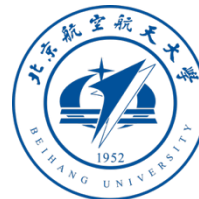


# Revisiting Audio-Visual Segmentation with Vision-Centric Transformer

Shaofei Huang<sup>1,2</sup> Rui Ling<sup>3</sup> Tianrui Hui<sup>1</sup> Hongyu Li<sup>4</sup> Xu Zhou<sup>5</sup> Shifeng Zhang<sup>5</sup> Si Liu<sup>3</sup> Richang Hong<sup>1</sup> Meng Wang<sup>1</sup>

<sup>1</sup>HFUT <sup>2</sup>IIE, CAS <sup>3</sup>SCSE, BUAA <sup>4</sup>SAI, BUAA <sup>5</sup>Sangfor



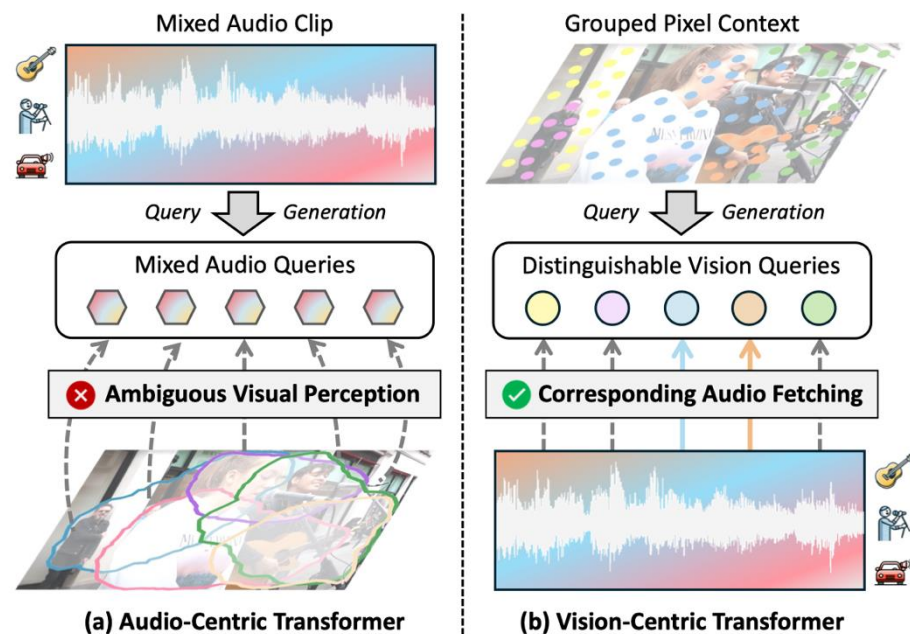
# Task

- Predict the pixel-level mask of the objects which make sounds in the given video clip according to the associated audio information



