

# **Audio-Visual Speech Codecs: Rethinking Audio-Visual Speech Enhancement by Re-Synthesis**

Karren Yang<sup>1</sup> Dejan Marković<sup>2</sup> Steven Krenn<sup>2</sup> Vasu Agrawal<sup>2</sup> Alexander Richard<sup>2</sup>  
<sup>1</sup>MIT    <sup>2</sup>Meta Reality Labs Research

karren@mit.edu    {dejanmarkovic, stevenkrenn, vasuagrawal, richardalex}@fb.com

CVPR 2022, <https://arxiv.org/pdf/2203.17263.pdf>

# Introduction

- Denoising Speech Signal

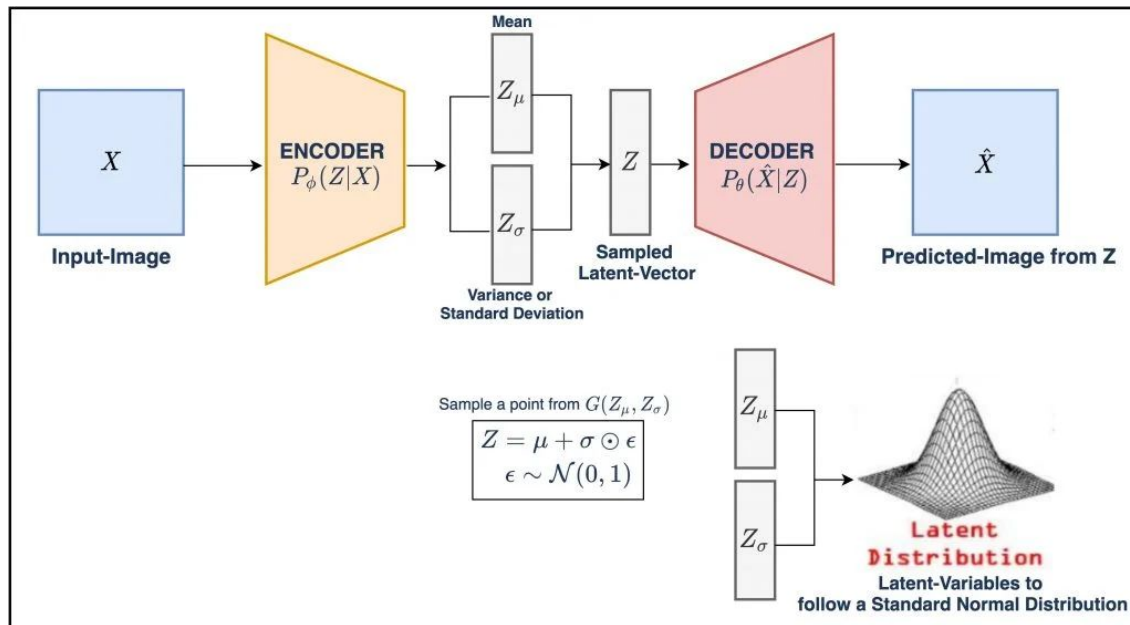
$$x(t) + n(t) = y(t)$$



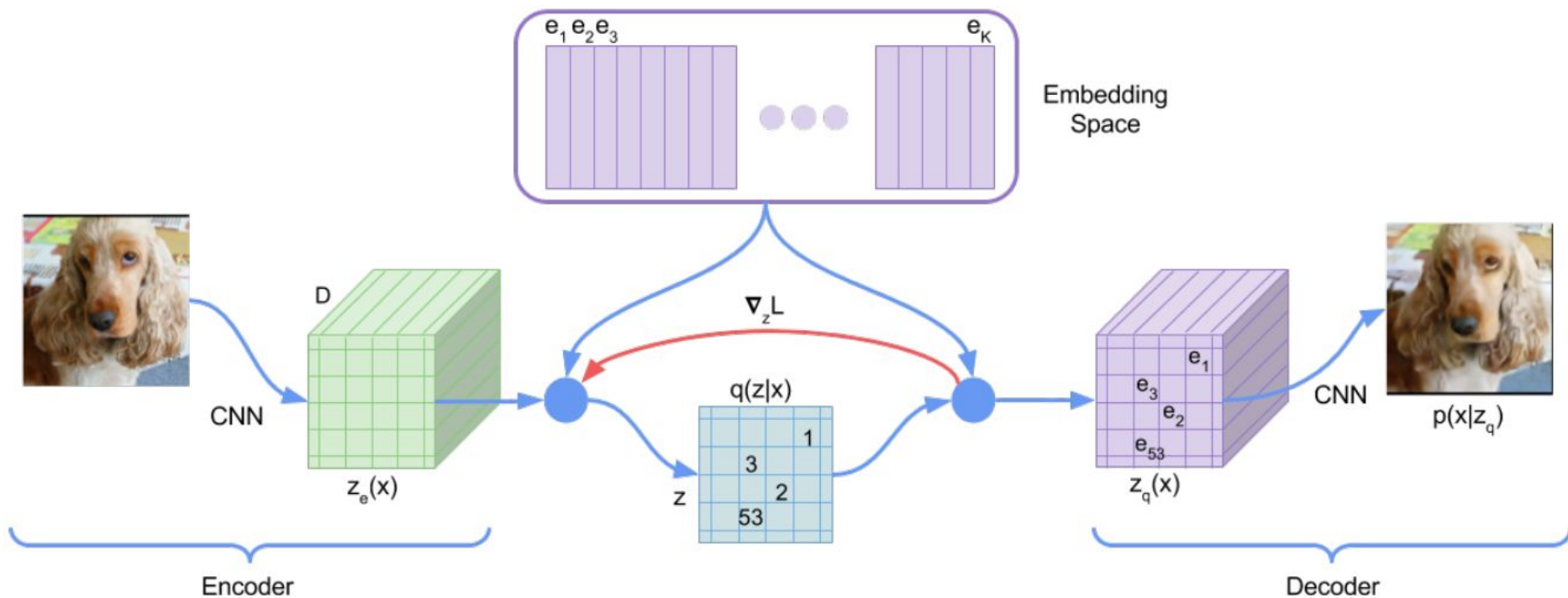
- Using Visual Information as auxiliary signal
- Used Generative approach (VAE)

