

# Dataset Distillation for Audio-Visual Datasets

## (CVPR'24 Sight and Sound Workshop)

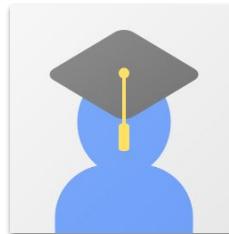
**Saksham Singh Kushwaha<sup>1</sup>, Siva Sai Nagender Vasireddy<sup>1</sup>, Kai Wang<sup>2</sup>,Yapeng Tian<sup>1</sup>**

<sup>1</sup>University of Texas at Dallas, USA, <sup>2</sup>National University of Singapore, Singapore

<sup>1</sup>{sakshamsingh.kushwaha, sivasainagender.vasireddy, yapeng.tian}@utdallas.edu, <sup>2</sup>kai.wang@nus.edu.sg



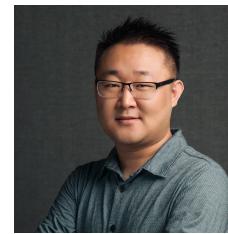
**Saksham Singh Kushwaha**  
Ph.D. Candidate



**Siva Sai Nagender  
Vasireddy**  
Ph.D. Candidate



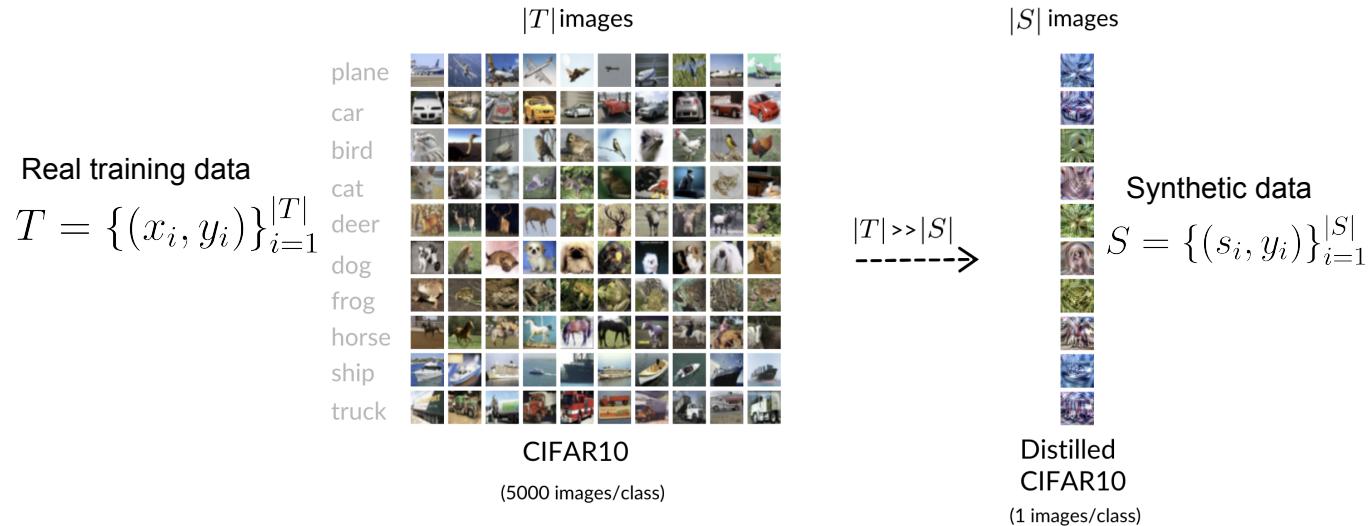
**Kai Wang**  
Ph.D. Candidate



**Yapeng Tian**  
Assistant Professor

# What is dataset distillation?

- Synthesize tiny and high-fidelity data
- Summarises most important knowledge from target dataset

