



# Visual Sound Localization in the Wild by Cross-Modal Interference Erasing

Xian Liu<sup>1,2\*</sup>, Rui Qian<sup>1,3\*</sup>, Hang Zhou<sup>1\*</sup>, Di Hu<sup>4</sup>, Weiyao Lin<sup>3</sup>, Ziwei Liu<sup>5</sup>, Bolei Zhou<sup>1</sup>, Xiaowei Zhou<sup>2</sup>

<sup>1</sup>The Chinese University of Hong Kong <sup>2</sup>Zhejiang University <sup>3</sup>Shanghai Jiao Tong University <sup>4</sup>Renmin University of China <sup>5</sup>Nanyang Technological University



## Motivation



- In the localization task, sound sources of **different volumes** should be evenly identified.
- The **off-screen sound** and **background noise** will compromise the procedure of audio-visual modality matching.

## Key Notations

Audio-visual Pairs:  $\{(a_i, v_i) | i = 1, 2, \dots, N\}$

Audio-visual Prototype:  $\mathcal{P} \in R^{K \times C} \quad y_i \in \{0, 1\}^K$

Cross-Distillation Loss:  $D_{KL}(p^{va} || p^{av}) + D_{KL}(p^{av} || p^{va})$

## Experiments

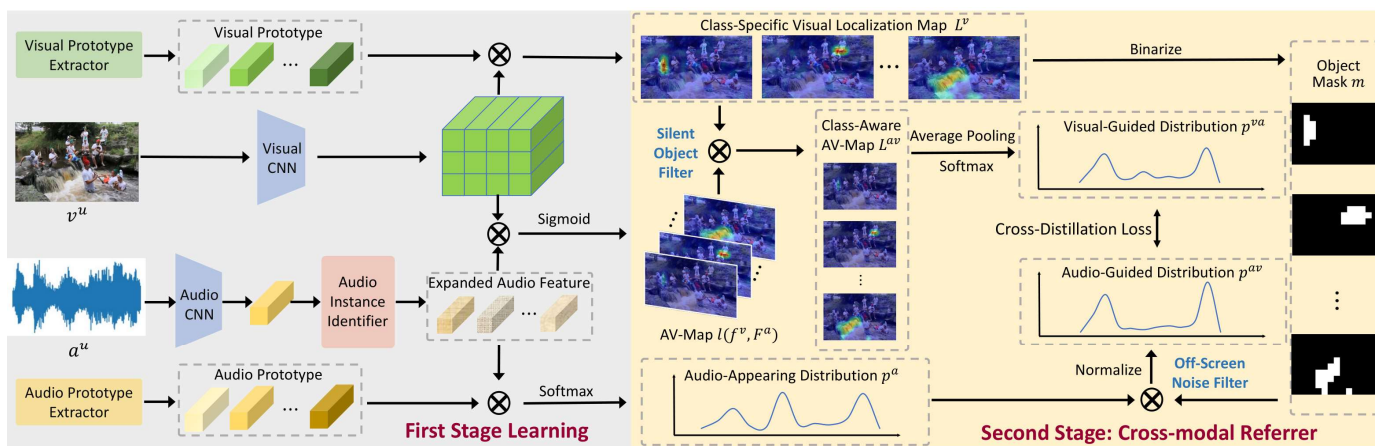
Scenario	Single Sound Scene (a)				General in-the-Wild Scene (b)							
	MUSIC		VGGSound		MUSIC-Syn.		MUSIC-Duet		MUSIC-Un.		VGG-Un.	
Dataset	IoU	AUC	IoU	AUC	CIoU	AUC	CIoU	AUC	CIoU	AUC	CIoU	AUC
Object-of-sound	26.1	35.8	48.4	46.1	3.7	10.2	13.2	18.3	0.1	6.8	7.8	15.1
Sound-of-pixel	40.5	43.3	42.5	45.1	8.1	11.8	16.8	16.8	7.5	11.6	7.9	14.4
DSOL	51.4	43.6	49.3	45.8	32.3	23.5	30.2	22.1	3.2	7.3	8.1	12.2
Interference Eraser	<b>53.9</b>	<b>50.7</b>	<b>51.3</b>	<b>46.9</b>	<b>47.6</b>	<b>29.8</b>	<b>52.9</b>	<b>33.8</b>	<b>15.6</b>	<b>15.3</b>	<b>12.8</b>	<b>17.6</b>

Visual Sound Localization results in (a) Single Sound Scene; (b) General in-the-Wild Scene.

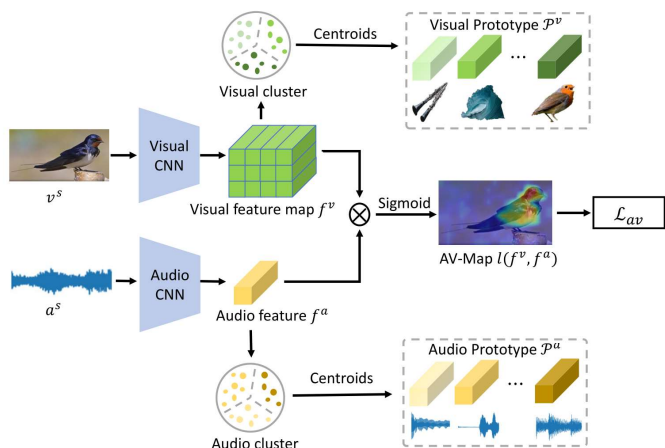
Dataset	Method	MUSIC				VGGSound			
		NMI	Prec.	Rec.	mAP	NMI	Prec.	Rec.	mAP
w/o Audio-Instance-Identifier		0.692	0.372	0.335	0.355	0.410	0.218	0.076	0.189
Audio-Instance-Identifier w/o curriculum		0.758	0.441	0.489	0.403	0.414	0.298	0.160	0.231
Audio-Instance-Identifier w/ curriculum		<b>0.809</b>	<b>0.461</b>	<b>0.715</b>	<b>0.433</b>	<b>0.436</b>	<b>0.346</b>	<b>0.232</b>	<b>0.283</b>

Ablation study on Audio-Instance-Identifier with curriculum learning strategy.

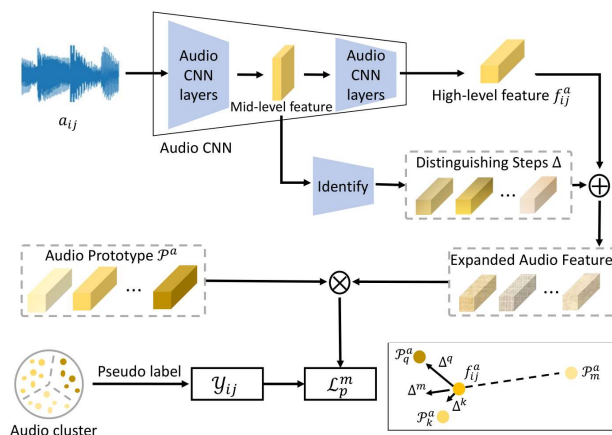
## Interference Eraser (IEr) Framework



### Left: Audio-visual Prototype Extraction



### Right: Audio-Instance-Identifier



Localization results on realistic and synthetic videos. The green box indicates target sounding object area, while the red box means this class of object is silent and its activation value should be low.

## Conclusions

- We introduce a novel framework Interference Eraser (IEr) to enhance robust visual sound localization for **in-the-wild scenes**.
- Our proposed Audio-Instance-Identifier learns the distinguishing-step to achieve volume agnostic **mixed sound perception**.
- We propose the Cross-modal Referrer to eliminate the interference of **visible but silent objects** and **audible but off-screen sounds**.