

## 1. Task & Motivation

### Audio-Visual Segmentation

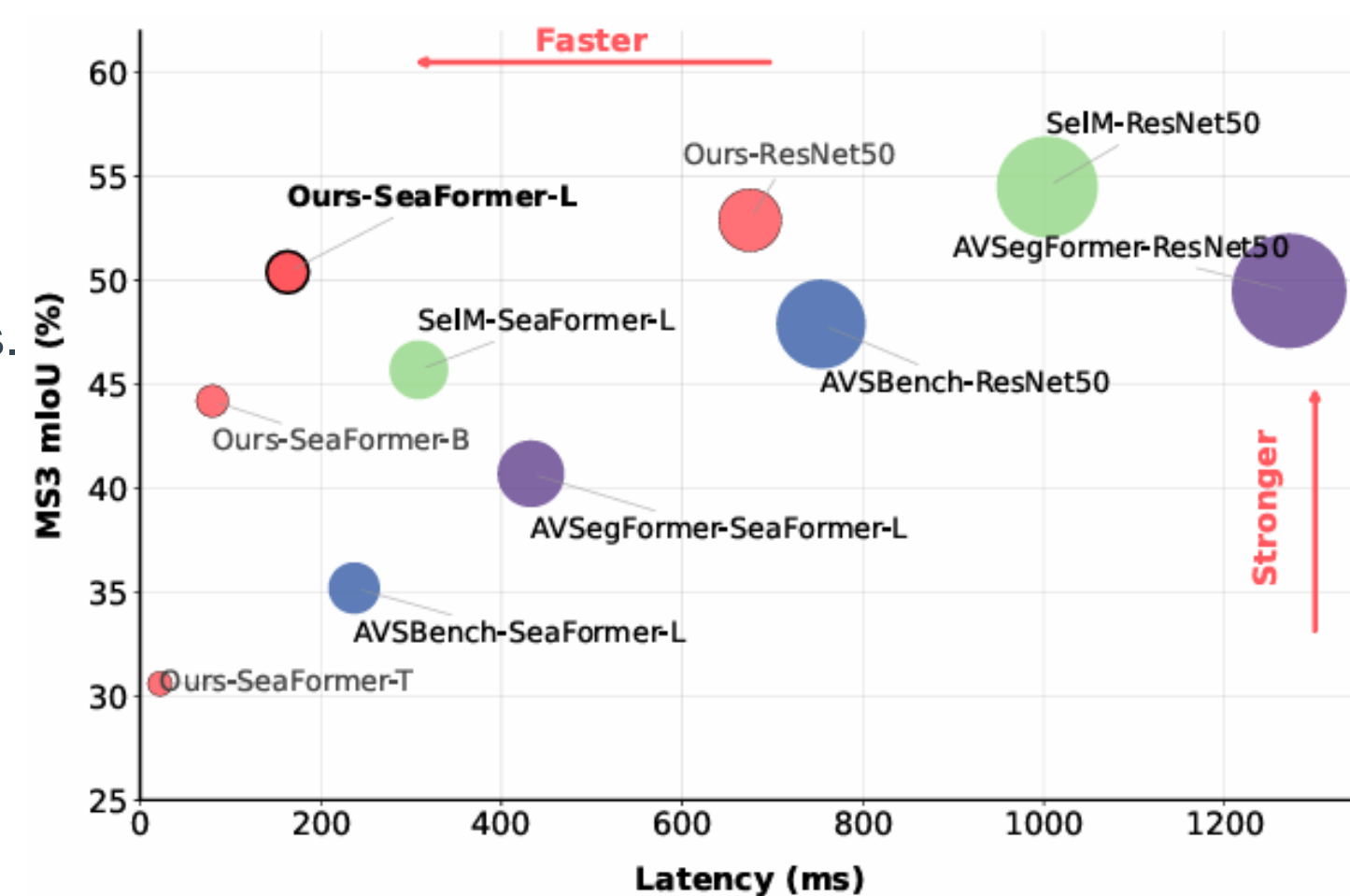
Pixel-level localization of sound-emitting objects in videos.

### Lightweight AVS

Enable accurate AVS on resource-constrained devices, e.g., mobile phones, VR headsets, and on-device video editing.