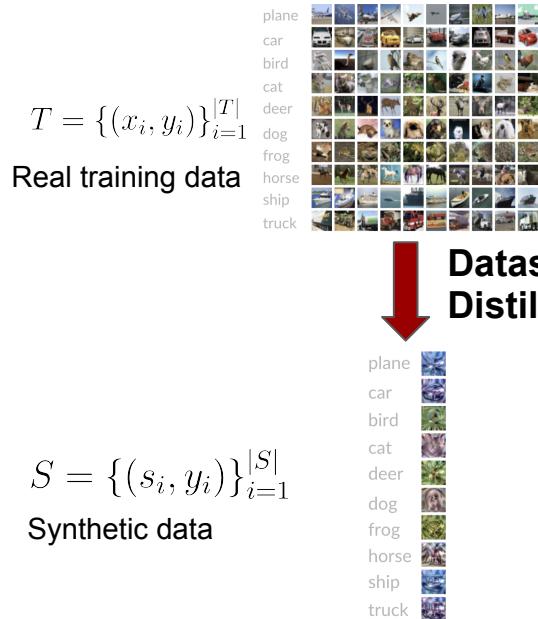


[1] Wang et. al, "Dataset Distillation", arXiv preprint, 2018.

[2] B. Zhao et. al, "Dataset Condensation with Distribution Matching", WACV, 2023.

# What is dataset distillation?

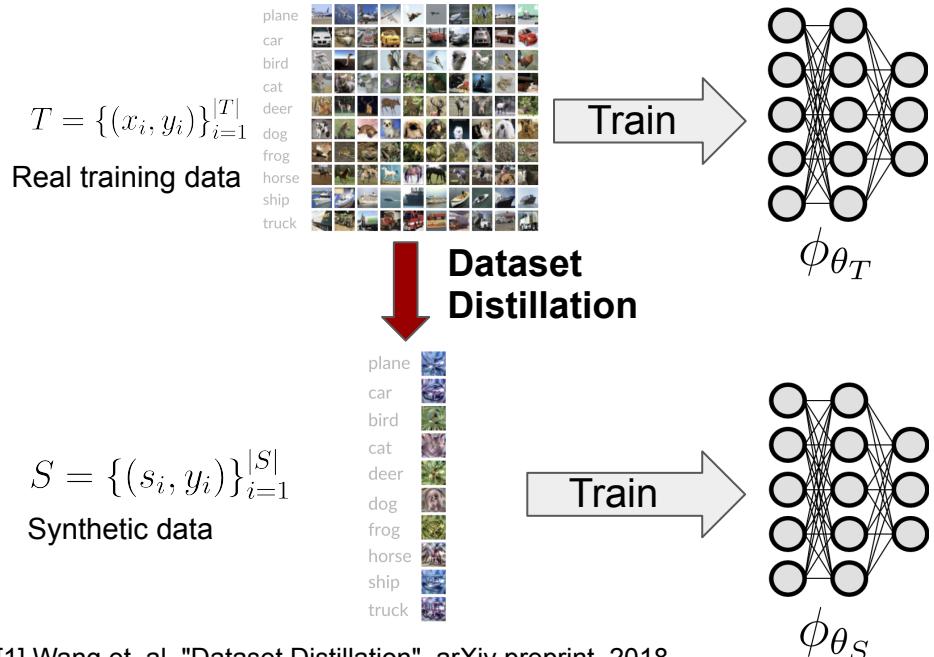
- Comparable performance on unseen real test data



[1] Wang et. al, "Dataset Distillation", arXiv preprint, 2018.

