# Conditional Generation of Audio from Video via Foley Analogies

Yuexi Du[1,2]  Ziyang Chen[1]  Justin Salamon[3]  Bryan Russell[3]  Andrew Owens[1]

[1]University of Michigan  [2]Yale University  [3]Adobe Research

## Introduction & Goal

*Make silent input video sounds like the conditional video*



Input silent video

Conditional video

Generated audio

## Pretext task: Foley analogies

**Idea**: Natural videos tend to contain **repeated events** that produce closely related sounds.



Given a *conditional audio-visual clip*

Predict the audio for *silent input clip*

Audio to be generated

Time

Foley analogies model $\mathcal{F}_\theta$

Silent input video: $\mathbf{v}_q$
Conditional audiovisual clip: $(\mathbf{a}_c, \mathbf{v}_c)$ } $\xrightarrow{\mathcal{F}_\theta}$ Output audio: $\mathcal{F}_\theta(\mathbf{v}_q, \mathbf{a}_c, \mathbf{v}_c)$

Learning Foley analogies:

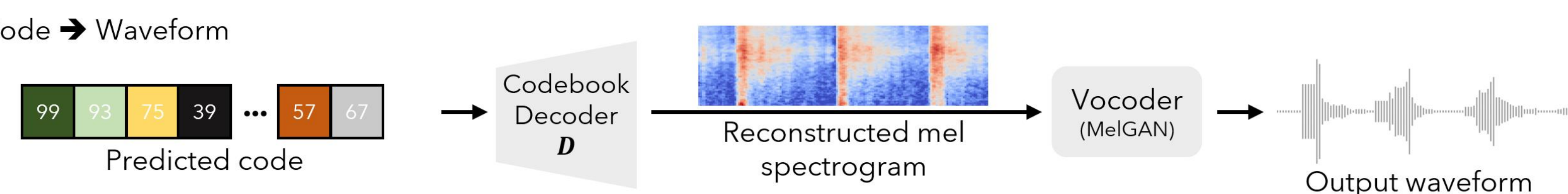$$\arg\min_\theta \mathcal{L}\left(\mathbf{a}_q, \mathcal{F}_\theta(\mathbf{v}_q, \mathbf{a}_c, \mathbf{v}_c)\right)$$

