

**SIGHT  
AND  
SOUND  
WORKSHOP**

# Precise Video-to-Audio Generation with Cross-Modal Alignment in Latent Space

by:

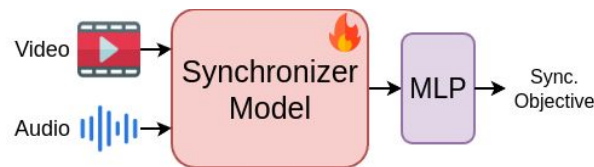
Thanh V. T. Tran<sup>1</sup>, Ngoc-Son Nguyen<sup>1</sup>, Luong Tran<sup>1</sup>, Long-Khanh Pham<sup>1</sup>,  
Paarth Neekhara<sup>2</sup>, Shehzeen Hussain<sup>2</sup>, Van Nguyen<sup>1</sup>



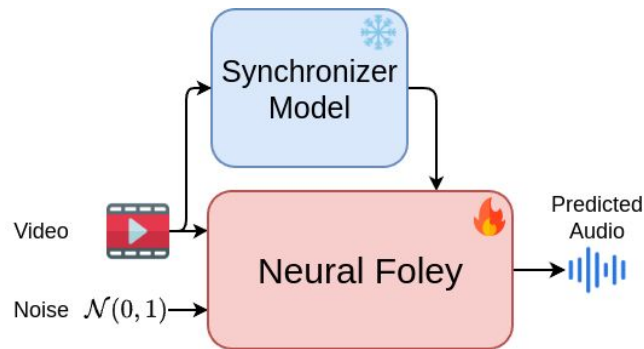
## Structural Limitations

Current V2A methods (right) rely on multi-stage training pipelines or external synchronization modules, leading to high computational overhead.

### Stage 1: Pre-training (Feature Extraction Model)



### Stage 2: Application (Foley Generation with Pre-trained Injection)









## Core Contribution

To remove the dependency on external modules

⇒ Introduce **Progressive Soft-masked Cross-Attention (PSCA)**

*Built directly into the attention mechanism to ensure high-fidelity synchronization in the latent space.*

## Mathematical Framework

Cross-attention mechanism

$$\text{PSCA}_\ell(Q_{\text{aud}}, K_{\text{vis}}, V_{\text{vis}}) = \text{softmax}\left(\frac{Q_{\text{aud}}K_{\text{vis}}^T}{\sqrt{d_k}} + \log(M^{(\ell)} + \epsilon)\right)V_{\text{vis}}$$

Progressive Masking

$$M_{ij}^{(\ell)} = \begin{cases} 1, & d_{ij} \leq \omega, \\ \beta_\ell \mathcal{F}(d_{ij} - \omega), & \omega < d_{ij} \leq \omega + \delta, \\ 0, & d_{ij} > \omega + \delta. \end{cases}$$

Decay Equation

$$\mathcal{F}(m) = \frac{1}{2} \left[ \cos\left(\frac{\pi m}{\delta}\right) + 1 \right]$$