# Approach

- Used Vector Quantized VAE (VQ-VAE)
- Discrete Representation

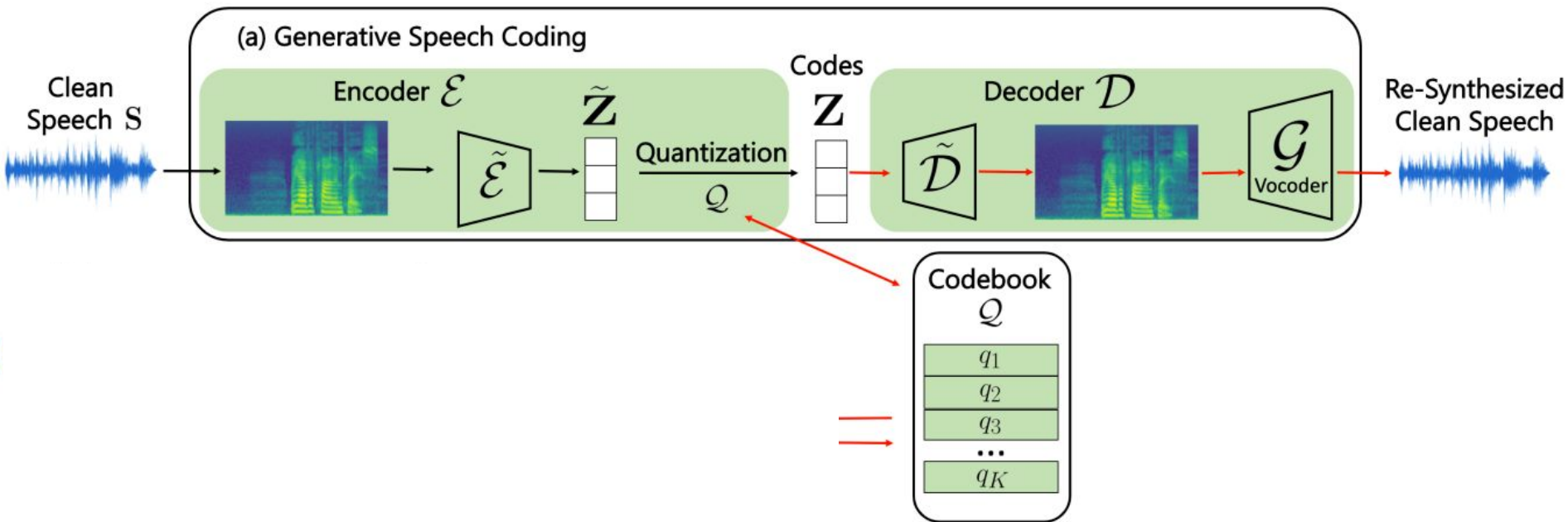


# VQ-VAE



# Approach

$$\mathbb{E}_{(\mathbf{S}, \mathbf{S}', \mathbf{V})} \|\mathbf{melspec}(\mathbf{S}) - \tilde{\mathcal{D}}(\mathbf{Z})\|_2^2$$



# Dataset

## FaceStar

- Audio-visual dataset containing 10 hours of speech data from two speakers

| <b>AV Dataset</b> | <b># Hours per Speaker</b> | <b>High-Quality Audio</b> | <b>Reliable Lip Motion</b> | <b>Unconstrained Natural Speech</b> |
|-------------------|----------------------------|---------------------------|----------------------------|-------------------------------------|
| <b>GRID</b>       | 0.8                        | ✓                         | ✓                          | ✗                                   |
| <b>TCD-Timit</b>  | 0.5                        | ✓                         | ✓                          | ✗                                   |
| <b>Lip2Wav</b>    | 20                         | ✗                         | ✗                          | ✓                                   |
| <b>Facestar</b>   | 5                          | ✓                         | ✓                          | ✓                                   |