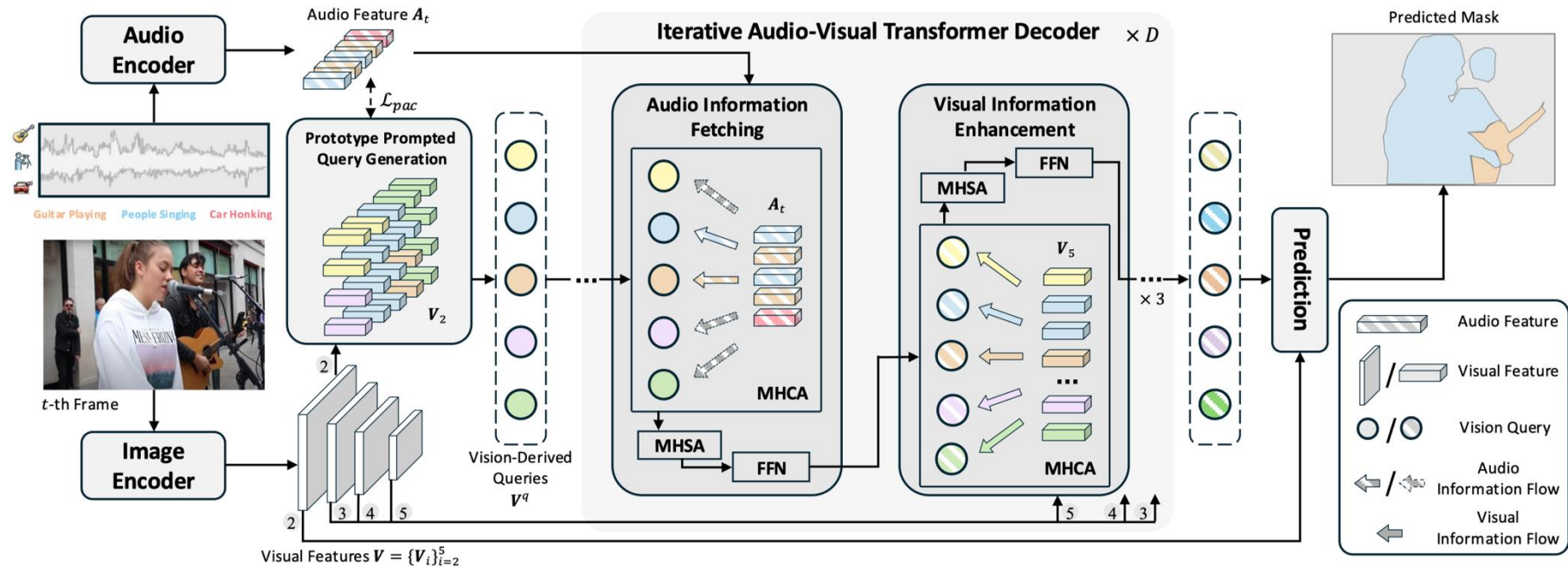
# Motivation

- **Audio-Centric Transformer (ACT)**
  - Audio-derived queries from mixed sound sources causes perception confusion for audio-derived queries
  - Delayed integration of visual information may lead to the loss of visual details
- **Visual-Centric Transformer (VCT)**
  - Vision-derived queries focus on semantically distinct visual regions, showing superior discriminative ability
  - Interacting with audio features excludes unrelated sounds and alleviates the perceptual ambiguity