# What is dataset distillation?

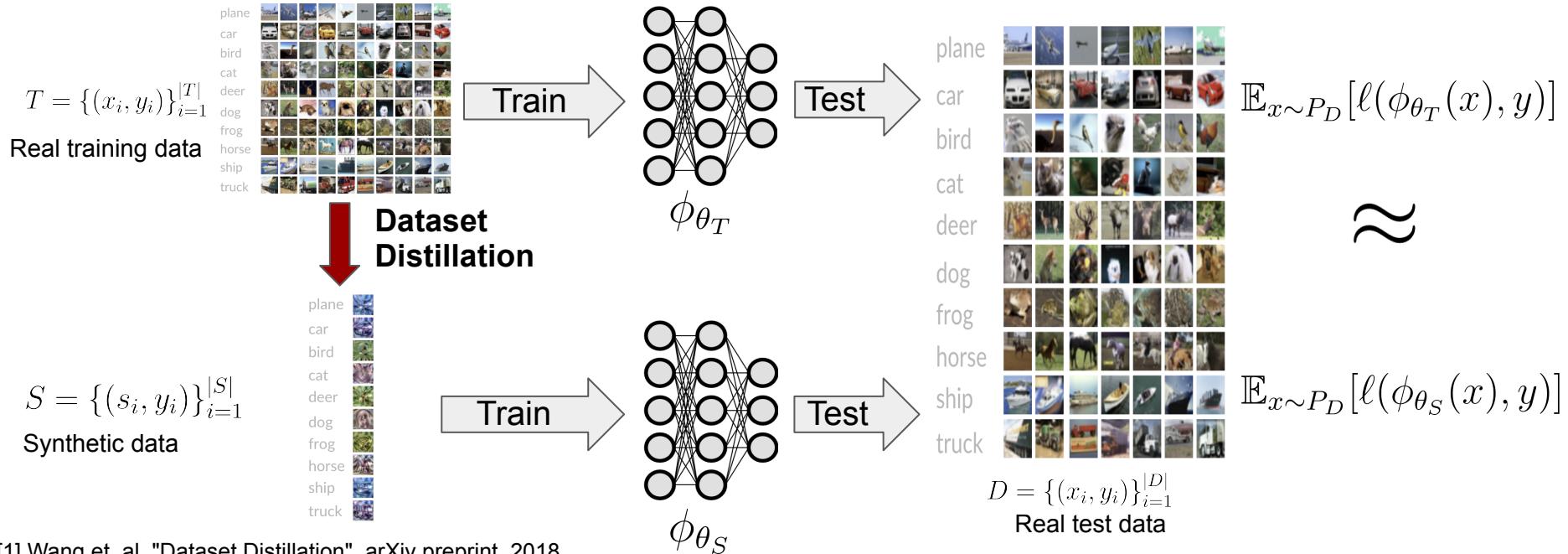
- Comparable performance on unseen real test data



[1] Wang et. al, "Dataset Distillation", arXiv preprint, 2018.

# What is dataset distillation?

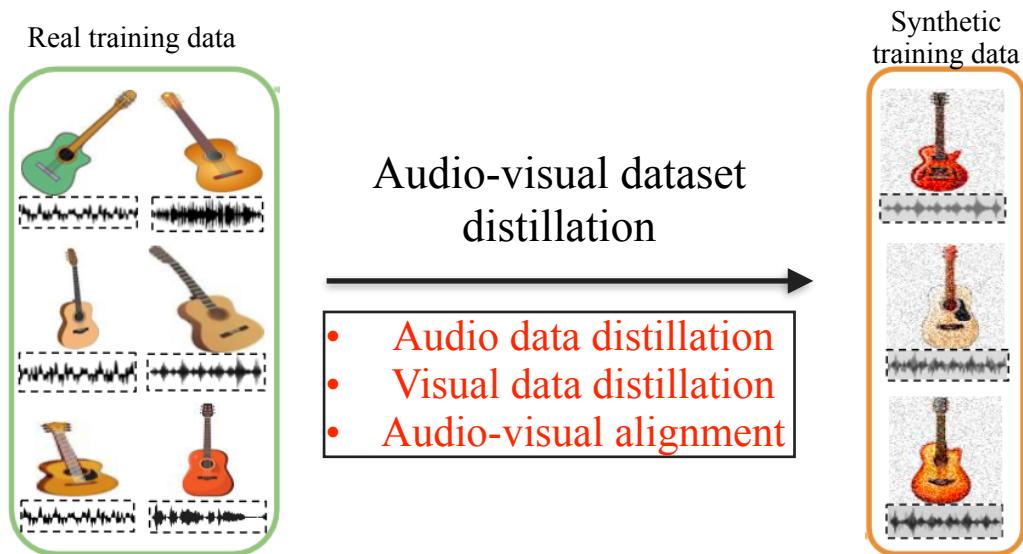
- Comparable performance on unseen real test data



[1] Wang et. al, "Dataset Distillation", arXiv preprint, 2018.