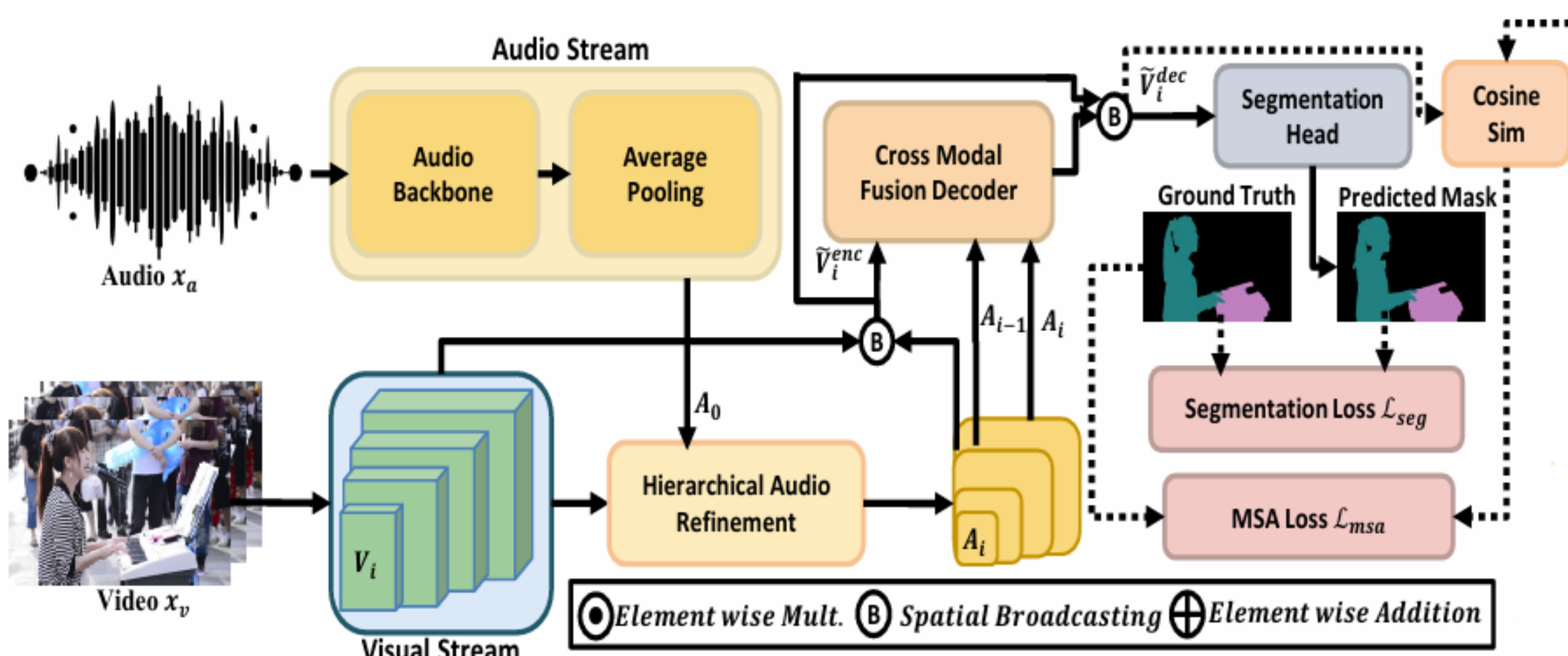
### Challenges

- Existing AVS methods rely on dense cross-modal attention.
- The interaction cost grows quadratically with spatial tokens.
- Backbone replacement alone does not solve the efficiency bottleneck.
- Lightweight models may learn spurious audio-visual correlations.