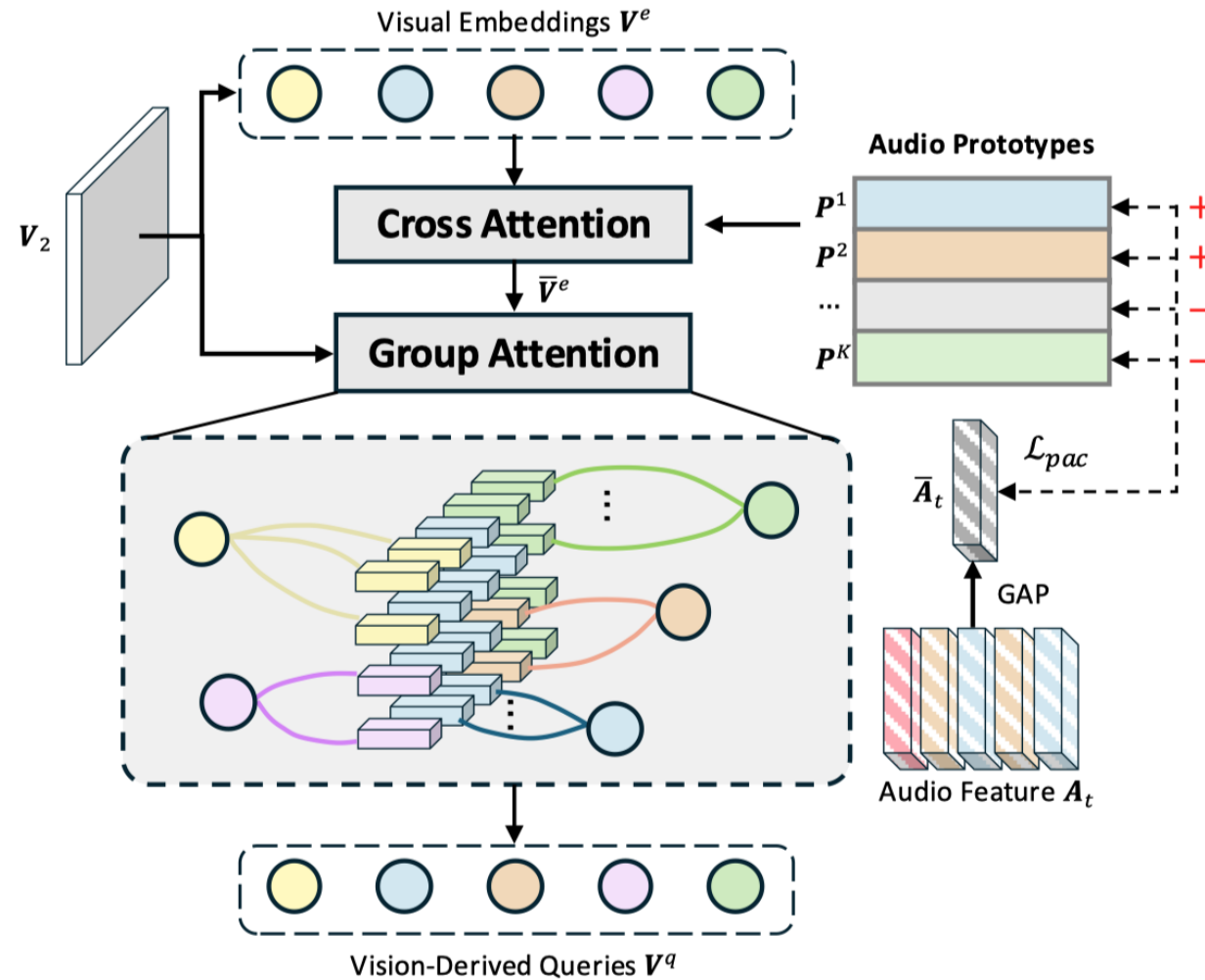
# Framework

- Our VCT framework leverages **vision-derived queries** to fetch audio information and visual details in the **iterative audio-visual Transformer decoder**
- **Audio information fetching**: obtain the corresponding sound information for each query region
- **Visual information enhancement**: captures multi-scale fine-grained visual features for accurate mask predictions



# Prototype Prompted Query Generation (PPQG)

- PPQG generates vision-derived queries that contain both rich visual details and audio semantics



# Prototype Prompted Query Generation (PPQG)

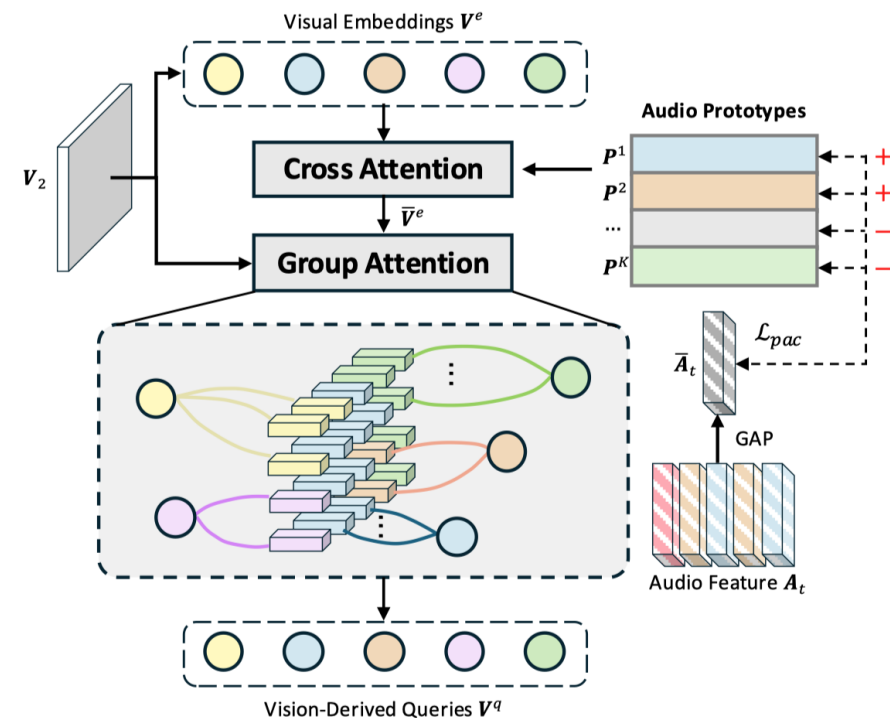
- I. **Visual embedding aggregation:** perform spatial information aggregation obtain a set of visual embeddings as initial vision-derived queries

$$\mathbf{V}^h = \text{Conv}_{1 \times 1}(\delta(\text{Conv}_{3 \times 3}(\delta(\text{Conv}_{1 \times 1}(\mathbf{V}_2))))),$$

$$\mathbf{V}^e = \text{Reshape}(\text{MLP}(\text{Reshape}(\mathbf{V}^h))),$$

- II. **Audio prototype prompting:** define audio prototypes to prompt queries with audio event categories present in the scene.