## Conditional audio prediction model

**Idea:** Predicted audio tokens auto-regressively with a GPT-2 transformer inspired by Spec-VQGAN (Iashin *et.al*, 2022) with silent input video and conditional audio-visual clip.
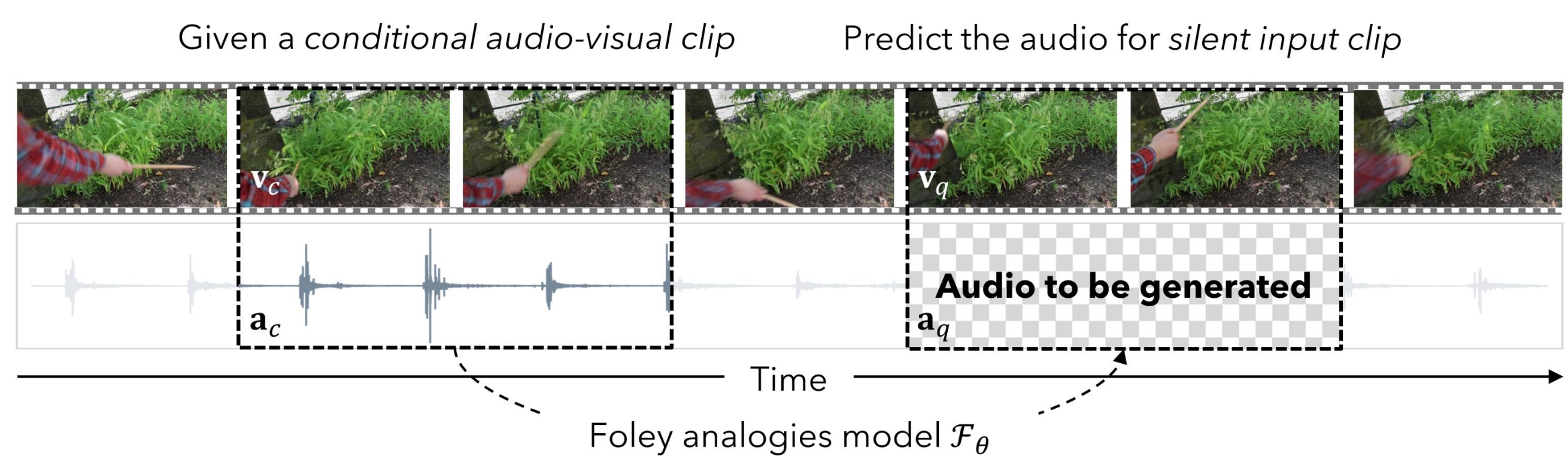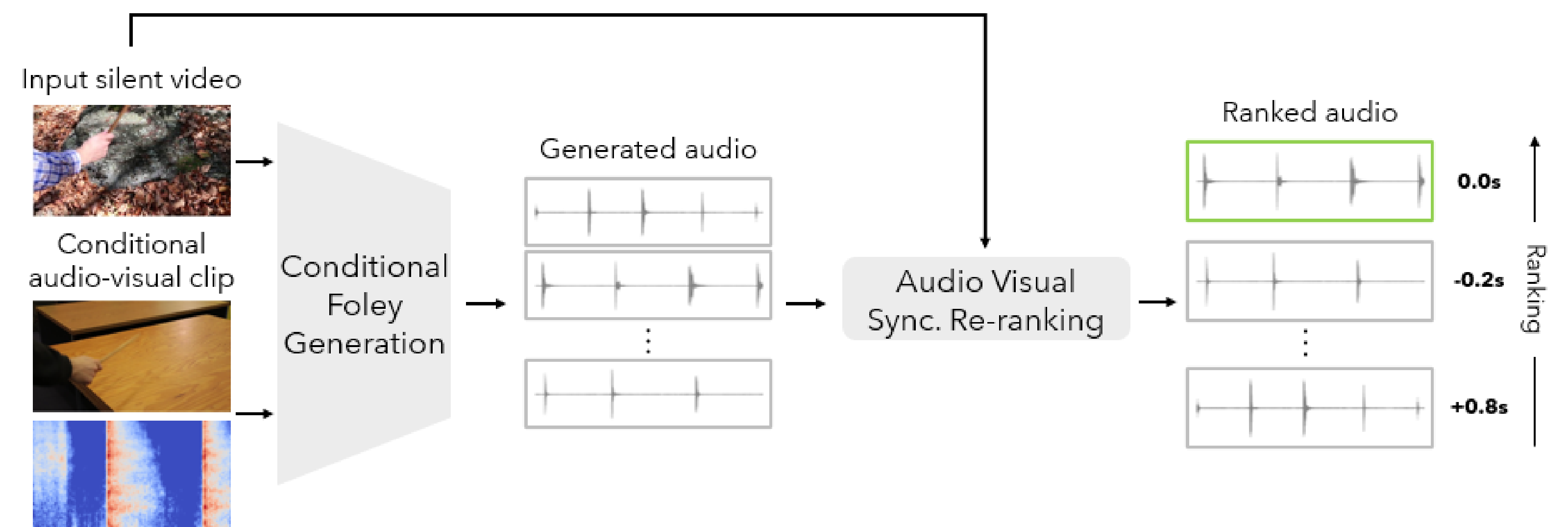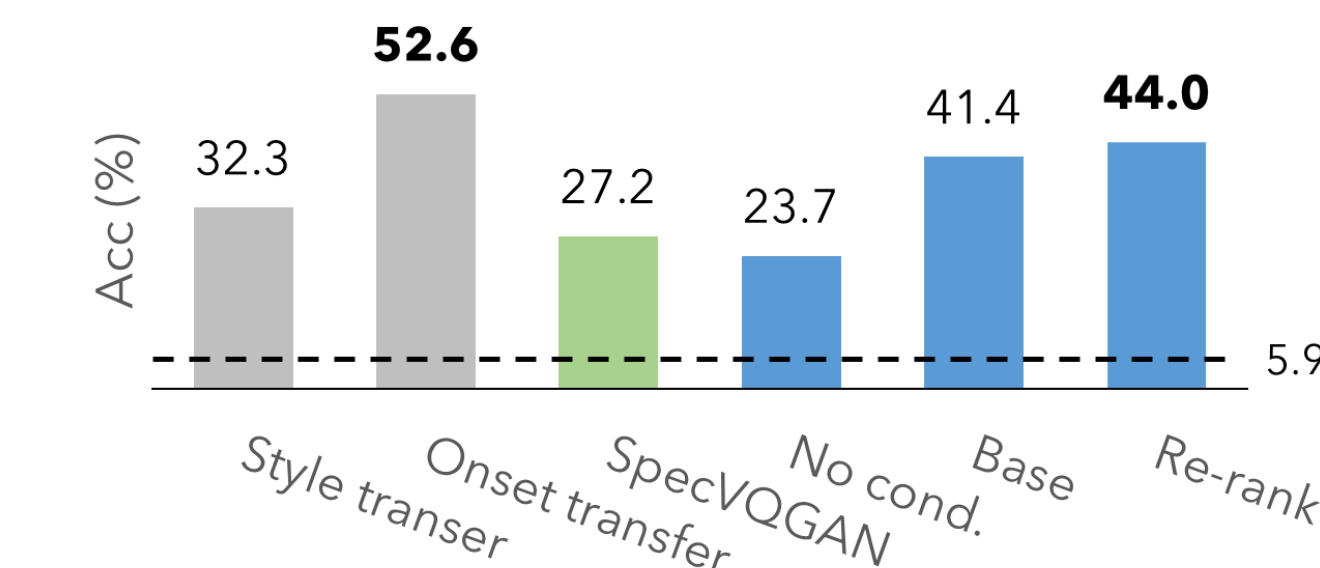
**A.** Inputs ➤ Code