# Audio-Visual dataset distillation

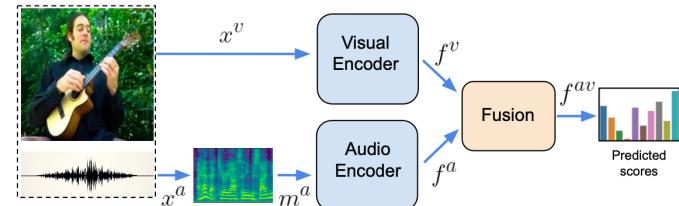
- Previous dataset distillation research focuses on image-only
- Explosion in large multimodal datasets



[1] Wang et. al, "Dataset Distillation", arXiv preprint, 2018.

# Audio-Visual dataset distillation

- Previous dataset distillation research focuses on image-only
- Explosion in large multimodal datasets
- Experimental Setup:
  - Task: **AV event recognition**
  - Datasets: **VGGSound-subset<sup>[1]</sup>** & **AVE<sup>[2]</sup>**
  - Samples: **1-sec video clips**



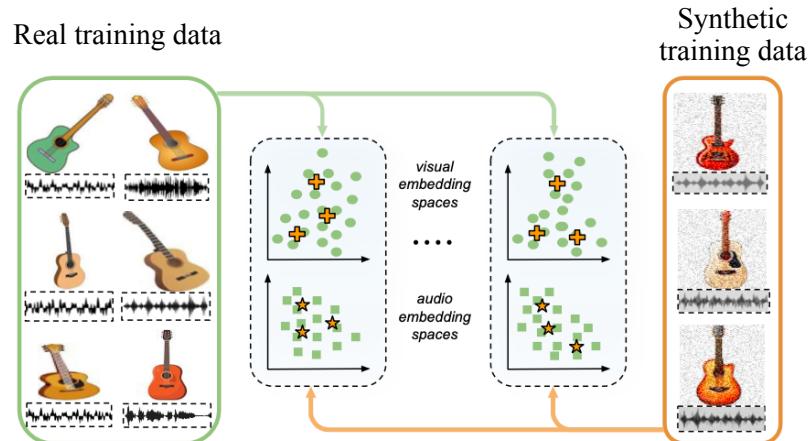
Audio Visual Event recognition

[1] H. Chen et. Al, "Vggsound: A large-scale audio-visual dataset.", ICCASP 2020.

[2] Y. Tian et. al, "Audio-visual event localization in unconstrained videos.", ECCV 2018.

# Does AV integration hold for distilled data?

- Vanilla Audio-visual data distillation
  - Naive extension of Distribution Matching approach<sup>[1]</sup>
  - **Audio-only & Visual-only** data distillation



[1] B. Zhao "Data Condensation using distribution matching" WACV 2023

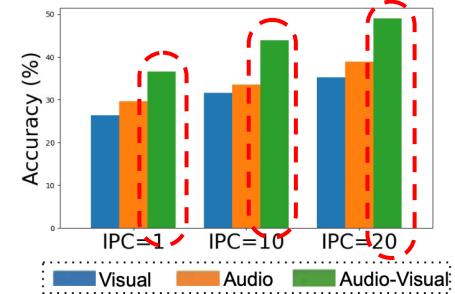
[2] G. Zhao "Improved Distribution Matching for Dataset condensation" CVPR 2023

[3] H. Zhang "M3D: Dataset condensation by minimizing maximum mean discrepancy" AAAI 2024

# Does integration hold for AV distilled data?



- Vanilla Audio-visual data distillation
    - Naive extension of Distribution Matching approach<sup>[1]</sup>
    - Audio-only & Visual-only data distillation
  - Distilled data: **Audio-visual integration helps**