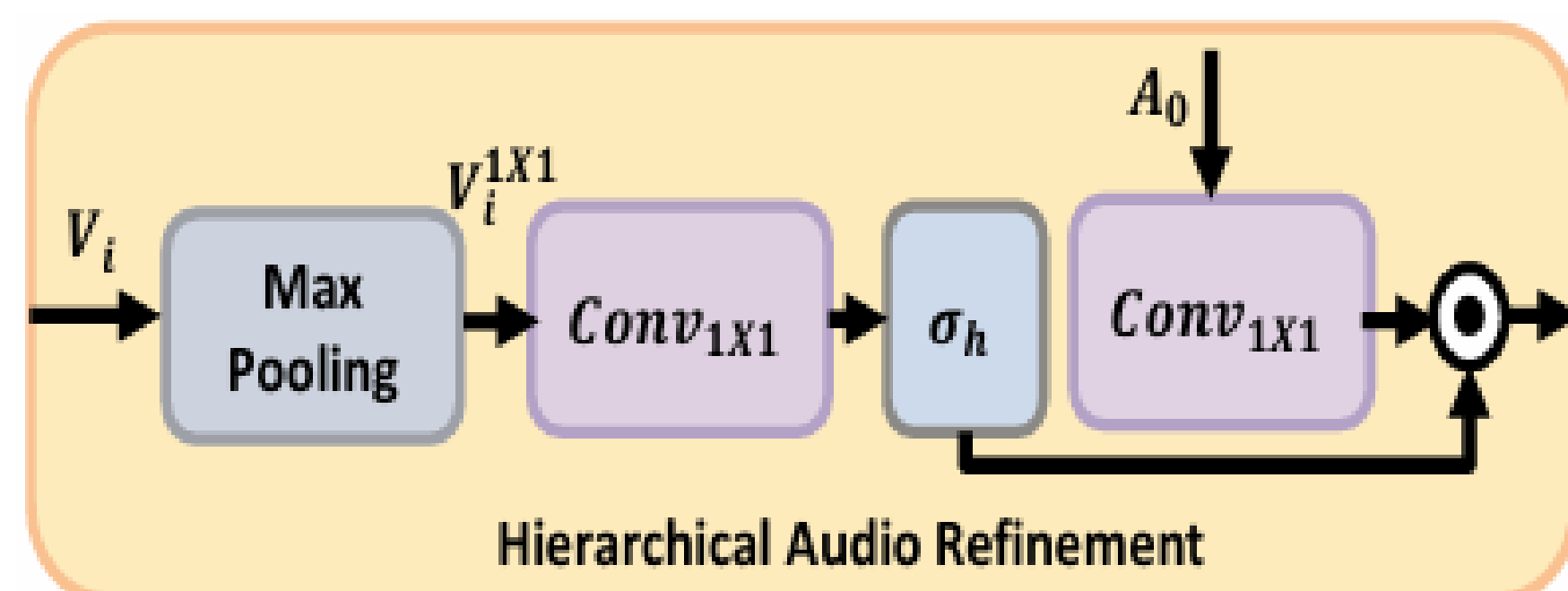
## 2. Approach

A lightweight AVS framework that decouples audio-visual interaction into semantic filtering and spatial grounding.