# Results

| Model           | Facestar        |                 |                  |                  |                            | Lip2Wav         |                 |                  |                  |                            |
|-----------------|-----------------|-----------------|------------------|------------------|----------------------------|-----------------|-----------------|------------------|------------------|----------------------------|
|                 | PESQ $\uparrow$ | STOI $\uparrow$ | F-SNR $\uparrow$ | MCD $\downarrow$ | Mel- $\ell_2$ $\downarrow$ | PESQ $\uparrow$ | STOI $\uparrow$ | F-SNR $\uparrow$ | MCD $\downarrow$ | Mel- $\ell_2$ $\downarrow$ |
| Demucs [8]      | 1.251           | 0.554           | 5.602            | 5.003            | 0.0106                     | 1.383           | 0.672           | 7.644            | 4.724            | 0.0109                     |
| AV-Masking [18] | 1.257           | 0.593           | 5.991            | 5.184            | 0.0093                     | 1.438           | 0.689           | 7.873            | 5.167            | 0.0093                     |
| AV-Mapping [15] | 1.332           | 0.626           | 2.802            | 4.885            | 0.0059                     | 1.417           | 0.661           | 6.892            | 4.643            | 0.0062                     |
| Ours            | <b>1.354</b>    | <b>0.661</b>    | <b>7.322</b>     | <b>3.815</b>     | <b>0.0056</b>              | <b>1.482</b>    | <b>0.740</b>    | <b>8.801</b>     | <b>4.072</b>     | <b>0.0055</b>              |

| Ours  | GT recordings      | Can not tell |
|-------|--------------------|--------------|
| 4.1%  | 44.5%              | 51.4%        |
| Ours  | AV Encoder Decoder | Can not tell |
| 73.3% | 6.0%               | 20.7%        |
| Ours  | AV Masking         | Can not tell |
| 78.5% | 5.7%               | 15.8%        |

Table 3. **Perceptual Evaluation.** Participants were presented two video clips and asked to tell which of the two sounds more natural.

| Model                     | only                                 |                |
|---------------------------|--------------------------------------|----------------|
|                           | reverb + noise<br>+ interfering spkr | reverb + noise |
| Vision-Only               | 0.0085                               |                |
| Audio-Only                | 0.0091                               | 0.0056         |
| No Auto-Regressive Module | 0.0051                               | 0.0036         |
| Full Model                | <b>0.0043</b>                        | <b>0.0033</b>  |

Table 4. **Ablation Results.** The values shown are the mean  $\ell_2$  errors between predicted and ground truth mel-spectrograms for ablation models trained on the Facestar dataset (Speaker 1); lower is better. See text for details.

# Multi-Speaker Model

- Extend framework to multispeaker setting by adding a speaker identity encoder

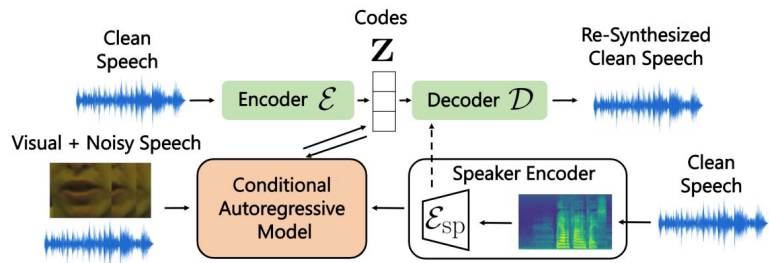


Figure 4. **Multi-speaker model**. A speaker encoder is added to the pipeline from Figure 2. Restricting the size of the codebook forces the model to disentangle speech content and speaker identity as shown in [48].

|  | GRID Speaker |           |            |            |
|--|--------------|-----------|------------|------------|
|  | Sp. 1 (M)    | Sp. 3 (M) | Sp. 11 (F) | Sp. 15 (F) |
| <b>Single-speaker model</b>  | 0.00509      | 0.00794   | 0.00746    | 0.00781    |
| <b>Multi-speaker model</b>   | 0.00657      | 0.00909   | 0.00960    | 0.01594    |
| <b>Multi-speaker model personalized to new speaker with <math>k</math> minutes of data</b> |              |           |            |            |
| <b>5 min</b>   | 0.00481      | 0.00682   | 0.00625    | 0.00681    |
| <b>12.5 min</b>  | 0.00457      | 0.00620   | 0.00589    | 0.00655    |
| <b>25 min</b>  | 0.00443      | 0.00595   | 0.00570    | 0.00621    |
| <b>50 min</b>  | 0.00425      | 0.00561   | 0.00553    | 0.00596    |



