

# AV-SUPERB: Audio-visual Representations and How to Evaluate Them



Shang-Wen (Daniel) Li  
FAIR  
shangwel@meta.com



Yuan (Roger) Tseng  
National Taiwan University  
r11942082@ntu.edu.tw

# Outline:

1. Why audio-visual representations & notable recent works
2. The AV-SUPERB benchmark (ICASSP 2024)
3. Some noteworthy findings
4. What's next?

# Computational systems should be able to perceive the world like we do, by combining auditory and visual inputs

- For audio event classification,  
audio events usually co-occur with visual actions

# Computational systems should be able to perceive the world like we do, by combining auditory and visual inputs

- For audio event classification,  
audio events usually co-occur with visual actions
- For speech recognition,  
visual input provides more context to what is being said

# Computational systems should be able to perceive the world like we do, by combining auditory and visual inputs

- For audio event classification,  
audio events usually co-occur with visual actions
- For speech recognition,  
visual input provides more context to what is being said



HYPOTHESIS: and then cut this piece on top of your shirt

# Computational systems should be able to perceive the world like we do, by combining auditory and visual inputs

- For audio event classification,  
audio events usually co-occur with visual actions
- For speech recognition,  
visual input provides more context to what is being said



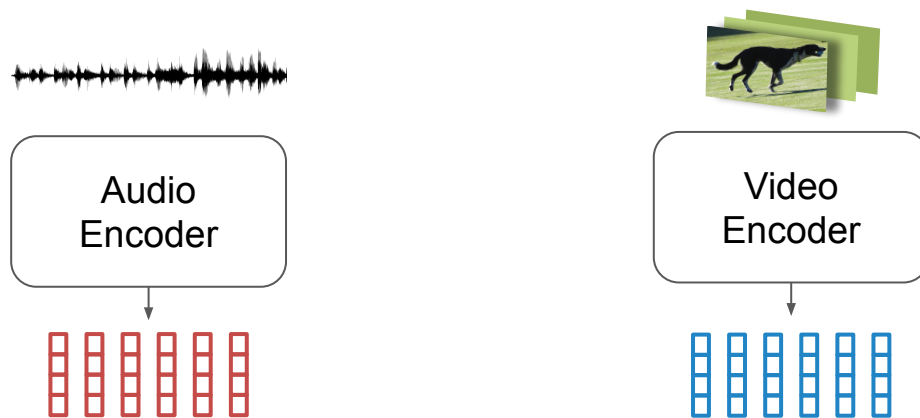
HYPOTHESIS: and then cut **the sleeves** on top of your shirt

# Computational systems should be able to perceive the world like we do, by combining auditory and visual inputs

- For audio event classification,  
audio events usually co-occur with visual actions
- For speech recognition,  
visual input provides more context to what is being said

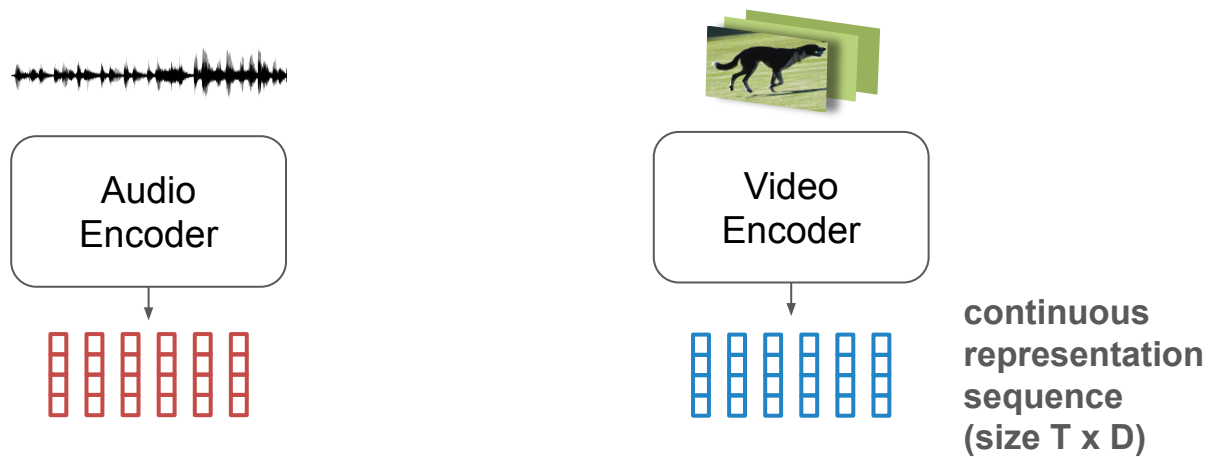
... which is why **audio-visual representation learning** is meaningful.

