# SonicVisionLM: Playing Sound with Vision Language Models
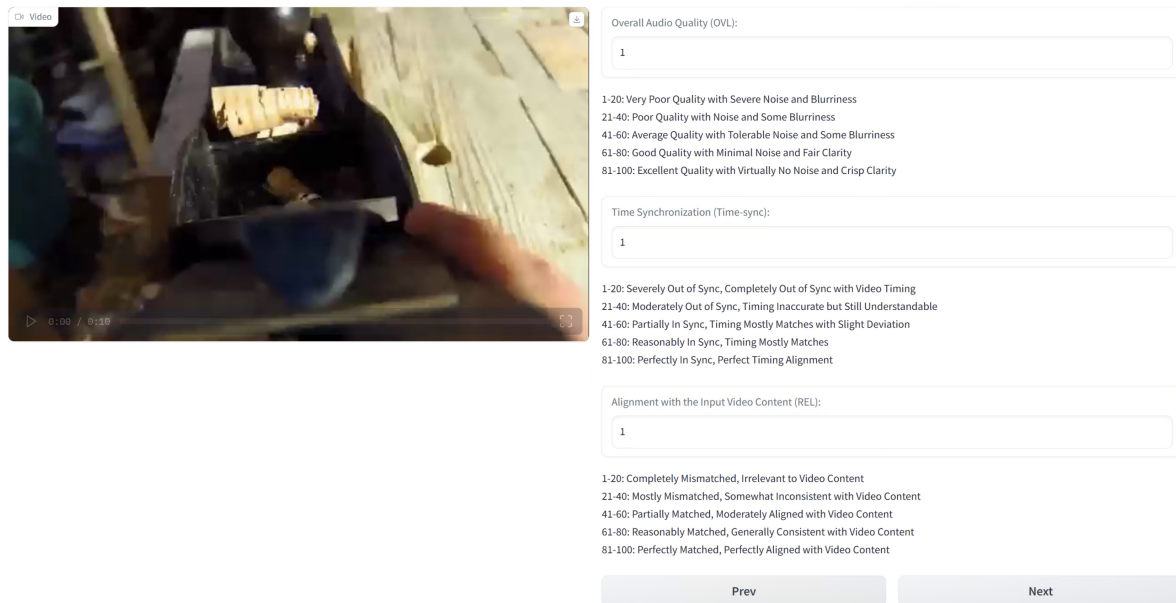
## Supplementary Material



Figure 1. The evaluation page for the unconditional generation task. We display a screenshot of the primary test interface that participants will encounter. Each participant is required to input three scores for the presented content. Upon clicking the 'Next' button, they will be directed to the subsequent video in the test sequence.

## A.1. Subjective Results Details

We provide the screenshot of the main evaluation page the participant will see during the test in Fig.1.

## A.2. Timestamp Detection Module Precision Experiment

Our timestamp detection model shows promising performance, with an Average Precision(AP) of 0.80 and an accuracy of 0.72 on the Greatest Hits test set, and even higher results on the validation set with an AP of 0.92 and accuracy of 0.82. However, its performance on the CountixAV test set is comparatively lower, achieving an AP of 0.52 and an accuracy of 0.52. This discrepancy likely stems from the complexity of sound sources in our training dataset, leading to potential inaccuracies in ground truth. Such labelling challenges can adversely affect recognition accuracy, particularly in scenarios involving non-static footage.

## A.3. Additional Results

**Conditional Generation Results.** As shown in Fig.2, the left column represents the target video, the middle column showcases the control conditions for CondFoleyGen and
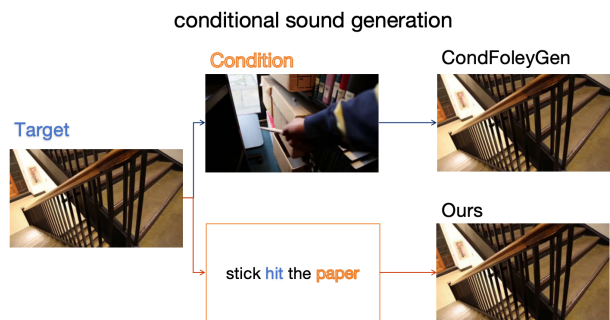


Figure 2. screenshot of the conditional sound generation task section in the demo video.

our model, and the right column displays the generated results. We provide 6 examples for the previously mentioned conditional generation task, and the corresponding audio outcomes can be observed in the demo video.

**Unconditional Generation Results.** As shown in Fig.3, we compare a comparison of the results generated by GT, SpecVQGAN, DIFF-FOLEY, and our model. We provide 6 examples from CountixAV, and the audio samples are available in the demo video.
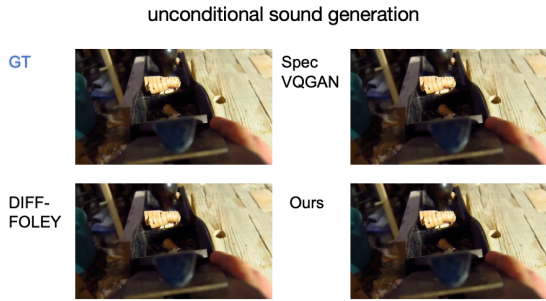
Figure 3. screenshot of the unconditional sound generation task section in the demo video.
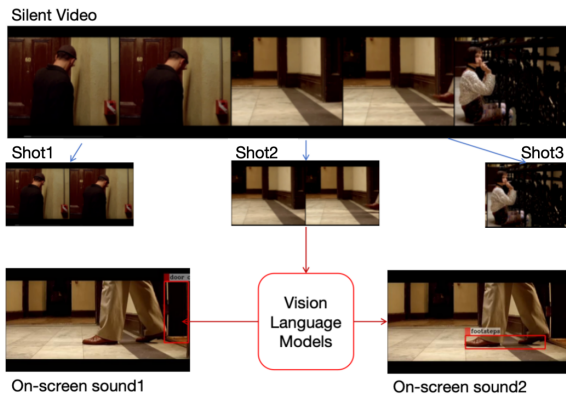


Figure 4. The process of video-to-text and text-based interaction components in multi-track generation tasks.

**Multi-soundtracks Generation Results.** As shown in Fig.4, we segment the video into different shots first. Shots with clear actions are directly processed by the VLMs to obtain corresponding sound effect descriptions and their spatial positioning within the video. These descriptions are used to generate on-screen sound. Other shots accept user editing and are then fed into the LDM to produce off-screen sounds. We provide a simple and a complex case respectively in the demo video.