# Voice Controllability

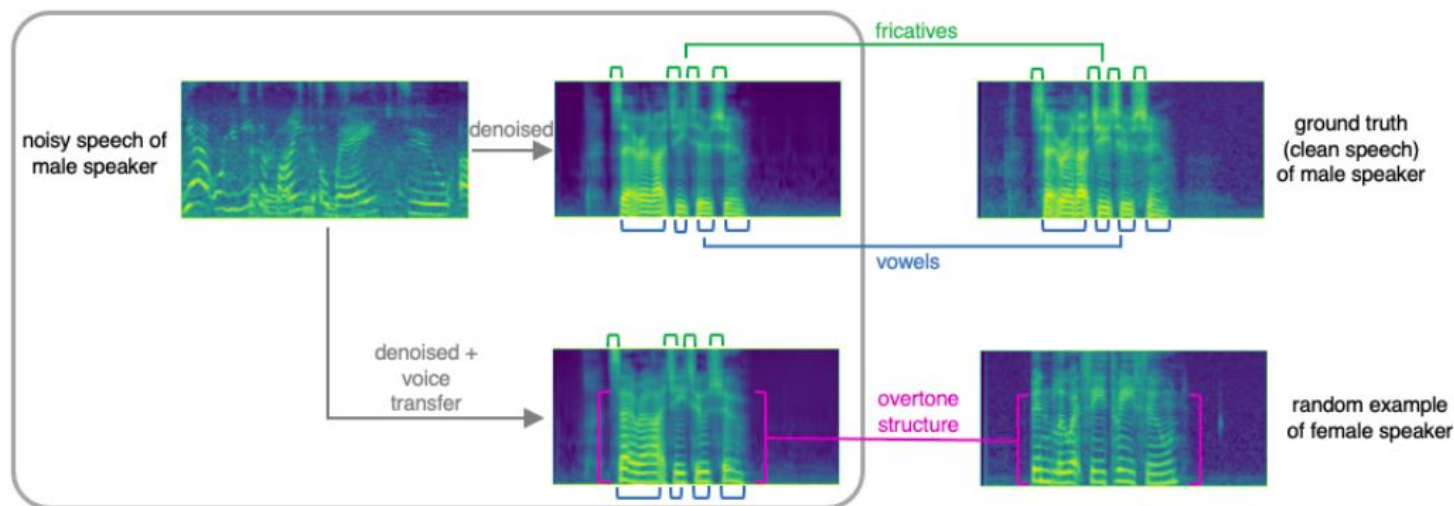


Figure 5. **Voice Transfer Examples.** By swapping the speaker code at the decoder stage, we can synthesize clean audio in a different target speaker's voice. Images shown are mel-spectrogram representations of audio. Note how the linguistic content (*i.e.*, vowels and fricatives) are carried over from the original male speaker, while the pitch and overtone structure are changed to that of the female speaker.

# Visual Acoustic Matching

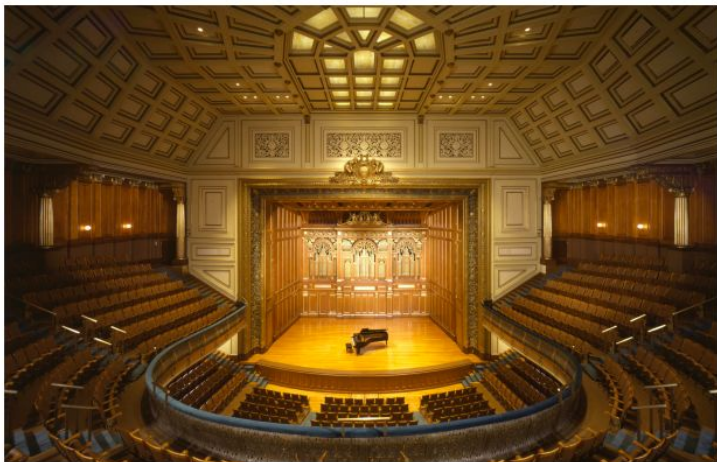
Changan Chen<sup>1,4</sup>   Ruohan Gao<sup>2</sup>   Paul Calamia<sup>3</sup>   Kristen Grauman<sup>1,4</sup>

<sup>1</sup>University of Texas at Austin   <sup>2</sup>Stanford University   <sup>3</sup>Facebook Reality Labs   <sup>4</sup>Facebook AI Research