IPC	Coreset Selection						Training Set Synthesis						Whole data		
	Random			Herding [7]			MTT [1]			DM [8]			A	V	AV
	A	V	AV	A	V	AV	A	V	AV	A	V	AV	A	V	AV
1	14.27 $\pm$ 0.97	11.65 $\pm$ 1.45	15.44 $\pm$ 1.87	26.32 $\pm$ 1.57	14.72 $\pm$ 2.87	20.77 $\pm$ 2.77	30.99 $\pm$ 1.48	24.15 $\pm$ 2.25	34.13 $\pm$ 3.62	29.60 $\pm$ 2.33	26.40 $\pm$ 1.10	36.54 $\pm$ 2.52			
10	32.01 $\pm$ 1.64	22.71 $\pm$ 1.57	32.50 $\pm$ 2.03	34.58 $\pm$ 1.98	28.9 $\pm$ 1.44	39.89 $\pm$ 1.64	36.57 $\pm$ 2.57	25.41 $\pm$ 1.58	36.79 $\pm$ 1.97	33.60 $\pm$ 1.35	31.63 $\pm$ 1.96	43.85 $\pm$ 1.75	62.07 $\pm$ 0.54	48.19 $\pm$ 0.54	68.24 $\pm$ 0.75
20	36.78 $\pm$ 2.88	31.05 $\pm$ 1.17	45.10 $\pm$ 2.31	44.11 $\pm$ 1.47	34.58 $\pm$ 0.84	50.20 $\pm$ 0.74	45.73 $\pm$ 1.03	29.52 $\pm$ 1.43	51.87 $\pm$ 1.26	38.93 $\pm$ 3.52	35.23 $\pm$ 1.16	49.01 $\pm$ 2.44			

[1] B. Zhao “Data Condensation using distribution matching” WACV 2023

[2] G. Zhao "Improved Distribution Matching for Dataset condensation" CVPR 2023

[3] H. Zhang "M3D: Dataset condensation by minimizing maximum mean discrepancy" AAAI 2024

# Does integration hold for AV distilled data?

- Vanilla Audio-visual data distillation
  - Naive extension of Distribution Matching approach<sup>[1]</sup>
  - Audio-only & Visual-only data distillation
- Distilled data: **Audio-visual integration helps**
- Multimodal Fusion: **Ensemble outperforms**

	Only-A	Only-V	Audio-Visual Fusion				
			Concat	Sum	Attention	Ensemble	
IPC	1	$29.60 \pm 2.33$	$26.40 \pm 1.10$	$33.77 \pm 1.65$	$34.72 \pm 1.27$	$9.97 \pm 0.83$	<b><math>36.54 \pm 2.52</math></b>
	10	$33.60 \pm 1.35$	$31.63 \pm 1.96$	$41.71 \pm 1.27$	$40.49 \pm 1.83$	$10.11 \pm 0.35$	<b><math>43.85 \pm 1.75</math></b>
	20	$38.93 \pm 3.52$	$35.23 \pm 1.16$	$46.59 \pm 1.34$	$46.05 \pm 1.74$	$11.10 \pm 1.88$	<b><math>49.01 \pm 2.44</math></b>