# Semantic Bottleneck

CVPR  
JUNE 3-7, 2026



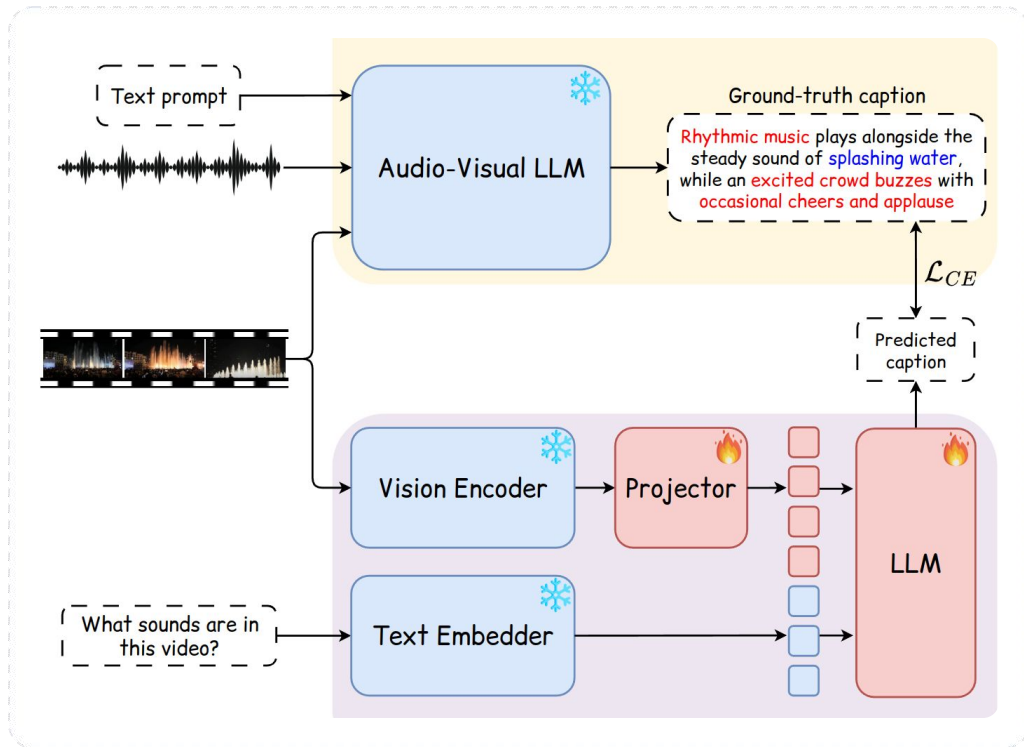
DENVER  
COLORADO

## Data Limitations

Existing V2A datasets provide only short labels, limiting the semantic granularity available to the model.

## SoundCap uses:

- **AV-LLM Teacher:** Generates detailed audio descriptions
- **VLM Student:** Learns to infer sounds from visual alone for inference
- **Noise robustness:** Instructional warnings filter irrelevant background noises (e.g., speech)



SoundCap: Training and Inference Process

# Experimental Results

CVPR  
JUNE 3-7, 2026



DENVER  
COLORADO

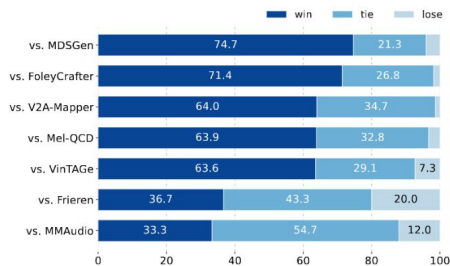
## Objective Benchmarking Results

Method	Params	Distribution Matching			Quality	Semantic Alignment( $\times 100$ )		Sync.
		KAD $\downarrow$	FAD $\downarrow$	KL $\downarrow$	IS $\uparrow$	IB-Score $\uparrow$	LB-Score $\uparrow$	Align Acc $\uparrow$
Frieren [16] <sup>†</sup>	159M	1.27	12.8	2.82	12.02	22.45	19.09	<b>97.13</b>
FoleyCrafter [17] <sup>†</sup>	1.22B	1.54	19.17	2.19	15.09	25.75	24.66	77.15
V2A-Mapper [14] <sup>‡</sup>	229M	1.34	11.73	2.50	12.43	22.38	22.32	79.08
MDSGen [12] <sup>†</sup>	131M	5.33	39.68	2.85	6.87	17.75	19.05	<u>91.70</u>
Mel-QCD [15] <sup>†</sup>	859M	1.53	19.17	2.09	10.32	23.79	23.80	73.85
VinTAGe [9] <sup>†</sup>	1.32B	1.08	17.88	2.15	17.34	21.10	21.51	67.11
MMAudio [3] <sup>◦</sup>	157M	0.57	7.89	1.91	12.68	28.09	21.98	89.73
+ SoundCap	157M	( $\uparrow 31.6\%$ ) <b>0.39</b>	( $\uparrow 10.1\%$ ) <b>7.09</b>	( $\uparrow 18.3\%$ ) <b>1.56</b>	( $\uparrow 15.8\%$ ) 14.68	( $\uparrow 2.7\%$ ) 28.85	( $\uparrow 3.0\%$ ) 22.64	( $\uparrow 0.8\%$ ) 90.53
Flowley	169M	<u>0.42</u>	7.65	<u>1.57</u>	<u>18.25</u>	<u>29.32</u>	<u>24.87</u>	89.37
+ SoundCap	169M	( $\uparrow 14.4\%$ ) <b>0.39</b>	( $\uparrow 1.7\%$ ) <u>7.52</u>	( $\uparrow 0.6\%$ ) <b>1.56</b>	( $\uparrow 7.8\%$ ) <b>19.68</b>	( $\uparrow 2.6\%$ ) <b>30.07</b>	( $\uparrow 1.8\%$ ) <b>25.33</b>	( $\uparrow 0.7\%$ ) 90.02