CVPR 2022, <https://arxiv.org/pdf/2202.06875.pdf>

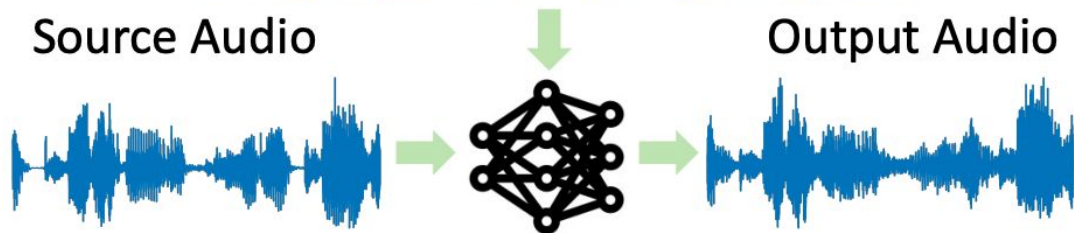
# Introduction

Target Space



$$A_r(t) = A_s(t) * R(t)$$

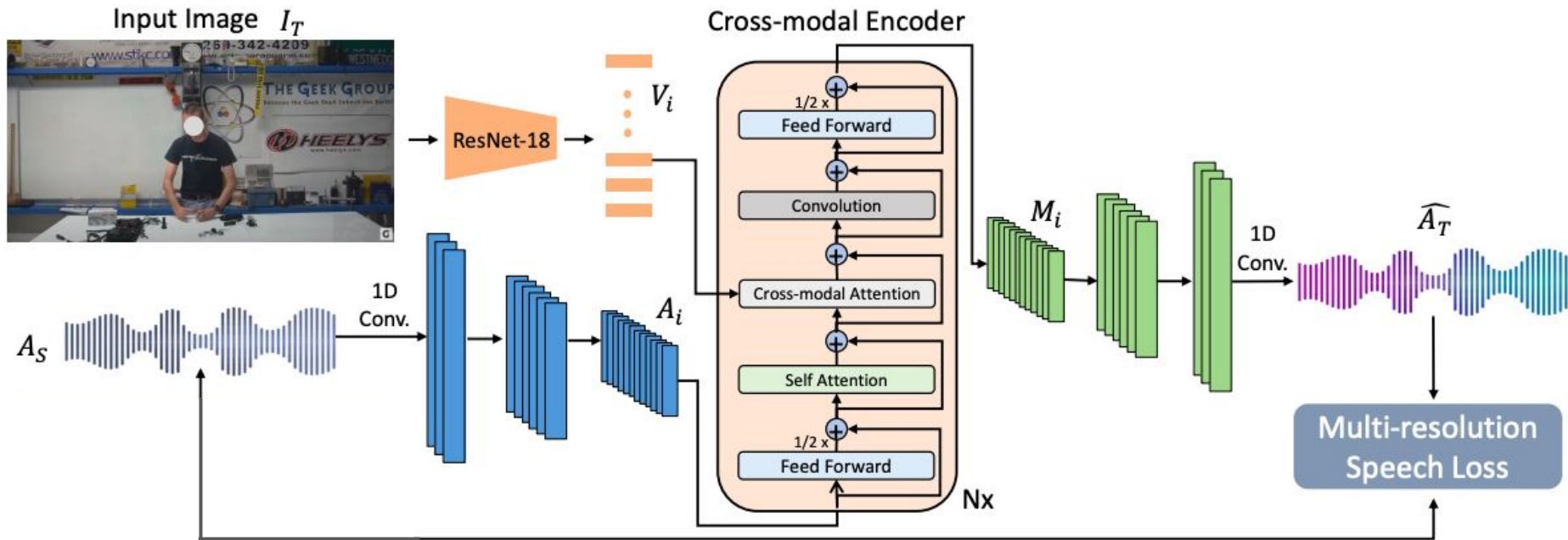
Target Audio      Source Audio      RIR



# Challenges

- Unpaired data
- Modelling different regions of the room
- Capture geometry and materials present in the scene