# Contrastive Audio-visual Learning

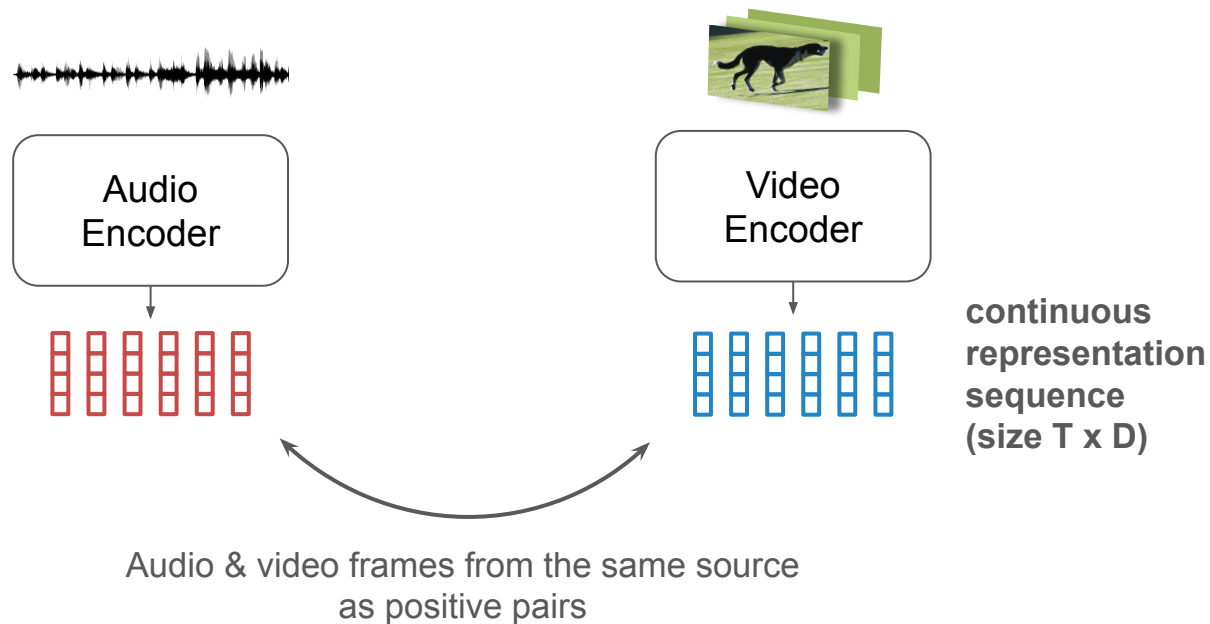




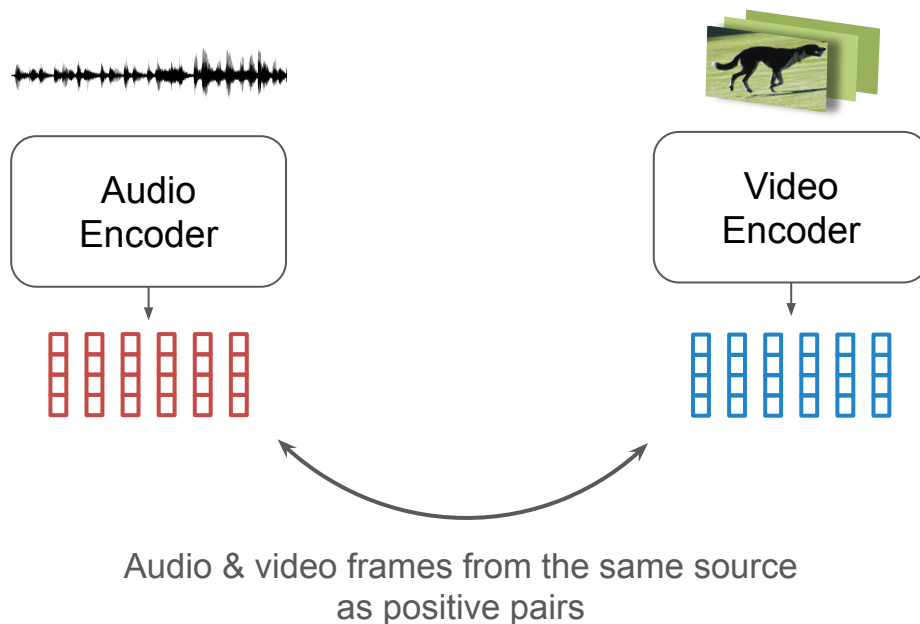
# Contrastive Audio-visual Learning



# Contrastive Audio-visual Learning

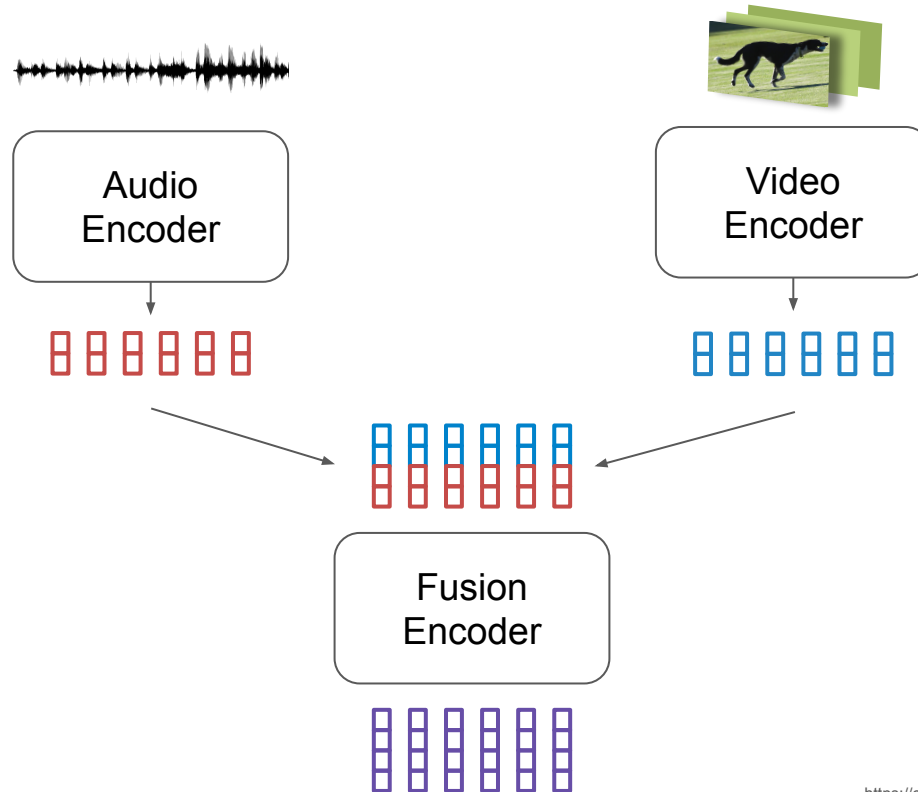


# Contrastive Audio-visual Learning

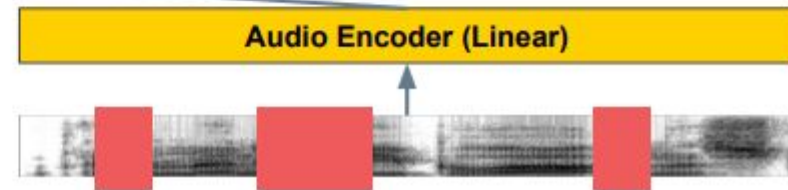
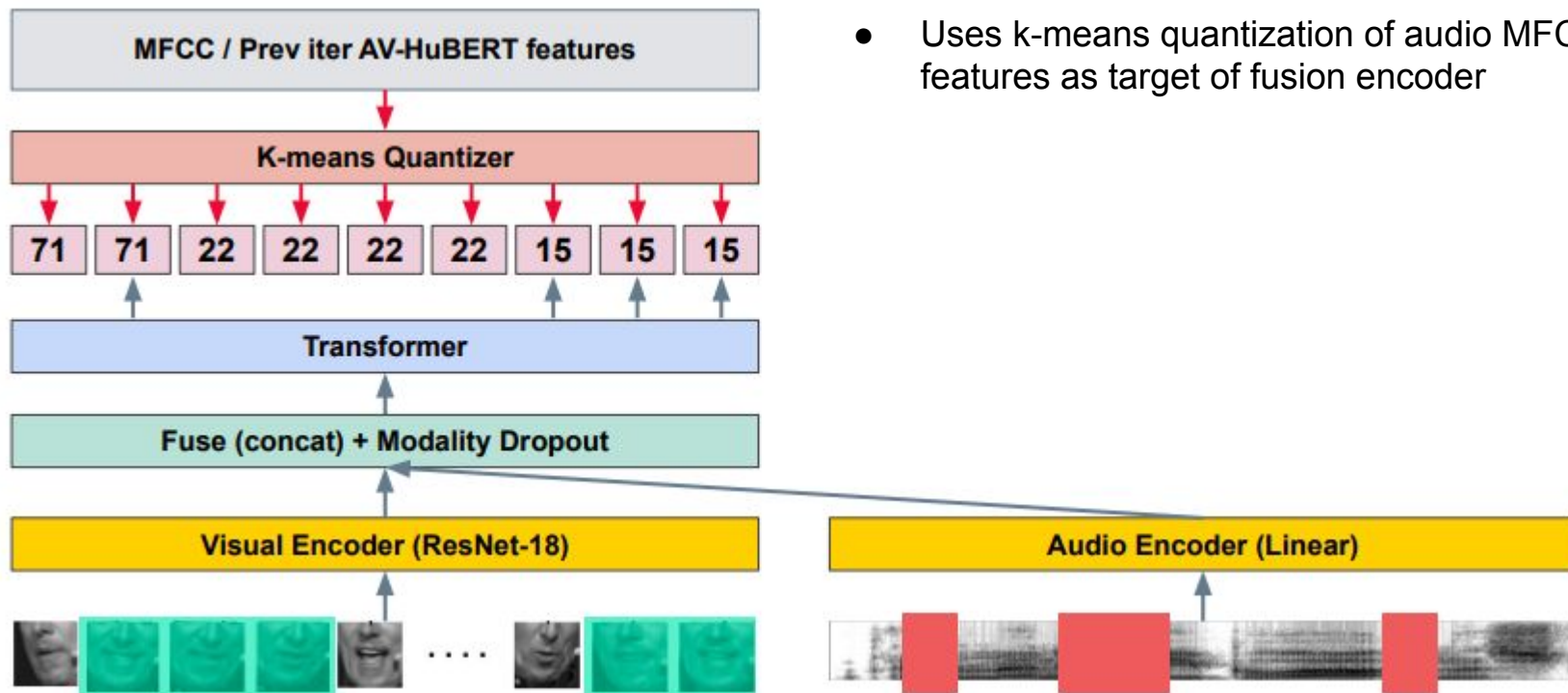