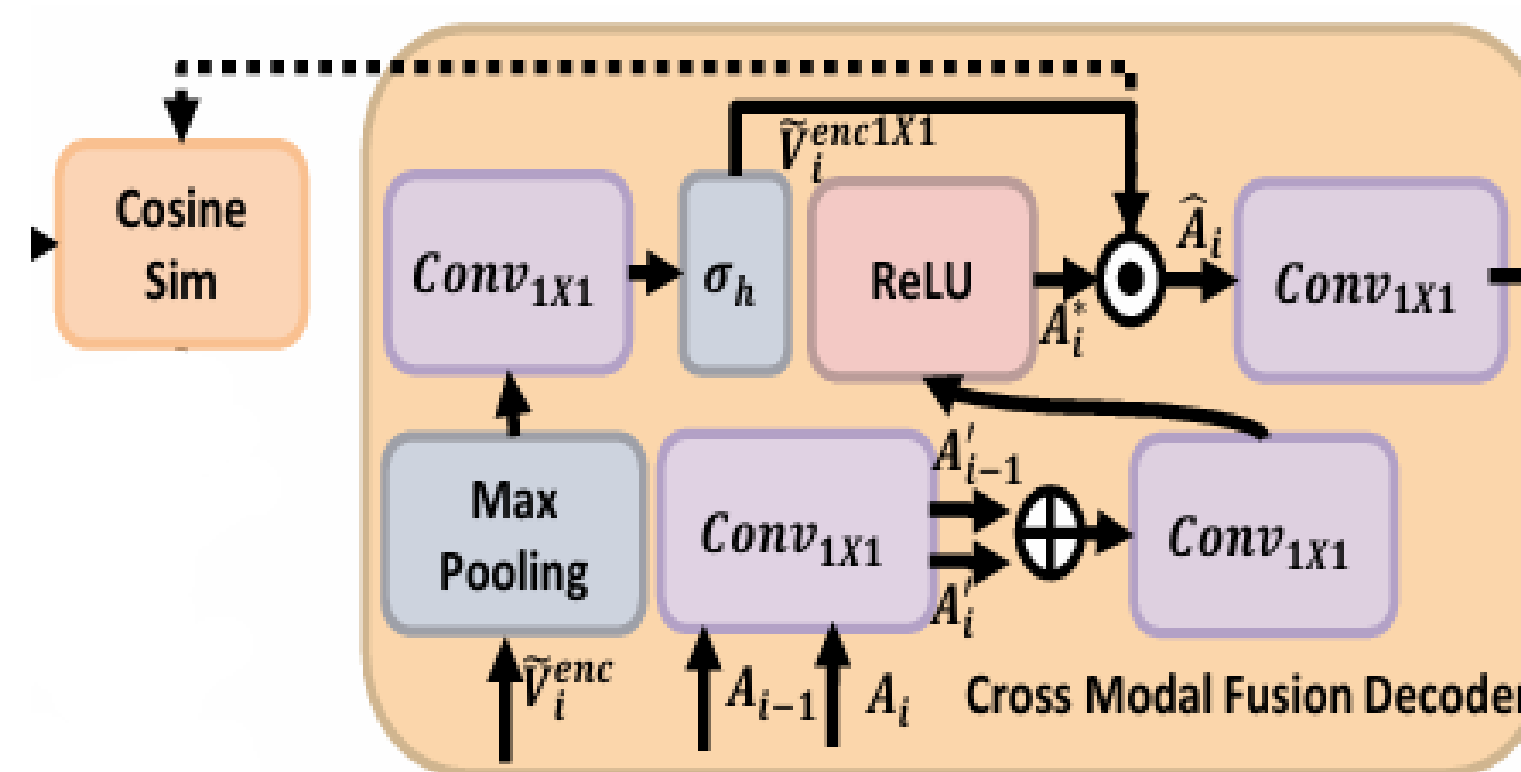
### Reciprocal Audio-Visual Encoder

- Refine the global audio state using visual context.
- Modulate visual features with refined audio cues.
- Replace dense attention with channel-wise interaction.
- Reduce interaction complexity from  $O(N^2)$  to  $O(N)$ .



### Cross-Modal Fusion Decoder

- Progressively recover spatial resolution.
- Maintain audio-guided semantic consistency.
- Inject refined audio cues into visual features.
- Generate the final sounding-object mask.



### Multi-Scale Audio-Visual Alignment Loss

- Align audio cues with valid visual regions.
- Guide the decoder to focus on sounding objects.  $\mathcal{L}_{msa} = \frac{1}{S} \sum_{i=1}^S \text{BCE}(\hat{s}_i, M)$
- Training-only supervision with zero inference overhead.

### Total Loss

$\mathcal{L}_{seg}$  supervises the final segmentation mask.  
 $\mathcal{L}_{msa}$  aligns audio cues with sounding visual regions.  
 Training-only loss with zero inference overhead.

$$\mathcal{L}_{seg} = \mathcal{L}_{dice} + \mathcal{L}_{bce}$$

$$\mathcal{L} = \mathcal{L}_{seg} + \lambda \mathcal{L}_{msa}$$

## 3. Results

### Performance

- New lightweight SOTA on S4, MS3, and AVSS.
- Achieves 75.6 MJ on S4 and 50.4 MJ on MS3.
- Surpasses AVSegFormer-R50 on MS3 with only 20.5M parameters.
- Runs at 163.4 ms on Snapdragon 8 Elite.