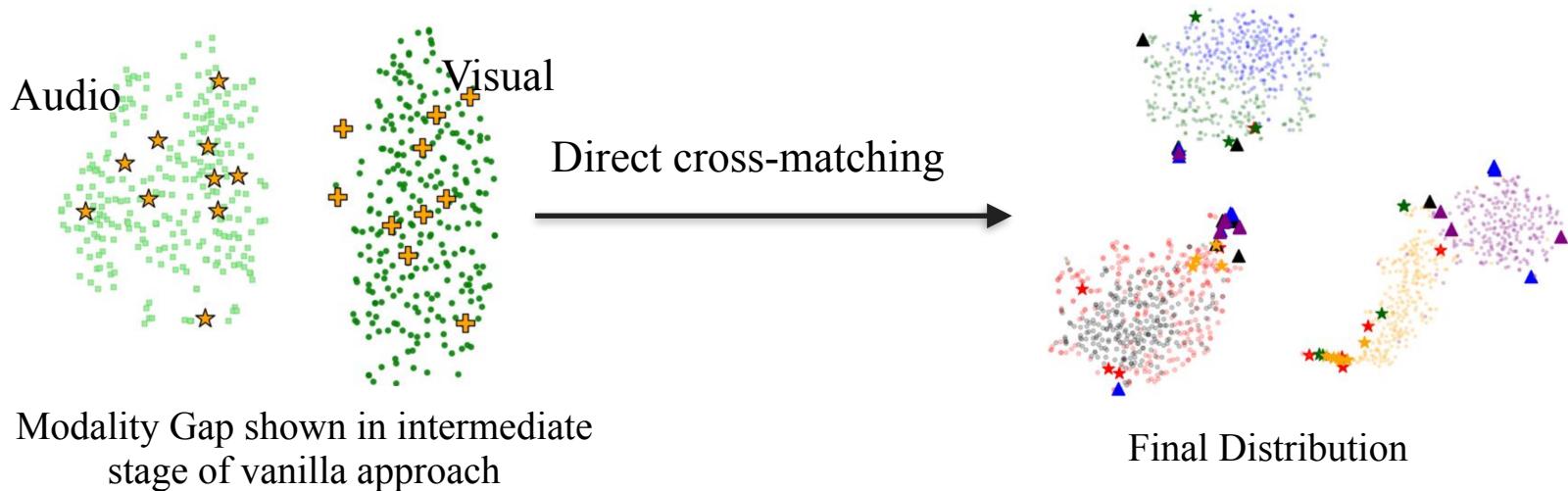
[1] B. Zhao “Data Condensation using distribution matching” WACV 2023

[2] G. Zhao “Improved Distribution Matching for Dataset condensation” CVPR 2023

[3] H. Zhang “M3D: Dataset condensation by minimizing maximum mean discrepancy” AAAI 2024

# Distilling cross-modal alignment.

- Vanilla approach separately distills audio-only and visual-only data
- Simple/direct cross-matching results in unstable learning
  - Modality Gap due to random feature extractors in DM method





# Distilling cross-modal alignment.

- Vanilla approach separately distills audio-only and visual-only data
- Simple/direct cross-matching results in unstable learning
  - Modality Gap due to random feature extractors in DM method
- Novel cross-matching losses
  - Joint Matching loss
  - Modality Gap Matching loss