Input video

Video Encoder

Input video $T_v(\mathbf{v}_q)$

Cond. video

Video Encoder

Cond. video $T_v(\mathbf{v}_c)$

Cond. mel spec.

Codebook Encoder $E$

Cond. audio $T_a(\mathbf{a}_c)$

Predicted code $\mathbf{s}_{<i}$

Tokenize

Transformer $p_\theta$

Ground truth code $\hat{\mathbf{s}}$

$\mathcal{L}_{CE}$

Output pred. code $\mathbf{s}$

**B.** Code ➤ Waveform



Predicted code

Codebook Decoder $D$

Reconstructed mel spectrogram

Vocoder (MelGAN)

Output waveform

## Inference-time audio re-ranking



Input silent video

Conditional audio-visual clip

Conditional Foley Generation

Generated audio

Audio Visual Sync. Re-ranking
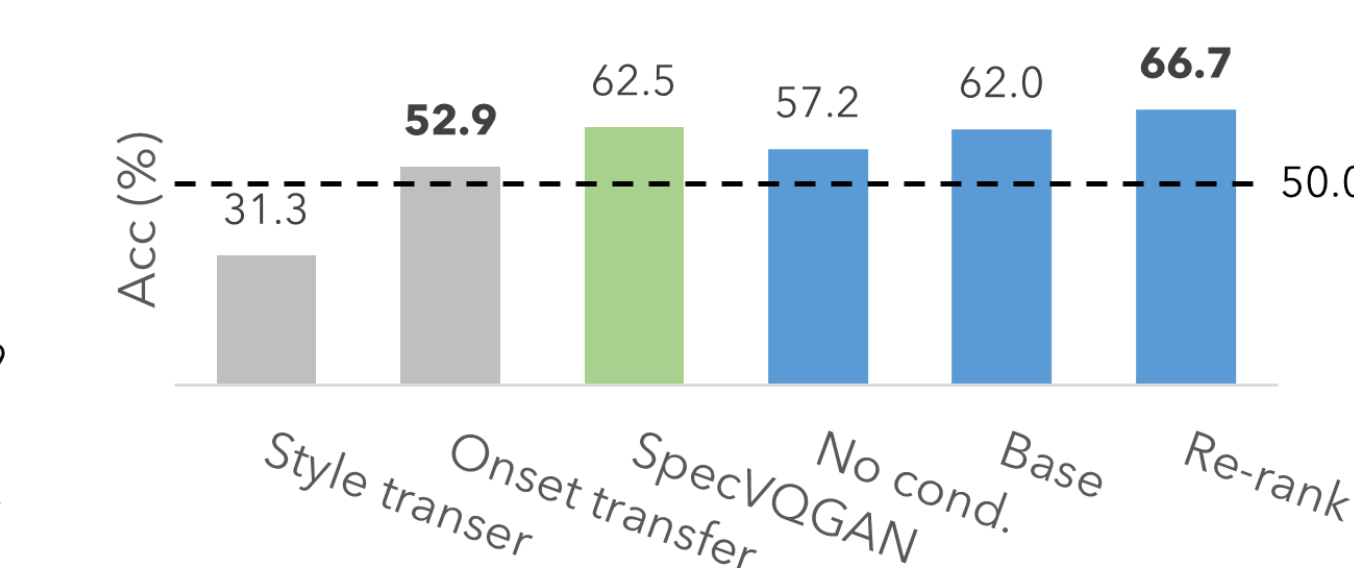
Ranked audio