# Approach



# Datasets

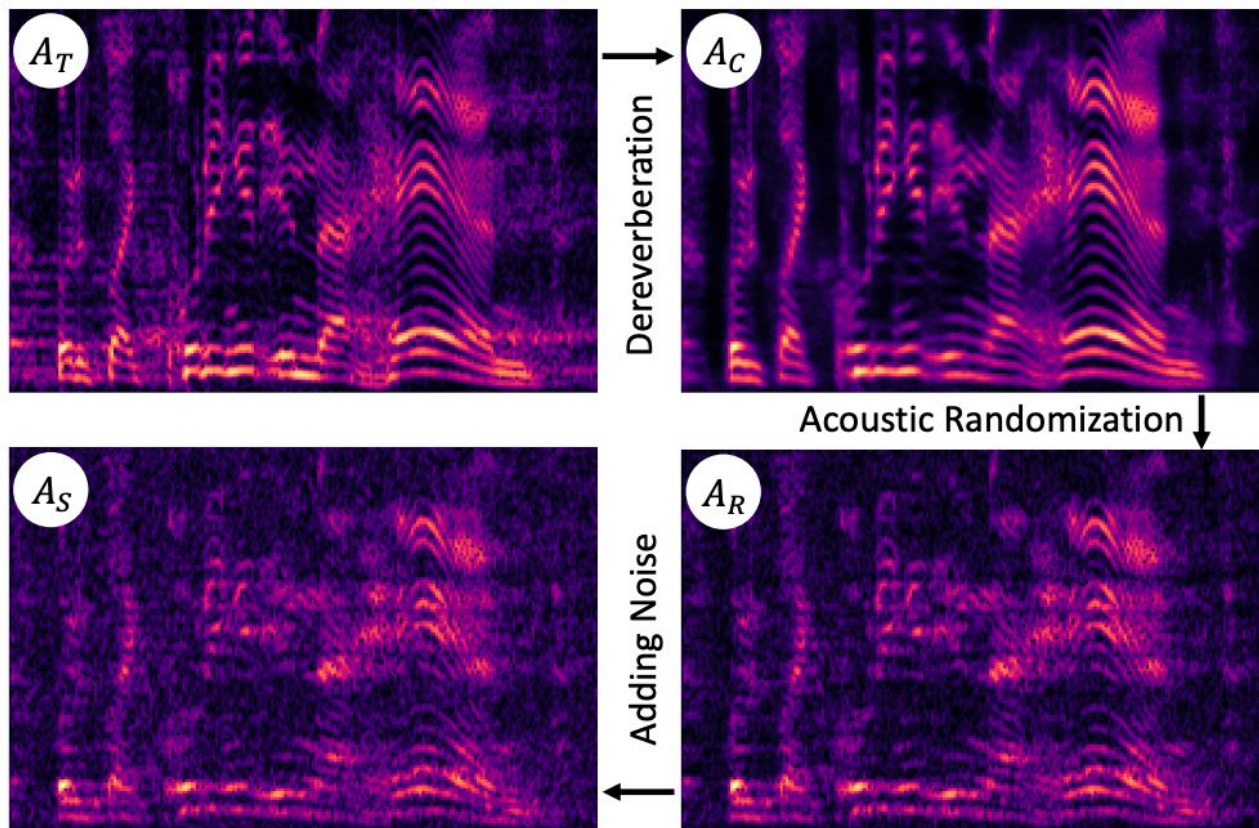
## 1. SoundSpaces-Speech Dataset

- Synthetic Data on Matterport3D
- RIR convolved with audio from LibriSpeech
- 3D humanoid inserted at speaker location

## 2. AVSpeech

- Subset of AVSpeech dataset
- 3-10 seconds videos with single visible human speaker

# Acoustic Alteration Process



# Loss Function

$$\mathcal{L}_G = \sum_{k=1}^K (\mathcal{L}_{Adv}(G; D_k) + \lambda_1 \mathcal{L}_{FM}(G; D_k)) + \lambda_2 \mathcal{L}_{Mel}(G).$$

$$\mathcal{L}_D = \sum_{k=1}^K \mathcal{L}_{Adv}(D_k; G)$$



# Evaluation Metric

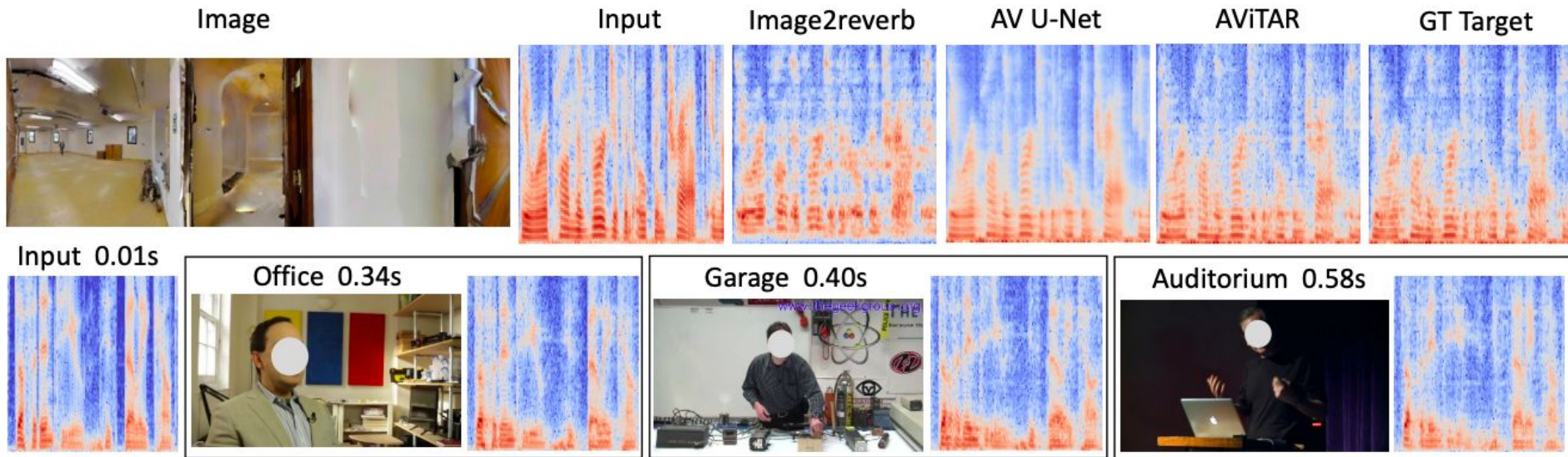
- STFT Distance
- RT60 Error (RTE)
- Mean Opinion Score Error (MOSE)
- User Study

# Results

|                         | <i>SoundSpaces-Speech</i> |              |              |               |              |              | <i>Acoustic AVSpeech</i> |              |               |              |
|-------------------------|---------------------------|--------------|--------------|---------------|--------------|--------------|--------------------------|--------------|---------------|--------------|
|                         | <i>Seen</i>               |              |              | <i>Unseen</i> |              |              | <i>Seen</i>              |              | <i>Unseen</i> |              |
|                         | STFT                      | RTE (s)      | MOSE         | STFT          | RTE (s)      | MOSE         | RTE (s)                  | MOSE         | RTE (s)       | MOSE         |
| Input audio             | 1.192                     | 0.331        | 0.617        | 1.206         | 0.356        | 0.611        | 0.387                    | 0.658        | 0.392         | 0.634        |
| Blind Reverberator [64] | 1.338                     | 0.044        | 0.312        | -             | -            | -            | -                        | -            | -             | -            |
| Image2Reverb [55]       | 2.538                     | 0.293        | 0.508        | 2.318         | 0.317        | 0.518        | -                        | -            | -             | -            |
| AV U-Net [22]           | <b>0.638</b>              | 0.095        | 0.353        | <b>0.658</b>  | 0.118        | 0.367        | 0.156                    | 0.570        | 0.188         | 0.540        |
| AViTAR w/o visual       | 0.862                     | 0.140        | 0.217        | 0.902         | 0.186        | 0.236        | 0.194                    | 0.504        | 0.207         | 0.478        |
| AViTAR                  | 0.665                     | <b>0.034</b> | <b>0.161</b> | 0.822         | <b>0.062</b> | <b>0.195</b> | <b>0.144</b>             | <b>0.481</b> | <b>0.183</b>  | <b>0.453</b> |