notable works include: [AVID-CMA](#) and [GDT](#) in action recognition, and [VisualVoice](#) in speech separation

# Audio-visual Fusion

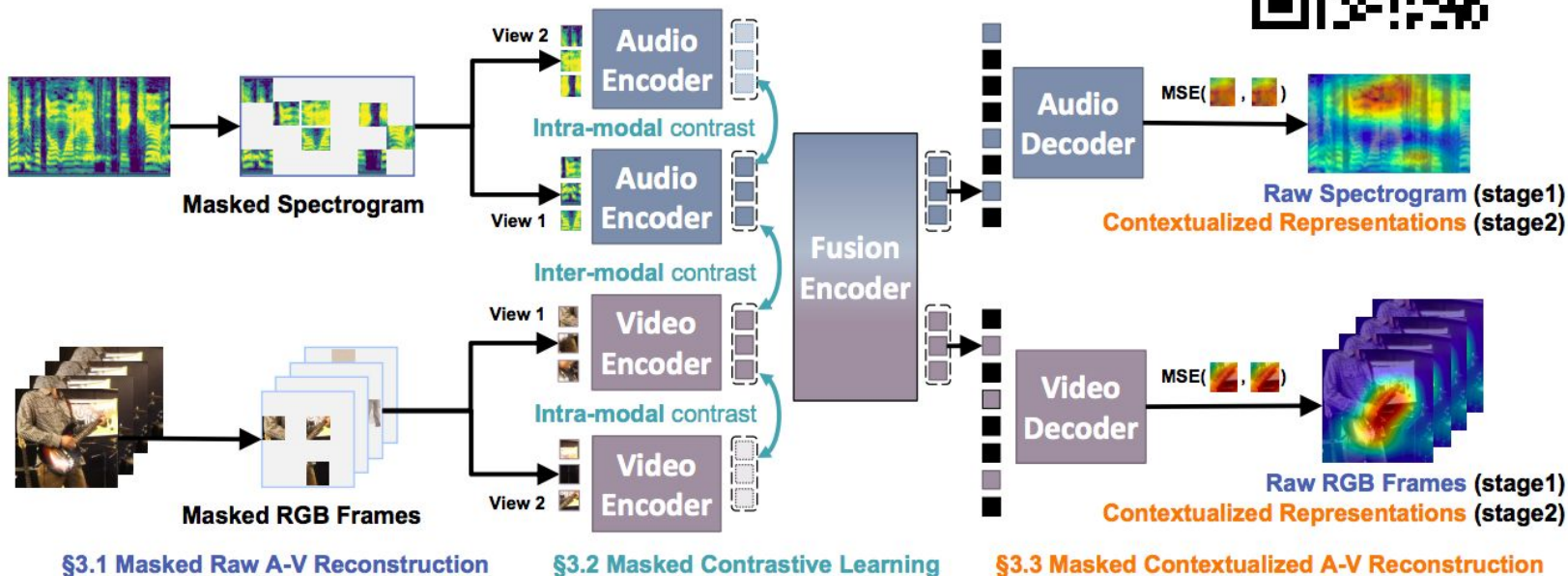


# AV-HuBERT



# MAViL (NeurIPS 2023)

Paper link:

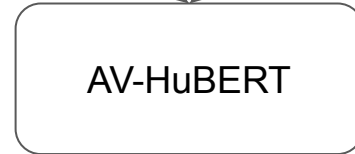


- new state-of-the-art on audio event classification & audio-to-video retrieval

# Existing audio-visual models are designed for different tasks



Audio Event:  
Electric Guitar



Speech Recognition:  
"We already use artificial intelligence..."

We can do all these tasks with one system:  
Our brains!

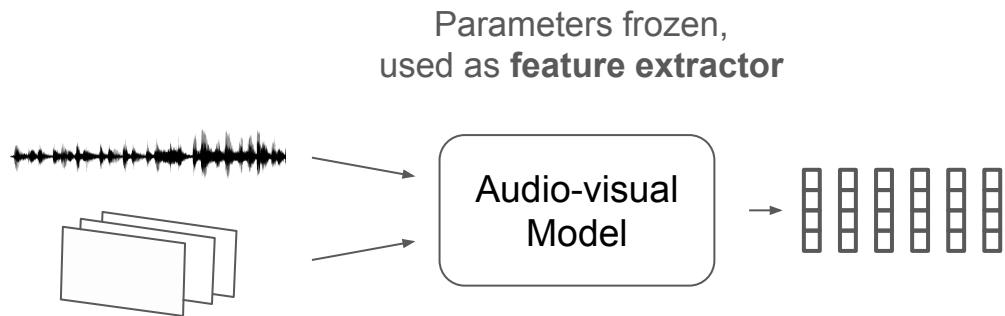
→ How far are we from a model  
that can similarly generalize?