$$\bar{\mathbf{V}}^e = \mathbf{V}^e + \text{Softmax}\left(\frac{(\mathbf{V}^e \mathbf{W}_1^q)(\mathbf{P} \mathbf{W}_1^k)^T}{\sqrt{C^h}}\right)(\mathbf{P} \mathbf{W}_1^v).$$



# Prototype Prompted Query Generation (PPQG)

We design a **prototype-audio contrastive loss** ( $\mathcal{L}_{pac}$ ) for audio prototype learning

$$\mathcal{L}_{bce}(\mathbf{M}_k, \mathbf{M}_k^*) = -\mathbf{M}_k^* \log(\mathbf{M}_k) - (1 - \mathbf{M}_k^*) \log(1 - \mathbf{M}_k),$$

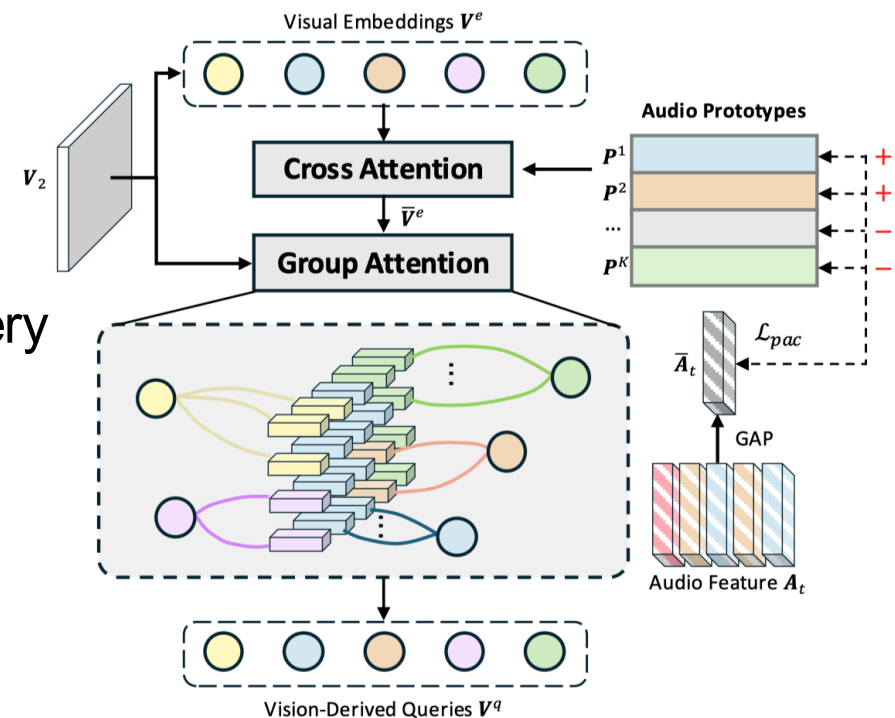
$$\mathcal{L}_{pac}(\mathbf{M}, \mathbf{M}^*) = \frac{1}{K} \sum_{k=1}^K \mathcal{L}_{bce}(\mathbf{M}_k, \mathbf{M}_k^*).$$

**III. Pixel context grouping:** group pixel context for each query through hard assignment to increase distinguishability

$$\mathbf{R} = \text{Softmax}((\bar{\mathbf{V}}^e \mathbf{W}_2^q)(\mathbf{V}^h \mathbf{W}_2^k)^T + \mathbf{G}),$$

$$\hat{\mathbf{R}} = \text{One-hot}(\arg \max_{\mathcal{N}}(\mathbf{R})) + \mathbf{R} - \text{sg}(\mathbf{R}),$$

$$\mathbf{V}^q = \bar{\mathbf{V}}^e + (\text{Norm}(\hat{\mathbf{R}})(\mathbf{V}^h \mathbf{W}_2^v)) \mathbf{W}^o,$$