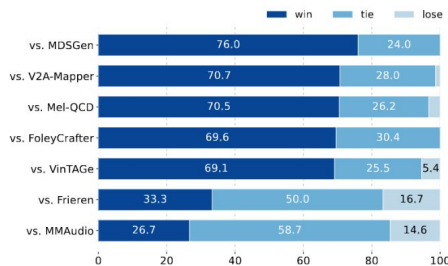
## Key Takeaways:

- **Flowley** demonstrates superior on both subjective and objective metrics.
- **SoundCap** enhances overall performance by addressing semantic gaps.

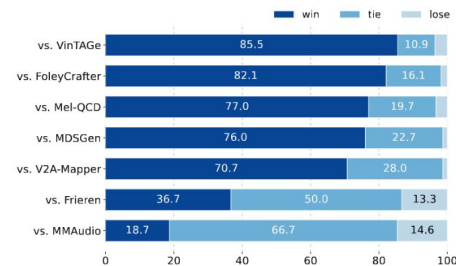
## Subjective Human Preference Results



(a) Audio Quality.



(b) Semantic Alignment.



(c) Temporal Alignment.

# Ablation Study

## Ablation on Single-Stream Block

#	CA <sub>t</sub>	CA <sub>v</sub>	PSCA <sub>v</sub>	KAD↓	IS↑	IB-Score↑	Align Acc↑
1	✗	✗	✗	0.44	16.99	27.38	86.26
2	✓	✗	✗	0.44	17.29	<u>28.93</u>	87.43
3	✗	✓	✗	<b>0.40</b>	17.58	28.02	87.04
4	✗	✗	✓	<b>0.40</b>	17.69	28.51	<u>88.72</u>
5	✓	✓	✗	<b>0.40</b>	<u>18.15</u>	28.74	87.61
6	✓	✗	✓	<u>0.42</u>	<b>18.25</b>	<b>29.32</b>	<b>89.37</b>

## Ablation on SoundCap's noise handling

Method	KAD↓	IS↑	IB-Score↑	Align Acc↑
Flowley	<u>0.42</u>	<u>18.25</u>	29.32	<u>89.37</u>
+ SoundCap w/o noise conditions	0.45	15.16	<u>29.61</u>	87.13
+ SoundCap w/ noise conditions	<b>0.39</b>	<b>19.68</b>	<b>30.07</b>	<b>90.02</b>

## Key Takeaways:

- **Significant Impact:** Removing both cross-attention streams results in a significant performance drop (#1 vs. #6).
- **Superiority:** PSCA is better than standard cross-attention (#4 vs. #6).
- **SoundCap's Robustness:** Noise awareness is important for in-the-wild datasets.

**THANK YOU!**

**CVPR**  
JUNE 3-7, 2026



**DENVER**  
**COLORADO**