# Distilling cross-modal alignment.

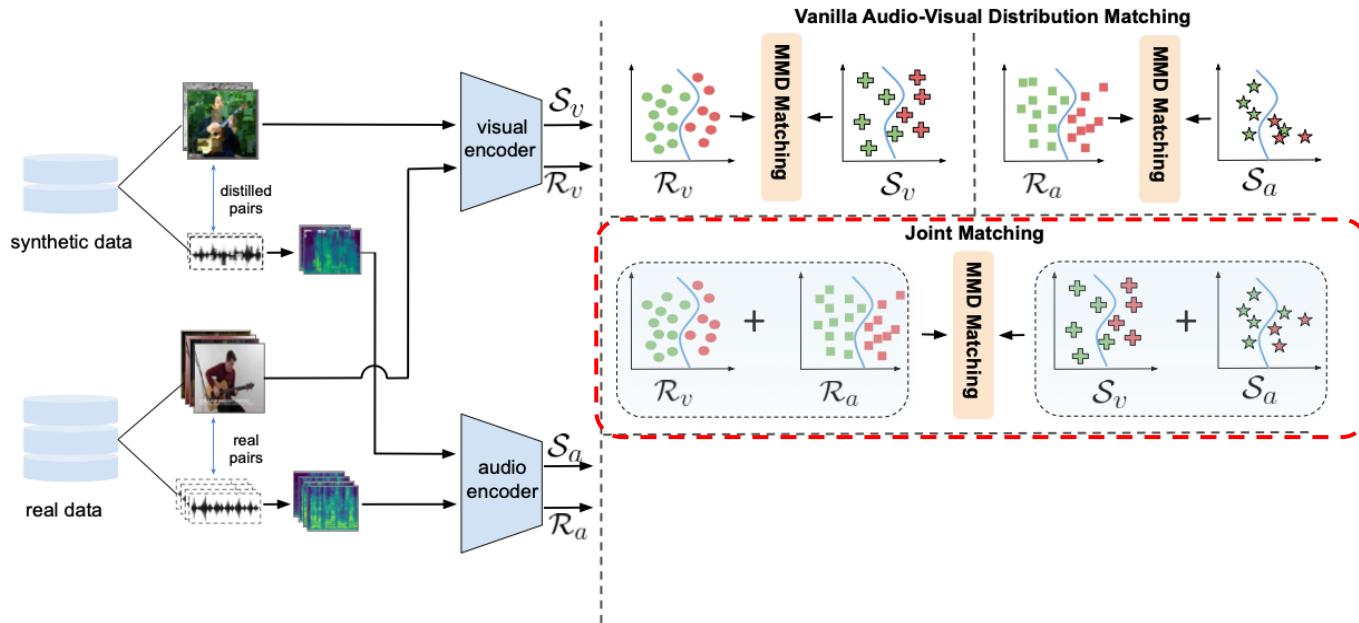
- Vanilla approach separately distills audio-only and visual-only data
- Simple/direct cross-matching results in unstable learning
  - Modality Gap due to random feature extractors in DM method
- Novel cross-matching losses
  - Joint Matching loss
  - Modality Gap Matching loss
- Additional improvements
  - Herding-based initialisation<sup>[1]</sup>
  - Factor technique<sup>[2]</sup>

[1] Max Welling, "Herding dynamical weights to learn", ICML 2009.

[2] G. Zhao et. Al., "Improved distribution matching", CVPR 2023.