|                   | SoundSpaces          | AVSpeech             |
|-------------------|----------------------|----------------------|
| Input Speech      | 42.1% / <b>57.9%</b> | 40.1% / <b>59.9%</b> |
| Image2Reverb [55] | 25.9% / <b>74.1%</b> | - / -                |
| AV U-Net [22]     | 29.8% / <b>70.2%</b> | 27.2% / <b>72.8%</b> |
| AViTAR w/o visual | 39.6% / <b>60.4%</b> | 46.3% / <b>53.9%</b> |

# Qualitative Results



<https://vision.cs.utexas.edu/projects/visual-acoustic-matching/>

# Self-supervised object detection from audio-visual correspondence

Triantafyllos Afouras<sup>1\*†</sup> Yuki M. Asano<sup>1\*</sup> Francois Fagan<sup>2</sup> Andrea Vedaldi<sup>2</sup> Florian Metze<sup>2</sup>

<sup>1</sup> Visual Geometry Group, University of Oxford

<sup>2</sup> Facebook AI

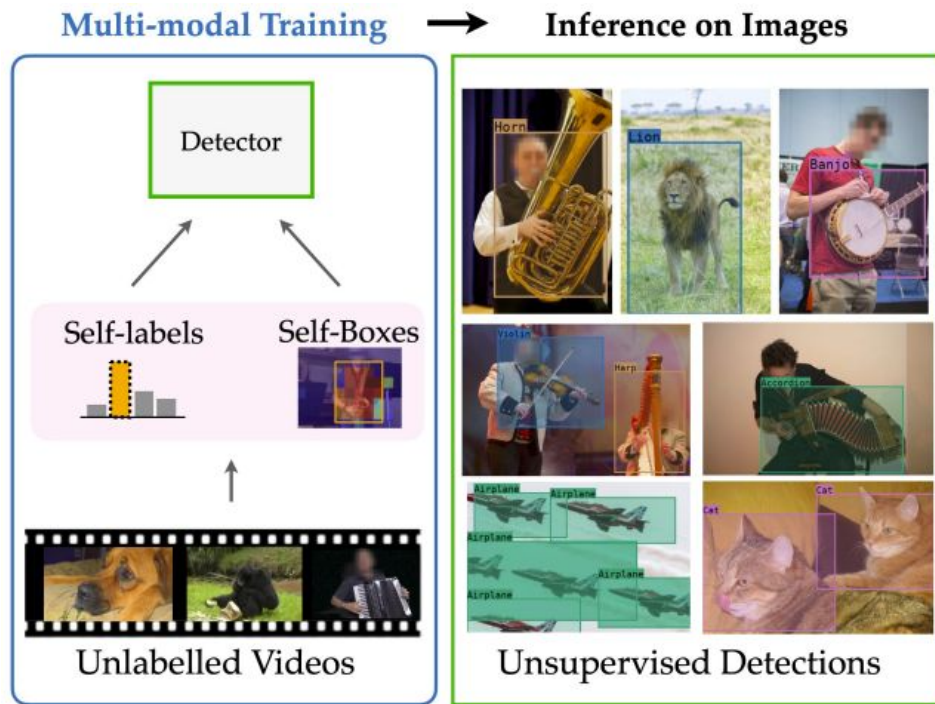
afourast@robots.ox.ac.uk

# Object Detection

- Supervised
- Weakly-Supervised
- Unsupervised

# Introduction

- Object Detection using Audio

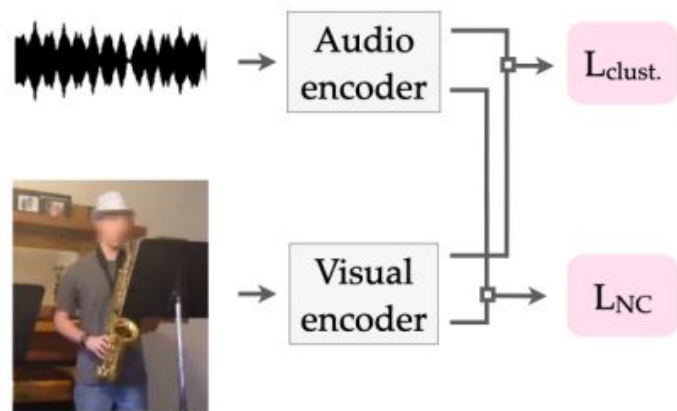


# Problems with audio-visual detection

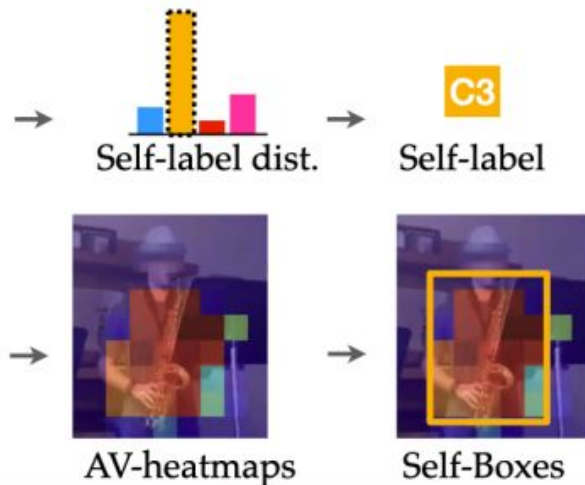
- Not applicable to silent videos
- Only heatmap, No bounding box

# Approach

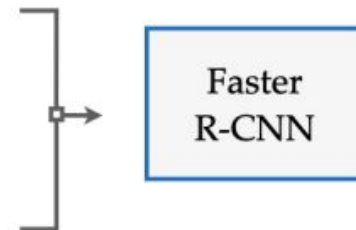
## 1. Representation Learning



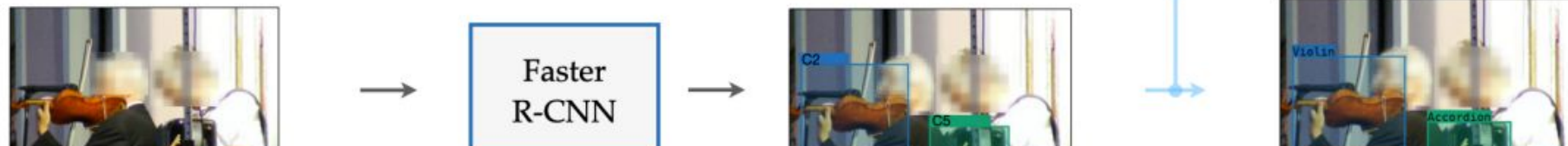
## 2. Self-detection extraction



## 3. Detector Training



## 4. Image-only Inference





# Results

| Method                | single-instr. |             | multi-instr. |
|-----------------------|---------------|-------------|--------------|
|                       | IoU-0.5       | AUC         | cIoU-0.3     |
| Sound of pixels [109] | 38.2          | 40.6        | 39.8         |
| Object t. Sound [7]   | 32.7          | 39.5        | 27.1         |
| Attention [82]        | 36.5          | 39.5        | 29.9         |
| DMC [44]              | 32.8          | 38.2        | 32.0         |
| DSOL [45]             | 38.9          | 40.9        | 48.7         |
| <b>Ours</b>           | <b>50.6</b>   | <b>47.5</b> | <b>52.4</b>  |

| Method                 | No labels? | VGGSound          |                   |                          | Audioset          |                   |                          | OpenImages        |                   |                          |