0.0s
-0.2s
+0.8s
Ranking

We re-rank multiple generated audio according to their **temporal alignment** with the input video predicted by off-the-shelf audio-visual synchronization model.
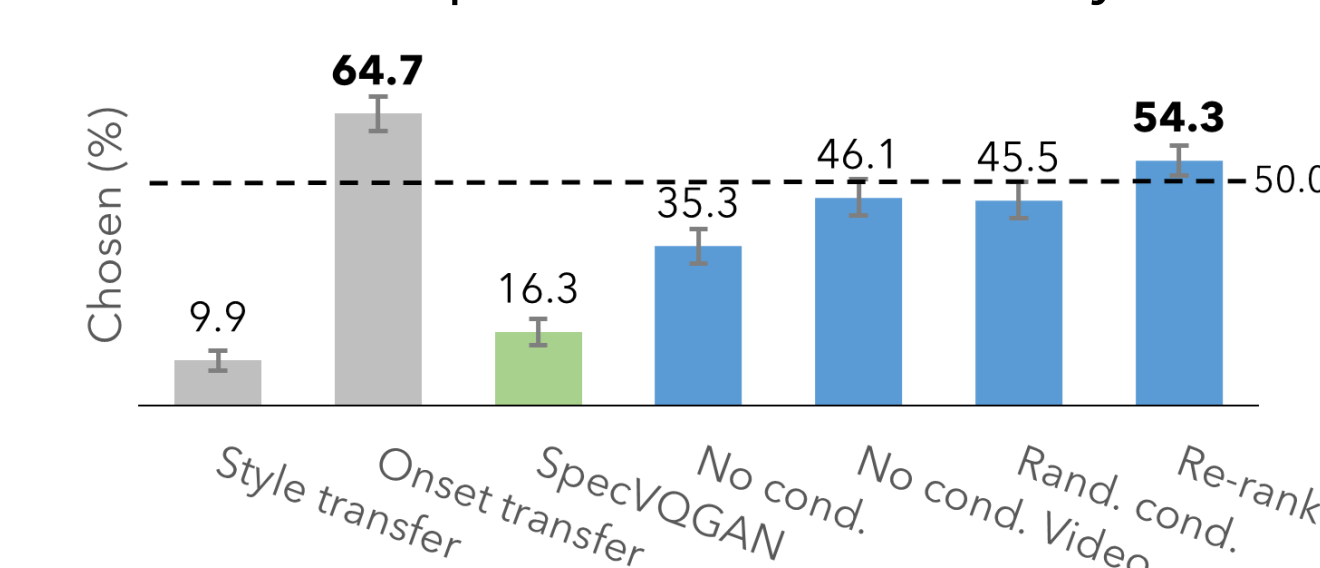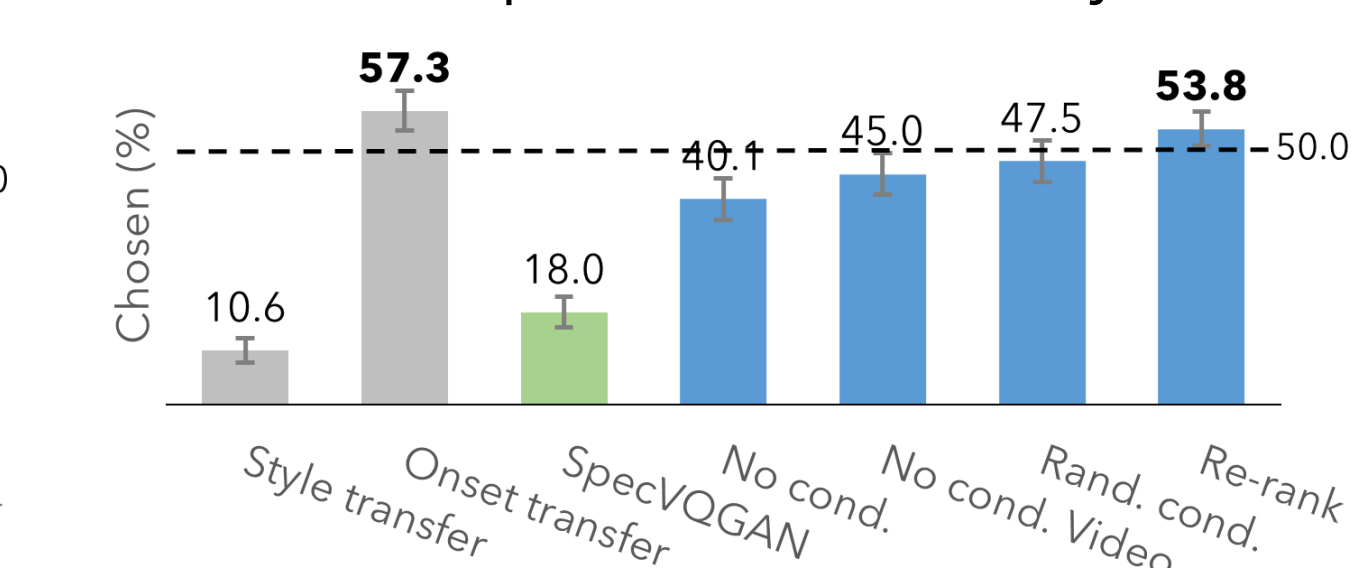
## Experiments

Non-generative Baseline
Generative Baseline
Ours



Automated Material Prediction

Acc (%)
32.3  52.6  27.2  23.7  41.4  44.0
Style transfer  Onset transfer  SpecVQGAN  No cond.  Base  Re-rank
5.9

Automated Action Prediction

Acc (%)
31.3  52.9  62.5  57.2  62.0  66.7
Style transfer  Onset transfer  SpecVQGAN  No cond.  Base  Re-rank
50.0

Perceptual Material Study

Chosen (%)
9.9  64.7  16.3  35.3  46.1  45.5  54.3
Style transfer  Onset transfer  SpecVQGAN  No cond.  No cond. Video  Rand. cond.  Re-rank
50.0

Perceptual Action Study

Chosen (%)
10.6  57.3  18.0  40.1  45.0  47.5  53.8
Style transfer  Onset transfer  SpecVQGAN  No cond.  No cond. Video  Rand. cond.  Re-rank
50.0

### Qualitative results (*Please check our website for video results*)



Input silent video

Conditional audio-visual example

Generated audio

Input sample's audio

In-the-wild Videos

Greatest Hits

Input silent video

Conditional audio-visual clips

Generated audio

Input sample's video