# Quantitative Results

Method	Reference	Backbone	Image Size	AVS-Semantic (AVSS)		Single-Source (S4)		Multi-Source (MS3)	
				$\mathcal{M}_{\mathcal{J}}$	$\mathcal{M}_{\mathcal{F}}$	$\mathcal{M}_{\mathcal{J}}$	$\mathcal{M}_{\mathcal{F}}$	$\mathcal{M}_{\mathcal{J}}$	$\mathcal{M}_{\mathcal{F}}$
TPAVI [42, 43]	[ECCV'22]	ResNet-50 PVT-v2	224×224	-	-	72.8	84.8	46.9	57.8
				29.8	35.2	78.7	87.9	54.0	64.5
AQFormer [15]	[IJCAI'23]	ResNet-50 PVT-v2	224×224	-	-	77.0	86.4	55.7	66.9
				-	-	81.6	89.4	61.1	<u>72.1</u>
ECMVAE [28]	[ICCV'23]	ResNet-50 PVT-v2	224×224	-	-	76.3	86.5	48.7	60.7
				-	-	81.7	90.1	57.8	70.8
CATR [20]	[ACMMM'23]	ResNet-50 PVT-v2	224×224	-	-	74.8	86.6	52.8	65.3
				32.8	38.5	81.4	89.6	59.0	70.0
AVSC [22]	[ACMMM'23]	ResNet-50 PVT-v2	224×224	-	-	77.0	85.2	49.6	61.5
				-	-	80.6	88.2	58.2	65.1
BAVS [23]	[TMM'24]	ResNet-50	224×224	24.7	29.6	78.0	85.3	50.2	62.4
		PVT-v2		32.6	36.4	82.0	88.6	58.6	65.5
		Swin-B		<u>33.6</u>	<u>37.5</u>	<u>82.7</u>	<u>89.8</u>	<u>59.6</u>	<u>65.9</u>
COMBO [40]	[CVPR'24]	ResNet-50 PVT-v2	224×224	33.3	37.3	<u>81.7</u>	90.1	54.5	66.6
				<u>42.1</u>	46.1	<u>84.7</u>	<u>91.9</u>	59.2	71.2
CPM [6]	[ECCV'24]	ResNet-50	224×224	<u>34.5</u>	<u>39.6</u>	81.4	<u>90.5</u>	<u>59.8</u>	<u>71.0</u>
TeSO [37]	[ECCV'24]	Swin-B	384×384	39.0	45.1	<u>83.3</u>	<u>93.3</u>	66.0	80.1
SelM [19]	[ACMMM'24]	ResNet-50 PVT-v2	224×224	31.9	37.2	76.6	86.2	54.5	65.6
				41.3	<u>46.9</u>	83.5	91.2	60.3	71.3
AVSBias [33]	[ACMMM'24]	Swin-B	384×384	<u>44.4</u>	<u>49.9</u>	<u>83.3</u>	93.0	<u>67.2</u>	<u>80.8</u>
<b>VCT (Ours)</b>	-	ResNet-50	224×224	<b>37.5</b>	<b>42.2</b>	<b>81.8</b>	<b>90.6</b>	<b>61.9</b>	<b>74.7</b>
		PVT-v2	224×224	<b>44.7</b>	<b>49.5</b>	<b>84.8</b>	<b>92.1</b>	<b>62.0</b>	<b>75.0</b>
		Swin-B	224×224	<b>47.9</b>	<b>52.9</b>	<b>84.7</b>	<b>92.3</b>	<b>67.5</b>	<b>79.3</b>
		Swin-B	384×384	<b>51.2</b>	<b>55.5</b>	<b>86.2</b>	<b>93.4</b>	<b>67.6</b>	<b>81.4</b>

# Quantitative Results

Framework	Queries	$\mathcal{M}_{\mathcal{J}}$	$\mathcal{M}_{\mathcal{F}}$
ACT	Audio-Derived Queries	33.2	37.0
VCT	Naive Vision Queries	35.2	39.3
	w/ Cross-Attention	35.8	39.8
	w/ Group-Attention	36.3	40.5
	w/ Audio Prototypes	<b>37.5</b>	<b>42.2</b>

Table 2. Ablation study of generating vision-derived queries in our PPQG module on the AVSS subset.

