

# SOUPLE: ENHANCING AUDIO-VISUAL LOCALIZATION AND SEGMENTATION WITH LEARNABLE PROMPT CONTEXTS



**Khanh Binh Nguyen<sup>1</sup>**  
Deakin University



**Chae Jung Park<sup>2\*</sup>**  
National Cancer Center

## BACKGROUND

### AUDIO-VISUAL SOUND SOURCE LOCALIZATION

Pinpointing origins of sound in intricate visual scenes remains challenging. It typically relies on aligning audio-visual representations as self-supervision signals in contrastive learning.

### CLIP IN AUDIO-VISUAL LEARNING

Large-scale pre-trained models exhibit robust multimodal representations. Extensions like AudioCLIP or Wav2CLIP incorporate audio.



How to effectively ground audio with visual representation without text labels?

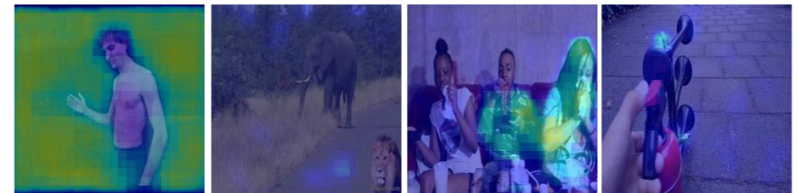
## 🎯 LIMITATION OF EXISTING WORK

Previous work like ACL-SSL adapts CLIP by replacing the classification token with an audio-embedded token:

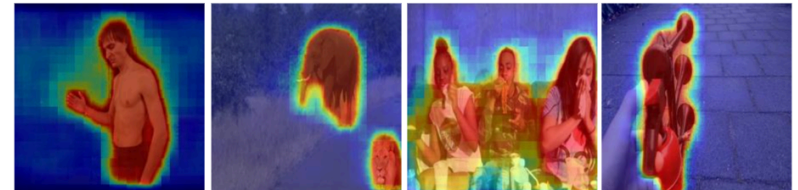
"a photo of a **[V<sub>a</sub>]**"

*This approach fails in some conditions because the fixed prompt "a photo of a" lacks semantic meaning relevant to the audio token [V<sub>a</sub>].*

ACL-SSL



Ours



**SEMANTIC MISALIGNMENT**



## PROMPT LEARNING FOR CLIP

Instead of using discrete handcrafted text, we propose **SouPLe (Sound-aware Prompt Learning)**, introducing instance-conditional, learnable context tokens.