|------------------------|------------|-------------------|-------------------|--------------------------|-------------------|-------------------|--------------------------|-------------------|-------------------|--------------------------|
|                        |            | mAP <sub>30</sub> | mAP <sub>50</sub> | mAP <sub>[50:95:5]</sub> | mAP <sub>30</sub> | mAP <sub>50</sub> | mAP <sub>[50:95:5]</sub> | mAP <sub>30</sub> | mAP <sub>50</sub> | mAP <sub>[50:95:5]</sub> |
| PCL (WSOD) [90]        | ✗          | 54.9              | 27.7              | 7.6                      | 39.0              | 17.5              | 4.4                      | 37.9              | 14.5              | 3.5                      |
| Ours - weak sup.       | ✗          | 67.6              | 42.9              | 14.2                     | 50.6              | 30.9              | 10.3                     | 48.9              | 33.7              | 9.5                      |
| Center Box*            | ✓          | 29.6              | 5.6               | 1.5                      | 15.1              | 3.5               | 0.7                      | 20.7              | 4.2               | 0.8                      |
| Selective Search* [94] | ✓          | 5.2               | 1.1               | 0.4                      | 2.8               | 0.4               | 0.1                      | 7.4               | 2.1               | 0.7                      |
| COCO-trained RPN*      | ✗          | 33.4              | 7.5               | 1.6                      | 19.0              | 4.1               | 0.8                      | 24.4              | 11.1              | 2.6                      |
| Ours - self-boxes*     | ✓          | 48.1              | 29.6              | 10.0                     | 27.8              | 14.1              | 4.8                      | NA                | NA                | NA                       |
| <b>Ours - full</b>     | ✓          | <b>52.3</b>       | <b>39.4</b>       | <b>14.7</b>              | <b>44.3</b>       | <b>28.0</b>       | <b>9.6</b>               | <b>39.9</b>       | <b>28.5</b>       | <b>7.6</b>               |