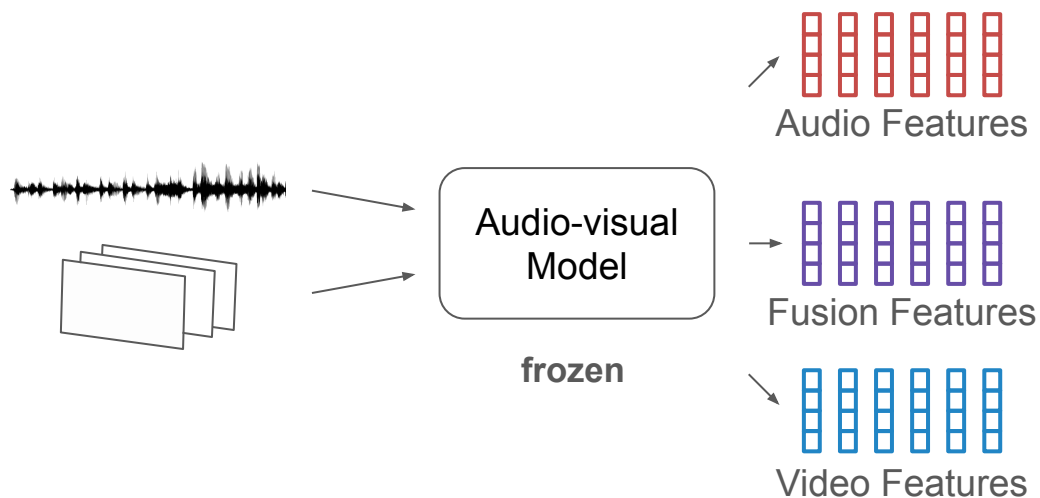
The AV-SUPERB benchmark:

# **Evaluation Protocol**

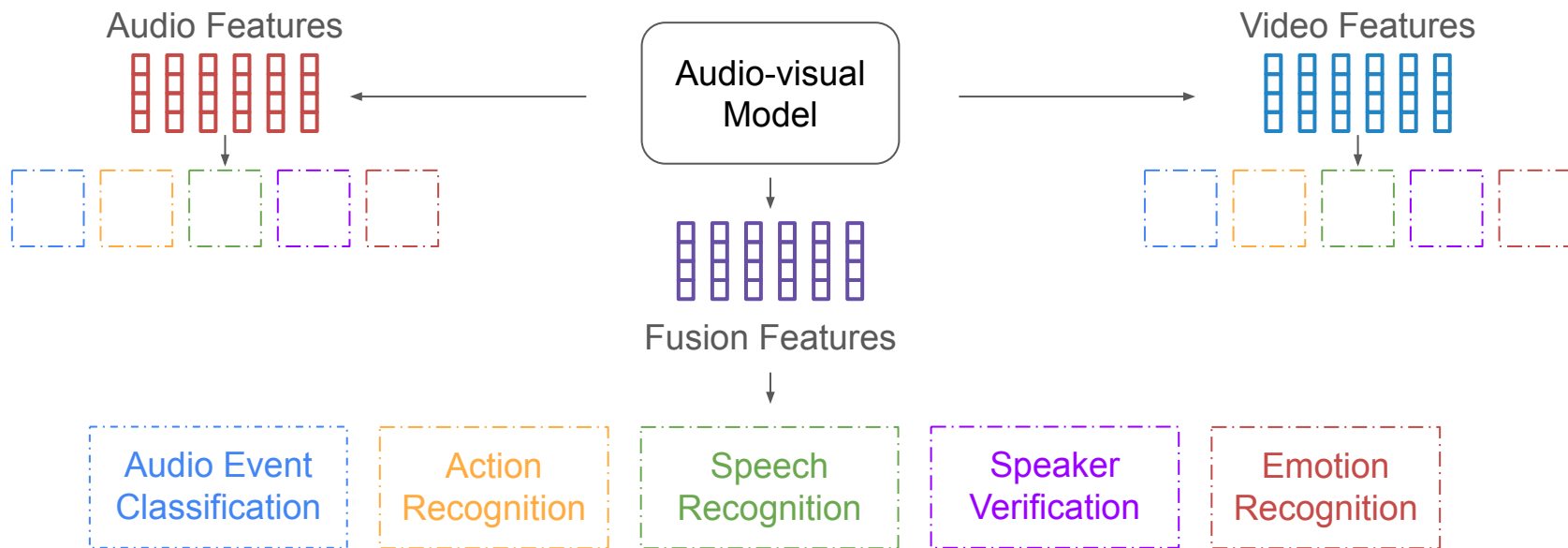
# View representations as the output of feature extractors:



# View representations as the output of feature extractors:



Each type of feature is evaluated on five tasks:



# Train a small prediction head for each task:



2-layer biLSTM



"and then cut the  
sleeves on top of  
your shirt"

Audio Event  
Classification

Action  
Recognition

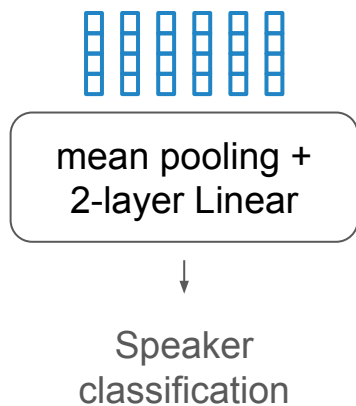
Speech  
Recognition

Speaker  
Verification

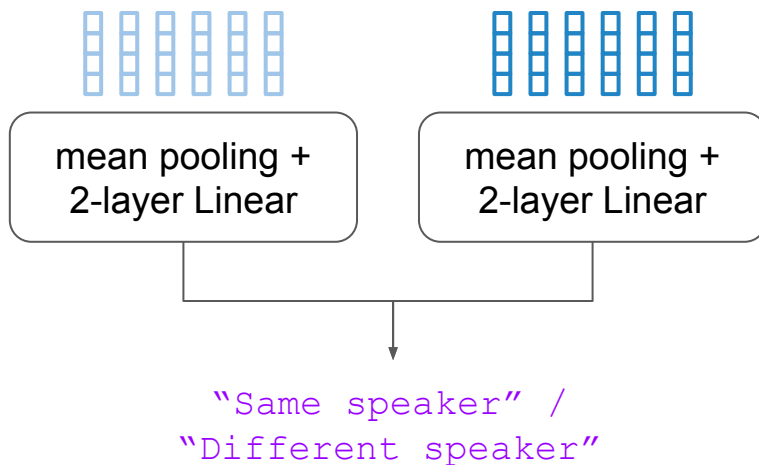
Emotion  
Recognition

# Train a small prediction head for each task:

Training:



Inference:



Audio Event  
Classification

Action  
Recognition

Speech  
Recognition

Speaker  
Verification

Emotion  
Recognition

The AV-SUPERB benchmark:  
**Some Noteworthy Findings**

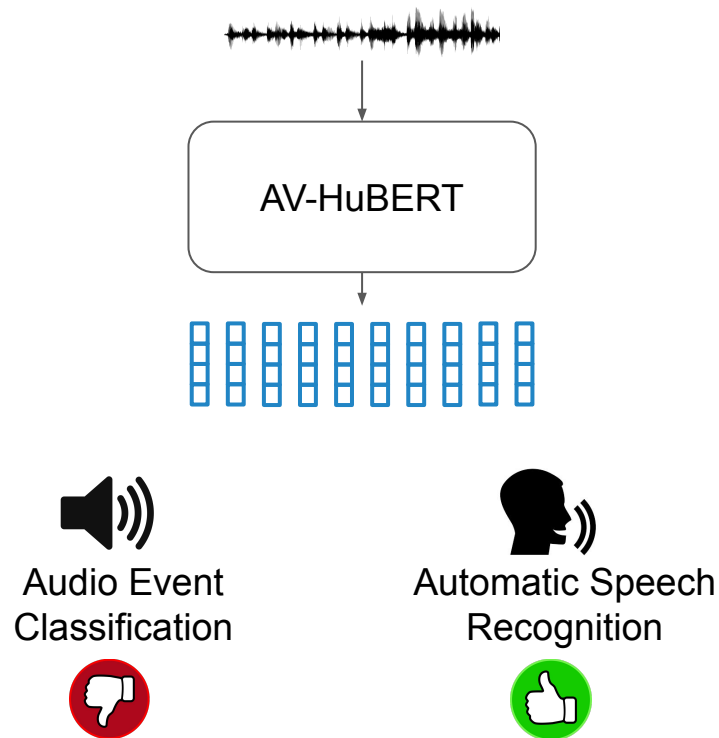
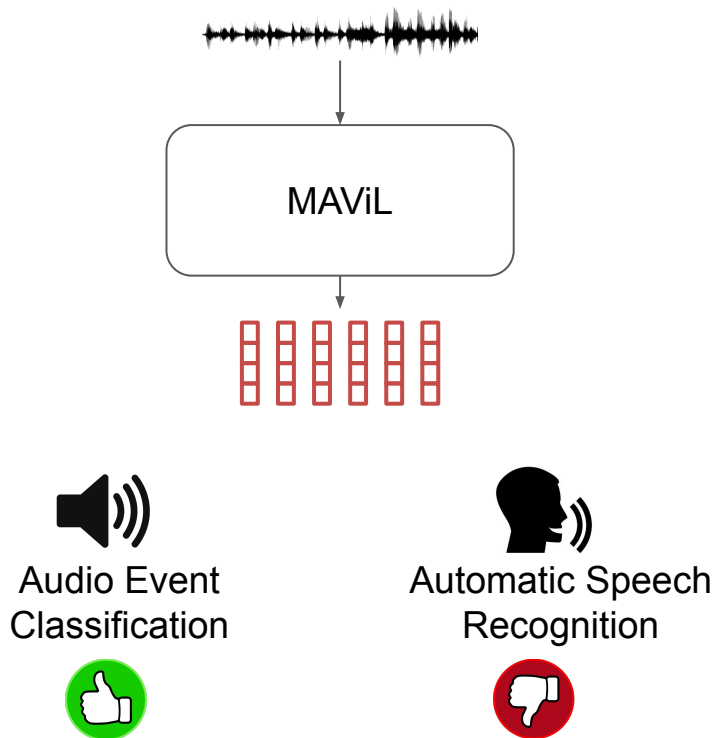
# 1. Existing pretrained SSL models do not generalize to all tasks

Representation Type	Params.	Overall Score	Audio-Visual				Speech-Visual		
			AEC		AR		ASR	ASV	ER
			AS-20K	VGGSound	Kinetics-Sounds	UCF101	LRS3-TED	VoxCeleb2	IEMOCAP
			(mAP ↑)	(Acc. ↑)	(Acc. ↑)	(Acc. ↑)	(CER ↓)	(EER ↓)	(Acc. ↑)
<i>Audio-only</i>									
FBANK	0	36.69	2.8	5.12	24.73	19.91	20.07	27.16	51.52
HuBERT								<u>15.58</u>	<b>62.14</b>
AV-HuBERT*								<b>14.45</b>	58.54
RepLAI								32.58	57.53
AVBERT								23.74	<u>60.94</u>
MAViL								20.71	59.46
<i>Video-only</i>									
HoG	0	23.57	1.3	3.81	16.76	23.87	71.46	36.32	35.83
AV-HuBERT*	103M	33.48	2.4	5.90	24.73	37.55	<b>50.91</b>	<b>11.90</b>	26.59
RepLAI	15M	36.40	5.5	13.5	46.68	56.69	<u>71.33</u>	36.95	40.72
AVBERT	37M	47.69	<u>11.5</u>	<u>28.73</u>	<u>62.67</u>	<u>77.42</u>	<u>72.29</u>	<u>20.00</u>	<b>45.8</b>
MAViL	87M	49.70	<b>18.0</b>	<b>32.08</b>	<b>74.01</b>	<b>79.37</b>	74.03	24.58	<u>43.03</u>
<i>Audio-visual fusion</i>									
AV-HuBERT	103M	53.42	13.3	32.69	52.23	41.46	<b>2.75</b>	<b>9.46</b>	46.45
AVBERT	43M	54.85	<u>22.9</u>	<u>44.54</u>	<u>71.31</u>	<u>71.76</u>	70.12	<u>18.31</u>	<b>61.87</b>
MAViL	187M	62.36	<b>26.7</b>	<b>47.22</b>	<b>79.51</b>	<b>77.98</b>	<u>30.18</u>	19.67	<u>54.94</u>