Method	$\mathcal{M}_{\mathcal{J}}$	$\mathcal{M}_{\mathcal{F}}$
w/o Audio Prototypes	36.3	40.5
w/ Audio Prototypes, w/o loss	36.3	40.4
w/ Audio Prototypes, w/ visual loss	36.5	40.8
w/ Audio Prototypes, w/ PAC loss	<b>37.5</b>	<b>42.2</b>

Table 3. Ablation study of audio prototype prompting in PPQG.

	Audio Feature	$\mathcal{M}_{\mathcal{J}}$	$\mathcal{M}_{\mathcal{F}}$
Visual Feature	Multiply	33.9	37.9
	Concatenation	35.3	39.3
	Addition	36.3	40.5
Audio Proto	Replace	36.6	40.8
VCT	Full Model	<b>37.5</b>	<b>42.2</b>

Table 4. Ablation study of audio feature incorporation.

# Qualitative Results

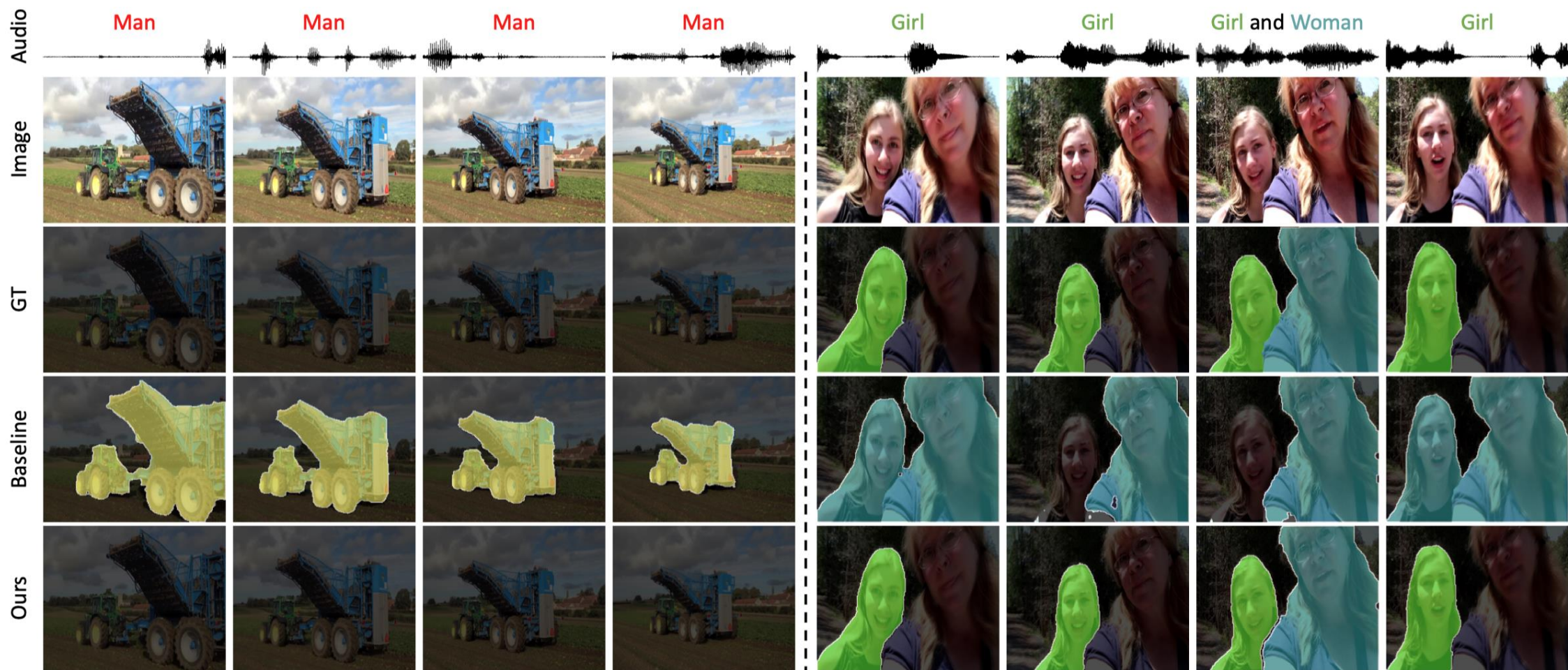


Figure 4. Qualitative comparison between our full model and the ACT baseline. Existing sounds correspond to masks of the same colors.

