. https://doi.org/10.1145/3586183.3606800

[75] Zheng Zhang, Zheng Ning, Chenliang Xu, Yapeng Tian, and Toby Jia-Jun Li. 2023. PEANUT: A Human-AI Collaborative Tool for Annotating Audio-Visual Data. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology* (San Francisco,CA,USA) *(UIST '23)*. Association for Computing Machinery, New York, NY, USA, Article 17, 18 pages. https://doi.org/10.1145/3586183.3606776

[76] Hang Zhou, Xudong Xu, Dahua Lin, Xiaogang Wang, and Ziwei Liu. 2020. Sep-stereo: Visually guided stereophonic audio generation by associating source separation. In *European Conference on Computer Vision*. Springer, 52–69.