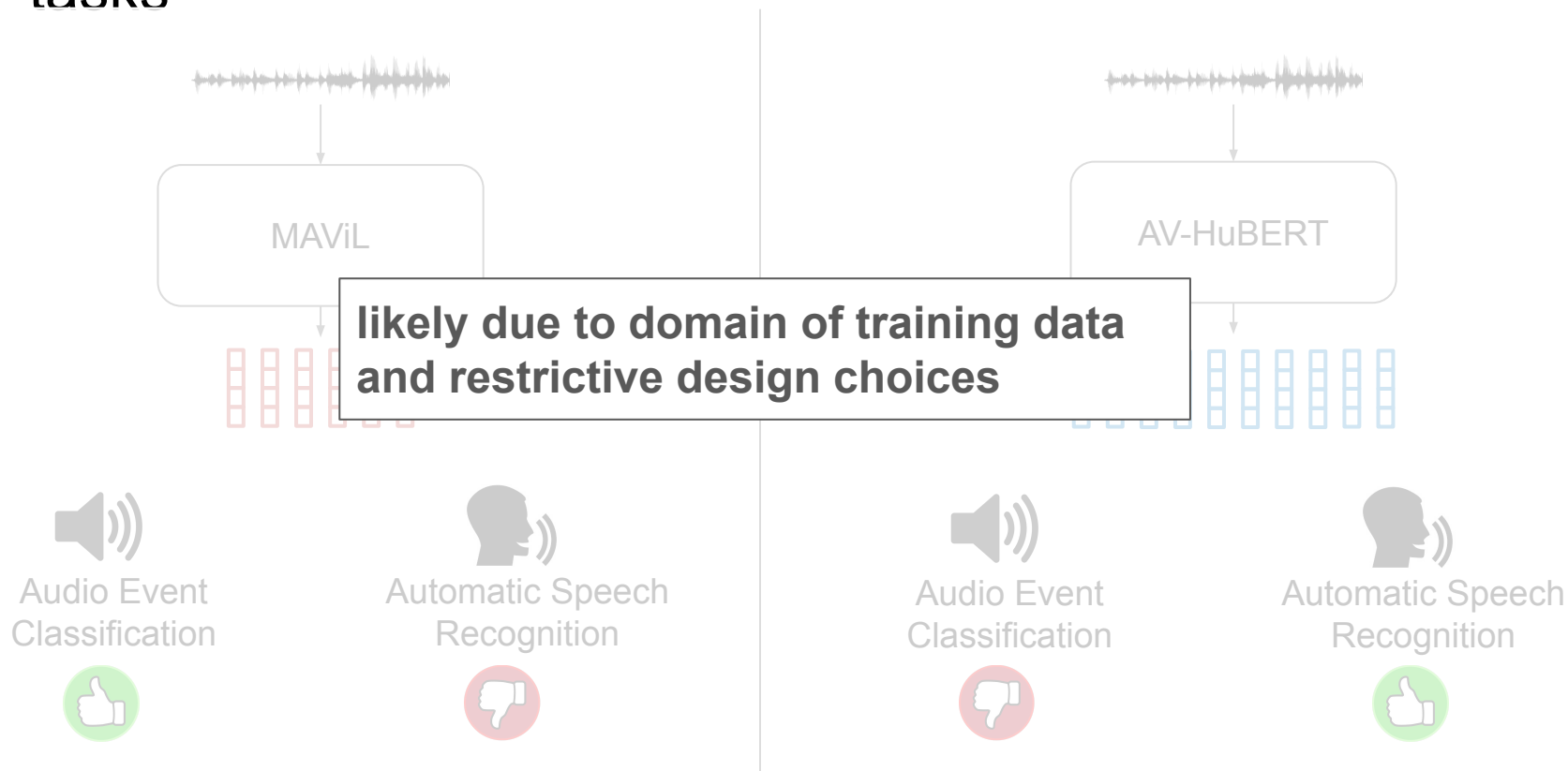
Evaluate 5 pretrained SSL models from different domains, with handcrafted features as baselines



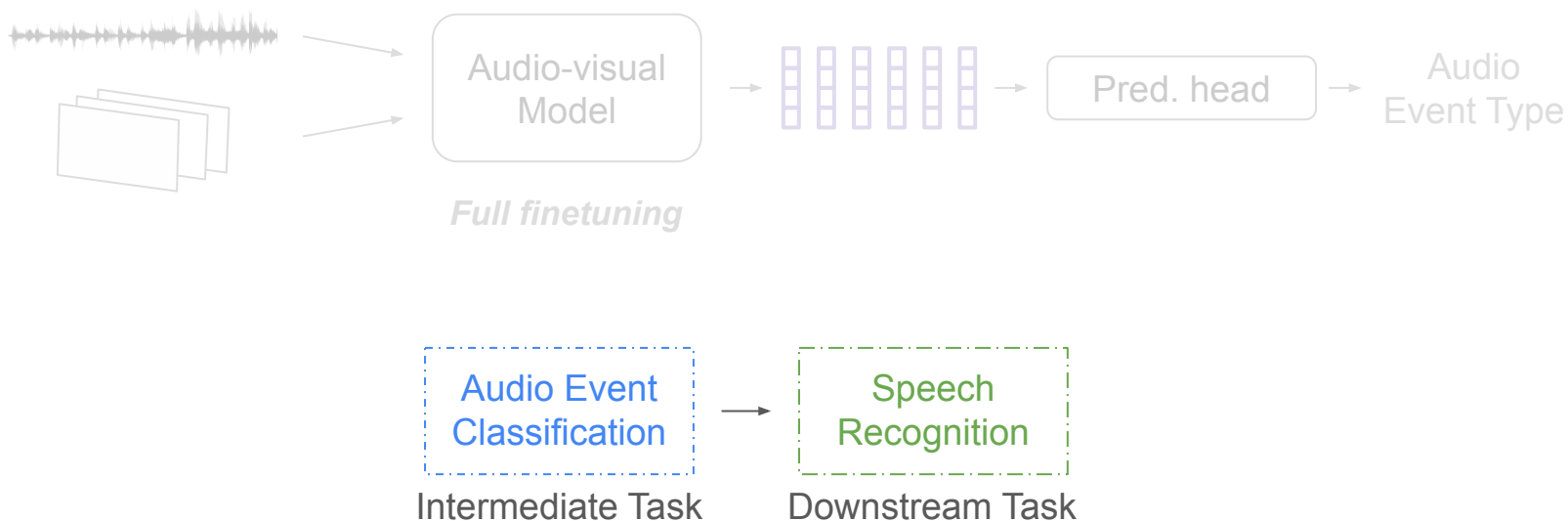
# 1. Existing pretrained SSL models do not generalize to all tasks



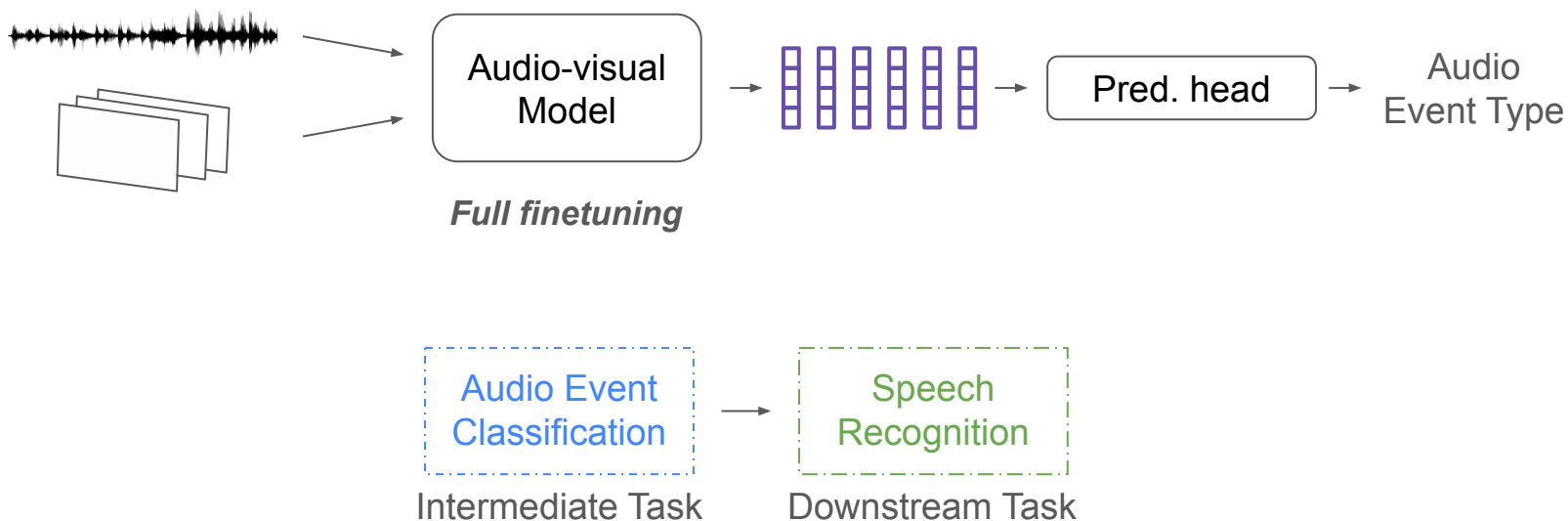
# 1. Existing pretrained SSL models do not generalize to all tasks