# Qualitative Results

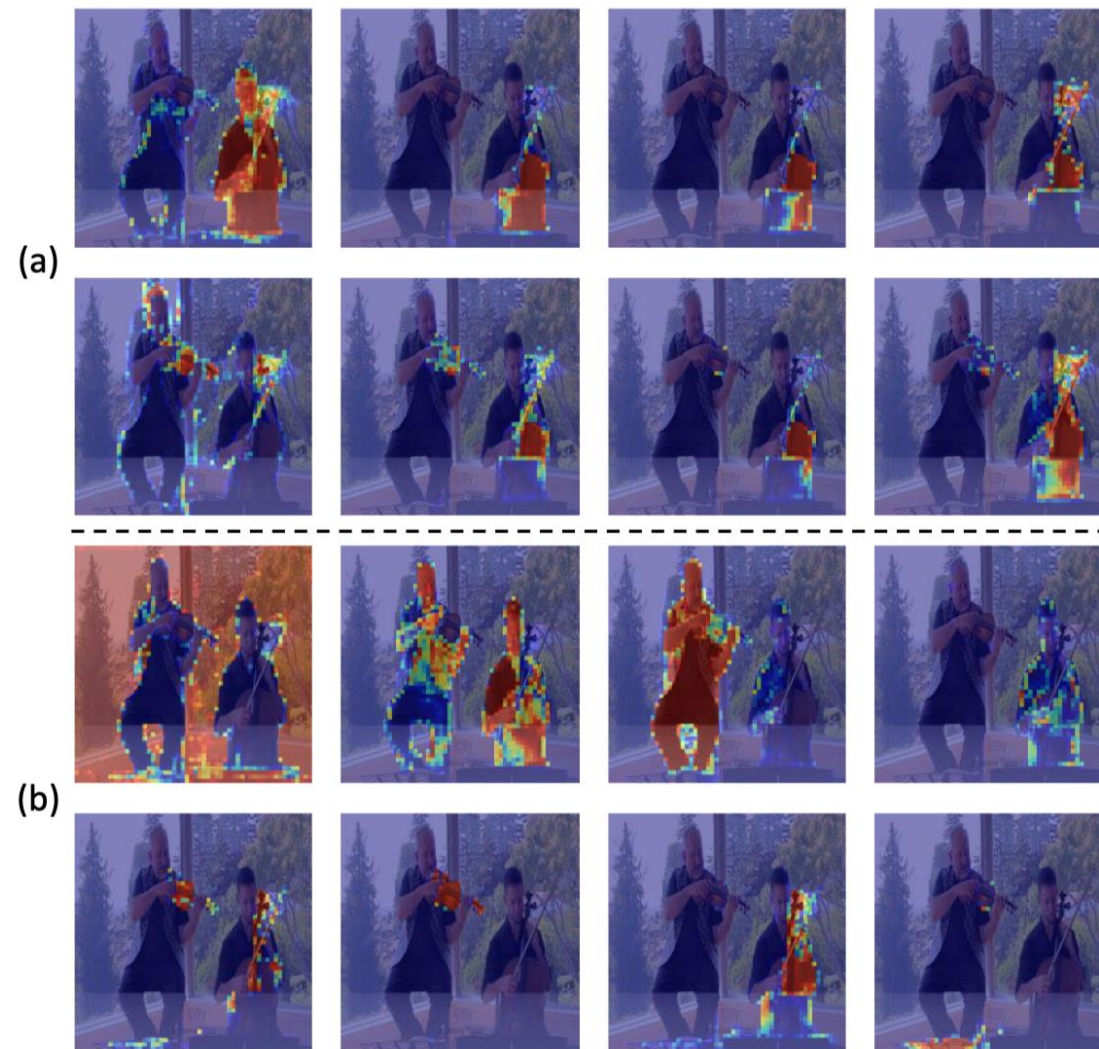


Figure 5. Visualization of logit maps from different types of queries. (a) Audio-derived queries. (b) Vision-derived queries.

# Thank You!