Method	Backbone		GPU	Mobile	S4		MS3		Params
	Visual	Audio	ms ↓	ms ↓	$M_J \uparrow$	$M_F \uparrow$	$M_J \uparrow$	$M_F \uparrow$	
AVSBench (Zhou et al., 2022)	R50	VGGish	21.2	753.5	72.8	84.8	47.9	57.8	91.4M
AVSegFormer (Gao et al., 2024)	R50	VGGish	29.0	1271.4	76.5	85.9	49.5	62.8	151.1M
SelM (Li et al., 2024)	R50	VGGish	25.3	1003.8	76.6	86.2	54.5	65.6	117.6M
AVSBench (Zhou et al., 2022)	Sea	MNetV2	19.5	237.1	47.9	64.5	35.2	44.1	30.2M
AVSegFormer (Gao et al., 2024)	Sea	MNetV2	22.2	432.6	53.8	71.4	40.7	50.7	51.0M
SelM (Li et al., 2024)	Sea	MNetV2	18.7	308.6	59.1	77.4	45.7	57.6	39.5M
<b>Ours</b>	Sea	MNetV2	<b>15.9</b>	<b>163.4</b>	<b>75.6</b>	<b>86.2</b>	<b>50.4</b>	<b>62.6</b>	<b>20.5M</b>

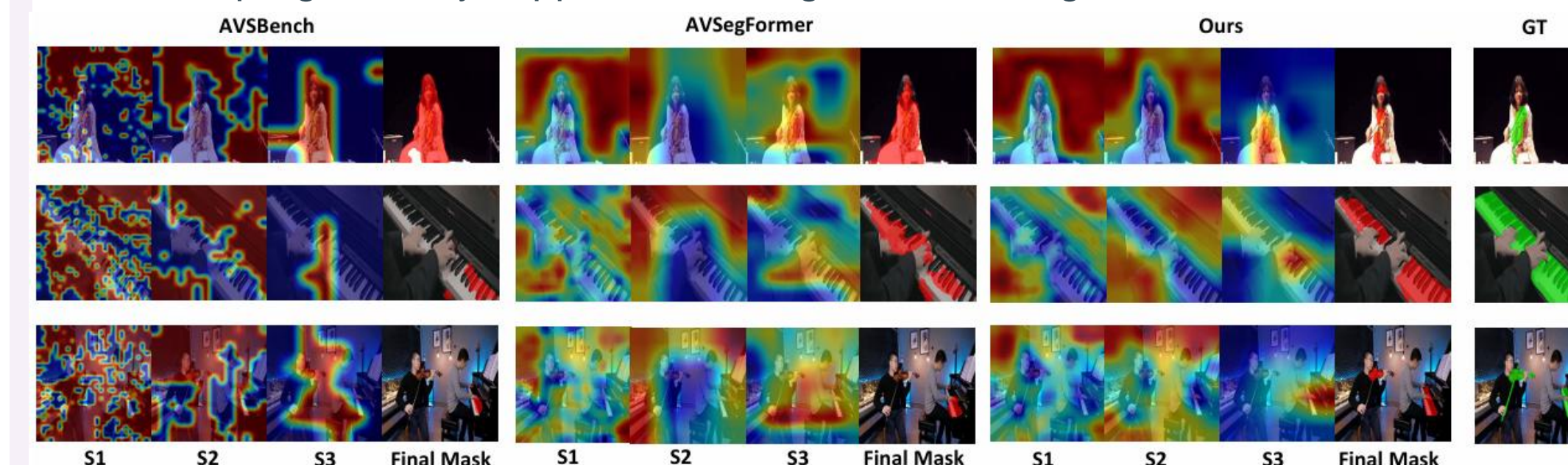
### Ablation Study

- Baseline static fusion: 44.4 MJ / 56.8 MF.
- Audio-guided visual enhancement improves to 46.8 MJ.
- Hierarchical audio refinement improves to 49.3 MJ.
- Full model with CMFD and Lmsa achieves 50.4 MJ / 62.6 MF.