ACL-SSL

~~"a photo of [V<sub>a</sub>]"~~

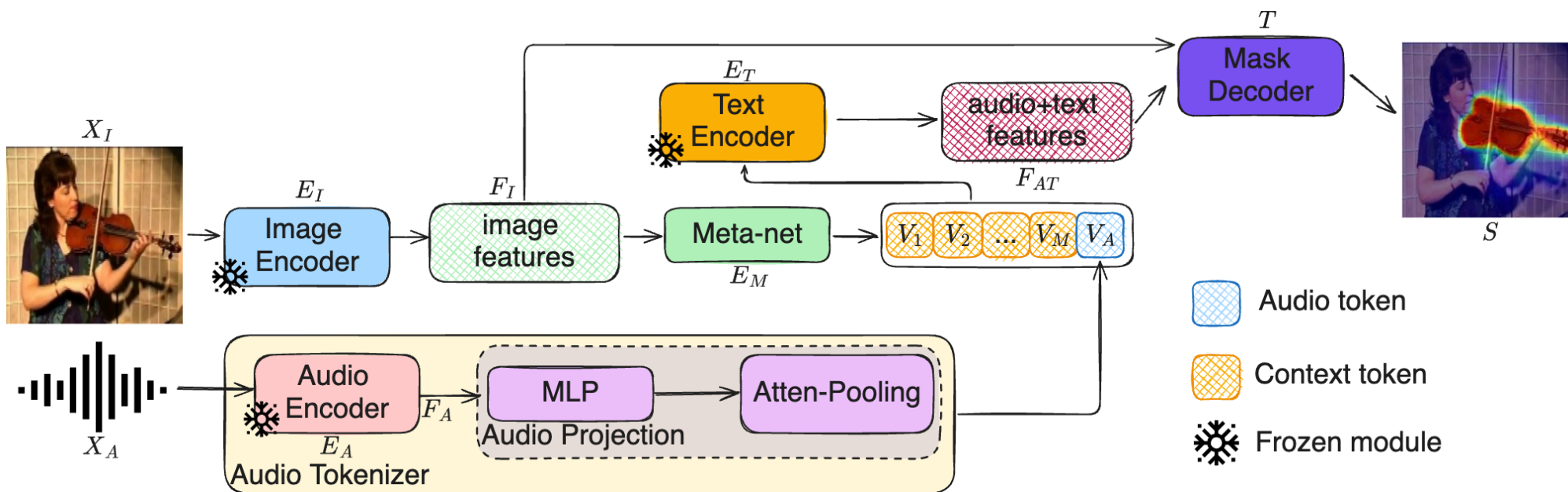


SouPLe

[V<sub>1</sub>][V<sub>2</sub>]...[V<sub>m</sub>][V<sub>a</sub>]

# SOUPLE: LEARNABLE PROMPT CONTEXTS

## METHODOLOGY: SOUPLE PIPELINE



## QUANTITATIVE RESULTS: VGG-SS

METHOD	cloU ↑	AUC ↑
Wav2CLIP	37.71	39.93
AudioCLIP	44.15	46.23
ACL-SSL (Baseline)	<b>49.46</b>	<b>46.32</b>
<b>SouPLe (Ours)</b>	<b>53.21 (+3.75)</b>	<b>48.15 (+1.83)</b>

## QUANTITATIVE RESULTS: SOUNDNET

METHOD	cloU ↑	AUC ↑
Wav2CLIP	26.00	29.60
AudioCLIP	47.20	45.22
ACL-SSL (Baseline)	<b>80.80</b>	<b>64.62</b>
<b>SouPLe (Ours)</b>	<b>84.80 (+4.00)</b>	<b>67.64 (+3.02)</b>

# SOUPLE: LEARNABLE PROMPT CONTEXTS

## ✓ OPEN-SET EVALUATION

110 HEARD CATEGORIES

### HEARD SET

cloU

# 54.76

+6.32% vs ACL-SSL

110 UNHEARD CATEGORIES

### UNHEARD SET

cloU

# 48.40

+6.42% vs ACL-SSL



PRESENTING AUTHOR

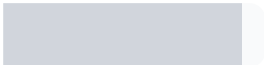
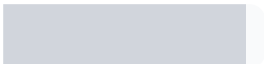
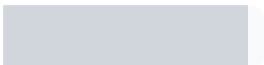
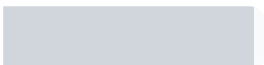

Khanh Binh Nguyen

## QUALITATIVE COMPARISONS



## ABLATION STUDIES

### TOKEN ORDER

$[V_A][V_1][V_2][V_3][V_4]$		49.91 cloU
$[V_1][V_A][V_2][V_3][V_4]$		50.71 cloU
$[V_1][V_2][V_A][V_3][V_4]$		51.08 cloU
$[V_1][V_2][V_3][V_A][V_4]$		52.41 cloU
$[V_1][V_2][V_3][V_4][V_A]$		<b>53.21</b> cloU

### CONTEXT LENGTH

<b>ctx=4</b>	<b>53.21</b> cloU
ctx=8	52.01 cloU
ctx=16	51.08 cloU





## CONCLUSION

We introduced **SouPLe**, utilizing prompt learning to bridge the semantic gap in self-supervised audio-visual localization.

By generating **instance-conditional** learnable context tokens via visual guidance, it gracefully surpasses current state-of-the-art baselines globally.

## THANK YOU

Questions?