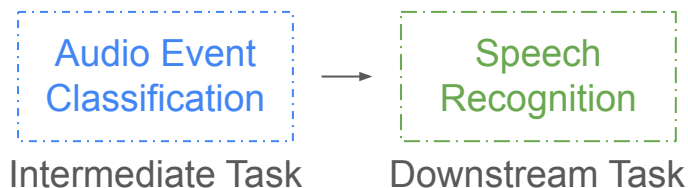
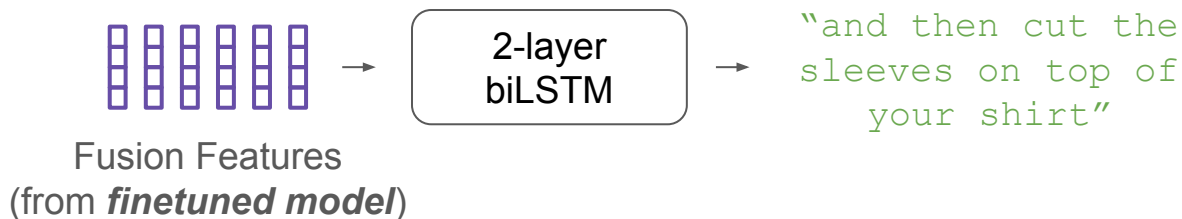
## 2. Intermediate-task finetuning can help but does not completely solve the problem



## 2. Intermediate-task finetuning can help but does not completely solve the problem



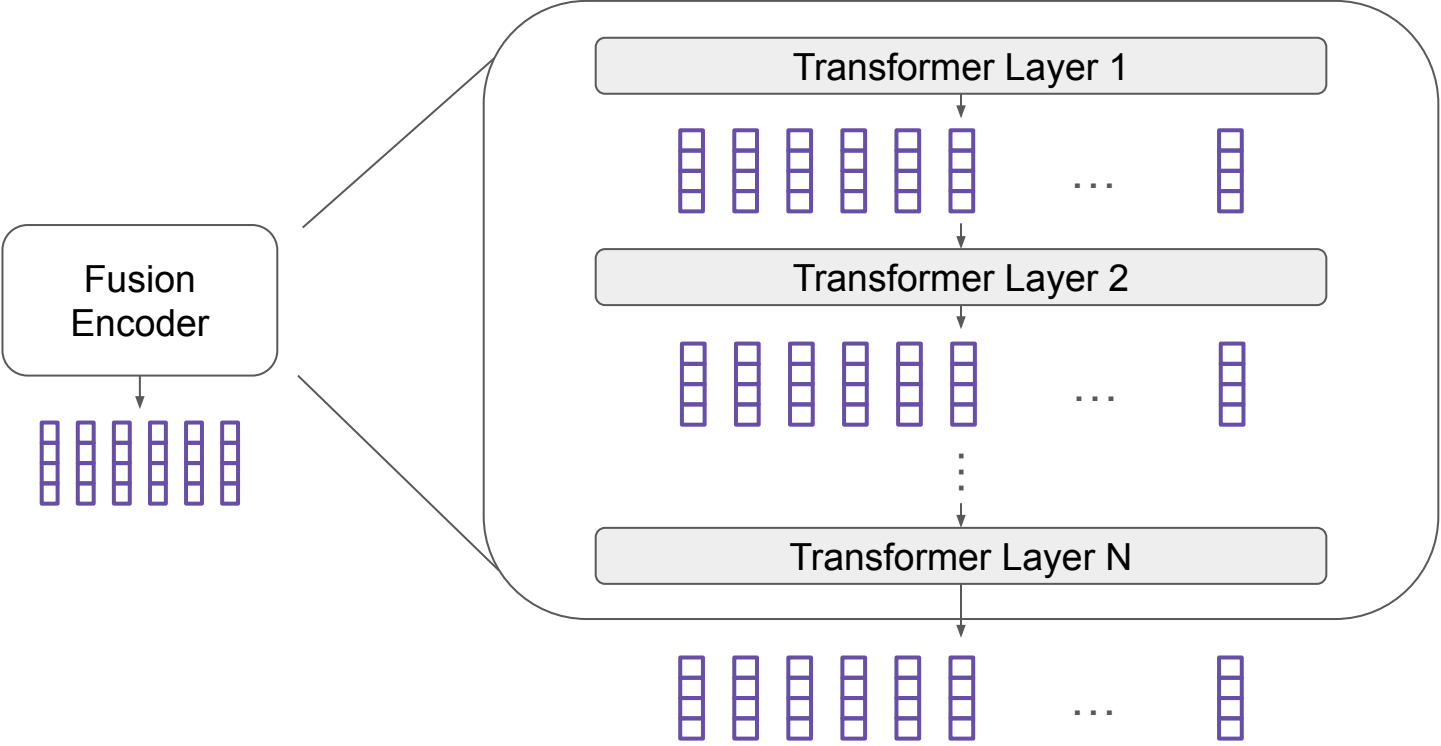
## 2. Intermediate-task finetuning can help but does not completely solve the problem