	AGVE	CMFD	HAR	$\mathcal{L}_{msa}$	$M_J \uparrow$	$M_F \uparrow$
					44.4	56.8
✓					46.8	58.1
✓	✓				48.3	59.1
✓			✓		49.3	61.2
✓	✓		✓	✓	50.4	62.6

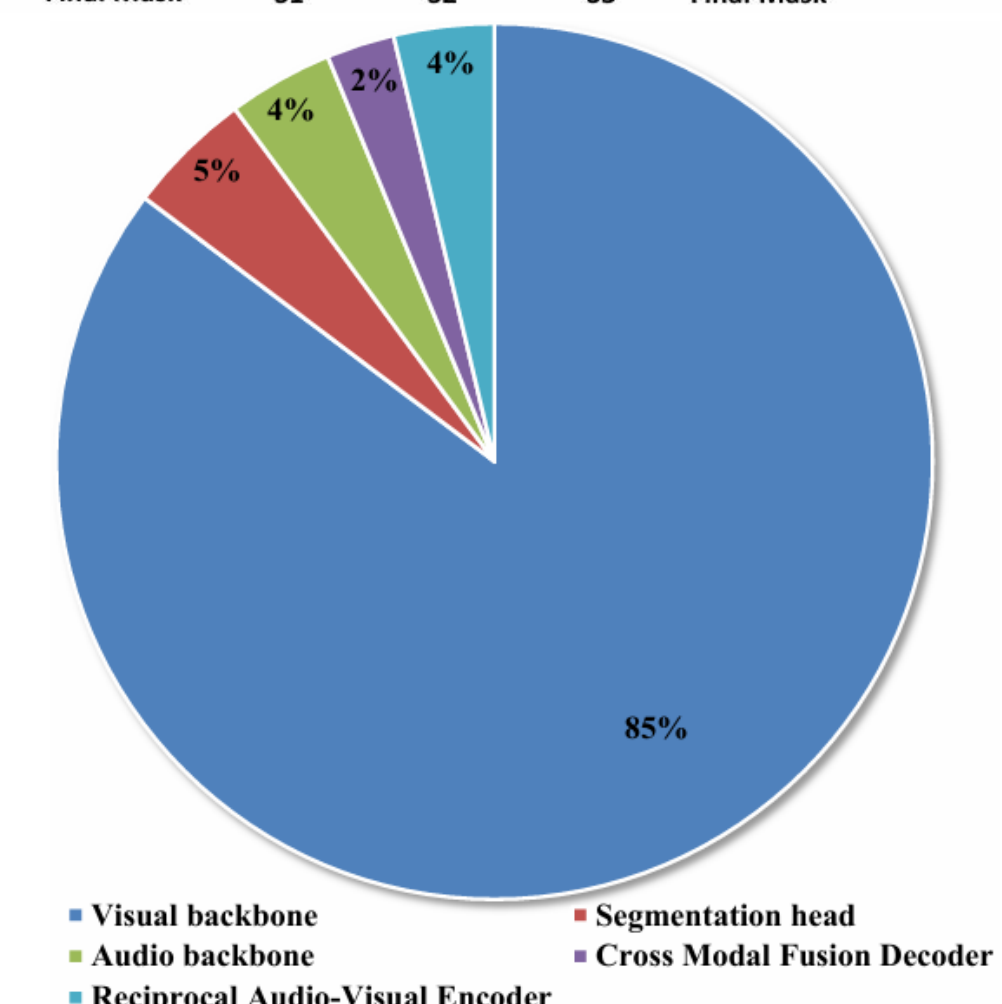
### Qualitative Result

- AVSBench is distracted by background noise.
- AVSegFormer focuses on targets but lacks clear boundaries.
- Ours shows coarse-to-fine refinement.
- Ours progressively suppresses background and aligns with GT.



### Latency Analysis

- Visual backbone dominates inference time.
- RAVE and CMFD take only 6.1 ms and 4.1 ms.
- Cross-modal interaction introduces negligible overhead.



## 4. Takeaways

### Efficient Interaction

Lightweight AVS requires efficient cross-modal interaction, not only a smaller backbone.

### What-Where Decoupling

Audio tells the model what is sounding, while visual features determine where to segment.

### Training-only Alignment

$\mathcal{L}_{msa}$  improves audio-visual consistency during training with zero inference overhead.

### Practical Mobile AVS

LightAVSeg achieves strong accuracy-efficiency trade-off, enabling high-performance AVS on mobile devices.