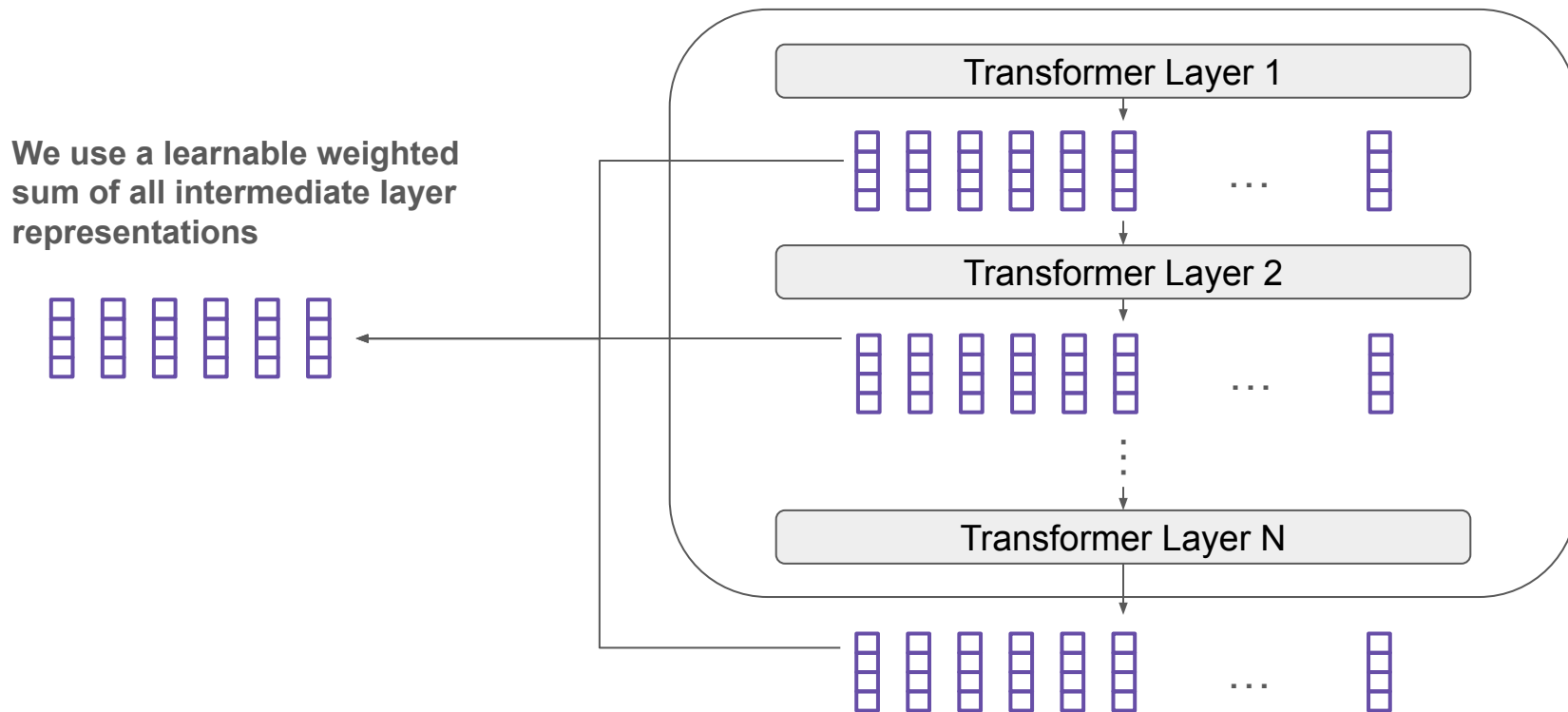
## 2. Intermediate-task finetuning can help but does not completely solve the problem

Intermediate Task Fine-tuning Data		Audio-Visual				Speech-Visual		
		AEC		AR		ASR	ASV	ER
		AS-20K (mAP $\uparrow$ )	VGGSound (Acc. $\uparrow$ )	Kinetics-Sounds (Acc. $\uparrow$ )	UCF101 (Acc. $\uparrow$ )	LRS3-TED (CER $\downarrow$ )	VoxCeleb2 (EER $\downarrow$ )	IEMOCAP (Acc. $\uparrow$ )
<i>MAViL</i>								
Audio	AudioSet-2M	28.3(+6.7)	44.79(+4.89)	62.93(+5.65)	50.10(+4.42)	23.99(+0.44)	21.77(-1.06)	58.17(-1.29)
Video		20.9(+2.9)	36.68(+4.58)	77.39(+3.38)	86.93(+7.56)	78.59(-4.56)	23.93(+0.65)	39.15(-3.88)
Fusion		39.1(+12.4)	55.94(+8.72)	84.93(+5.42)	88.07(+10.09)	30.65(-0.47)	18.61(+1.06)	46.35(-8.59)

### 3. Representations from the last layer may be suboptimal



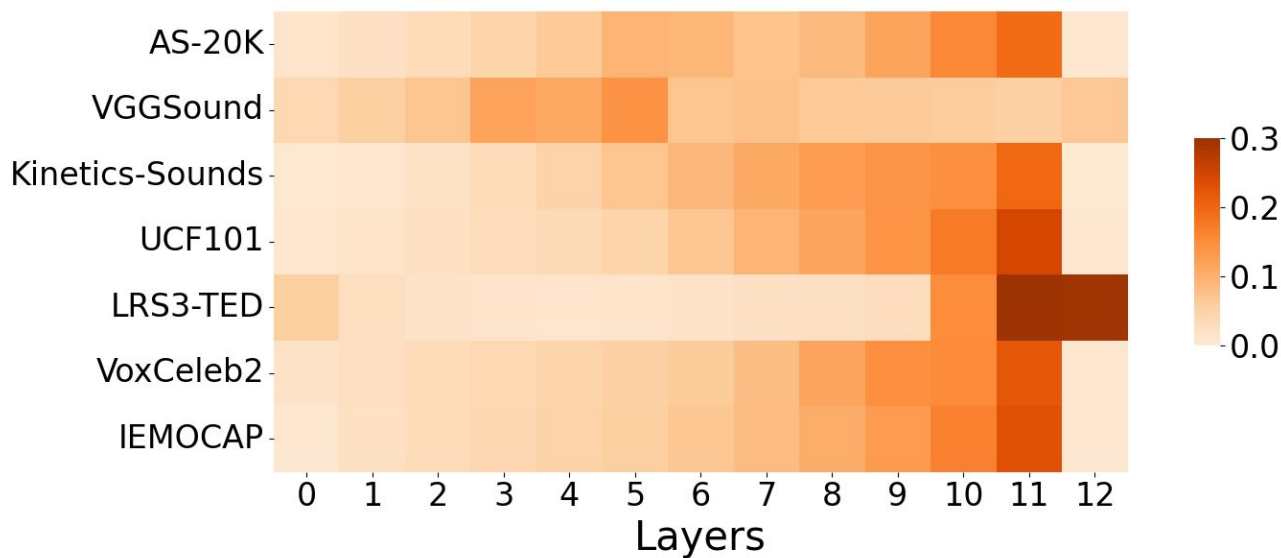
### 3. Representations from the last layer may be suboptimal





### 3. Representations from the last layer may be suboptimal

For AV-HuBERT fusion features, the penultimate layer often contributes more



Heatmap of learned weights for each downstream dataset

**What's next?**

- Including more realistic, useful, or fundamental tasks, such as retrieval, localisation, etc.
- Fairer comparison of models by unifying their training data, objective functions, or model architectures.

# Check out our preprint on arXiv:

- Accepted to ICASSP 2024!
- In progress: an evaluation platform for researchers to benchmark new models



[Paper link](#)

# Collaborators:



Yuan Tseng



Layne Berry



Yi-Ting Chen



I-Hsiang Chiu



Hsuan-Hao Lin



Max Liu



Puyuan Peng



Yi-Jen Shih



Hung-Yu Wang



Haibin Wu



Po-Yao Huang



Chun-Mao Lai



Shang-Wen Li



David Harwath



Yu Tsao



Shinji Watanabe



Abdelrahman Mohamed



Chi-Luen Feng



Hung-yi Lee

Thank you for listening!