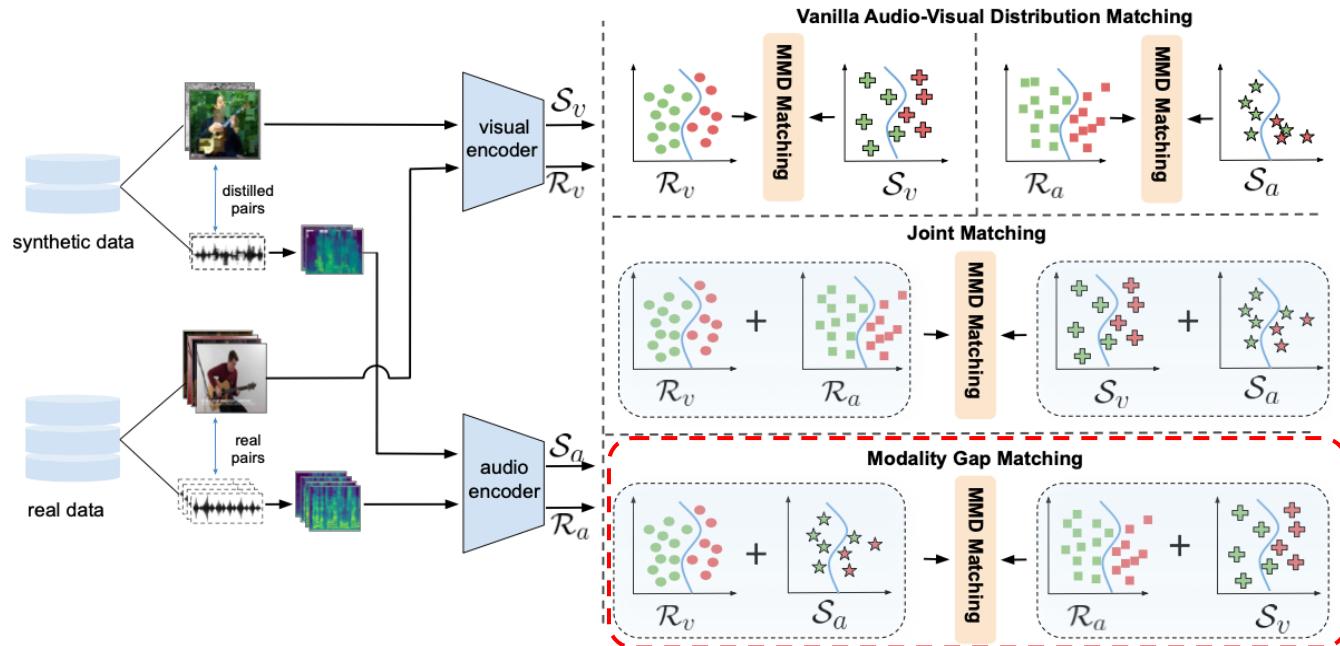
# Joint Matching Loss

- Implicitly distills cross-modal alignment



# Modality Gap Matching Loss

- Implicitly aligns the audio-visual gap between real and synthetic data



# Audio-visual event recognition

- Comparison with other baseline approaches
- Our approach performs better at different configurations
- Benefit of cross-modal alignment

	IPC	Ratio%	Coreset Selection		Training Set Synthesis		<b>Ours</b>	Whole data
			Random	Herding [7]	MTT [1]	DM [8]		
VGGS-10K	1	0.11	$15.44 \pm 1.87$	$20.77 \pm 2.11$	$34.13 \pm 3.62$	$36.54 \pm 2.52$	<b><math>40.41 \pm 1.81</math></b>	$68.24 \pm 0.75$
	10	1.13	$32.01 \pm 1.64$	$39.89 \pm 1.64$	$36.79 \pm 1.97$	$43.85 \pm 1.75$	<b><math>54.99 \pm 1.73</math></b>	
	20	2.26	$45.1 \pm 2.31$	$50.2 \pm 0.74$	$51.87 \pm 1.26$	$49.01 \pm 2.44$	<b><math>58.04 \pm 1.68</math></b>	
AVE	1	0.10	$10.07 \pm 1.16$	$11.84 \pm 0.4$	$12.13 \pm 0.41$	$21.70 \pm 1.46$	<b><math>23.00 \pm 1.37</math></b>	$52.20 \pm 0.38$
	10	1.0	$20.0 \pm 1.45$	$26.86 \pm 0.52$	$23.15 \pm 0.95$	$28.14 \pm 1.80$	<b><math>36.82 \pm 0.88</math></b>	
	20	2.0	$26.32 \pm 1.01$	$33.04 \pm 0.38$	-	$32.57 \pm 0.97$	<b><math>40.13 \pm 1.00</math></b>	

# Further experiments

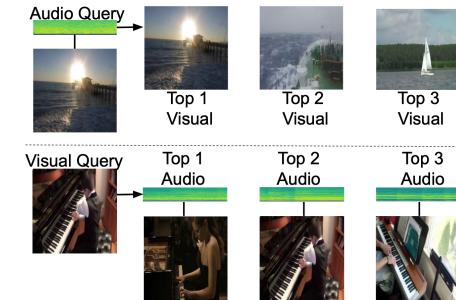
- Ablation study

Random	Herding	Factor	Base	JM	MGM	VGGS-10k	AVE
✓						$32.01 \pm 1.64$	$20.00 \pm 1.45$
	✓					$39.89 \pm 1.64$	$26.86 \pm 0.52$
✓		✓				$40.28 \pm 2.34$	$31.80 \pm 1.28$
✓	✓		✓			$45.31 \pm 2.68$	$34.80 \pm 1.68$
✓	✓		✓	✓		$49.07 \pm 1.97$	$35.13 \pm 1.14$
✓	✓		✓	✓	✓	<b><math>54.99 \pm 1.73</math></b>	<b><math>36.82 \pm 0.88</math></b>

# Further experiments

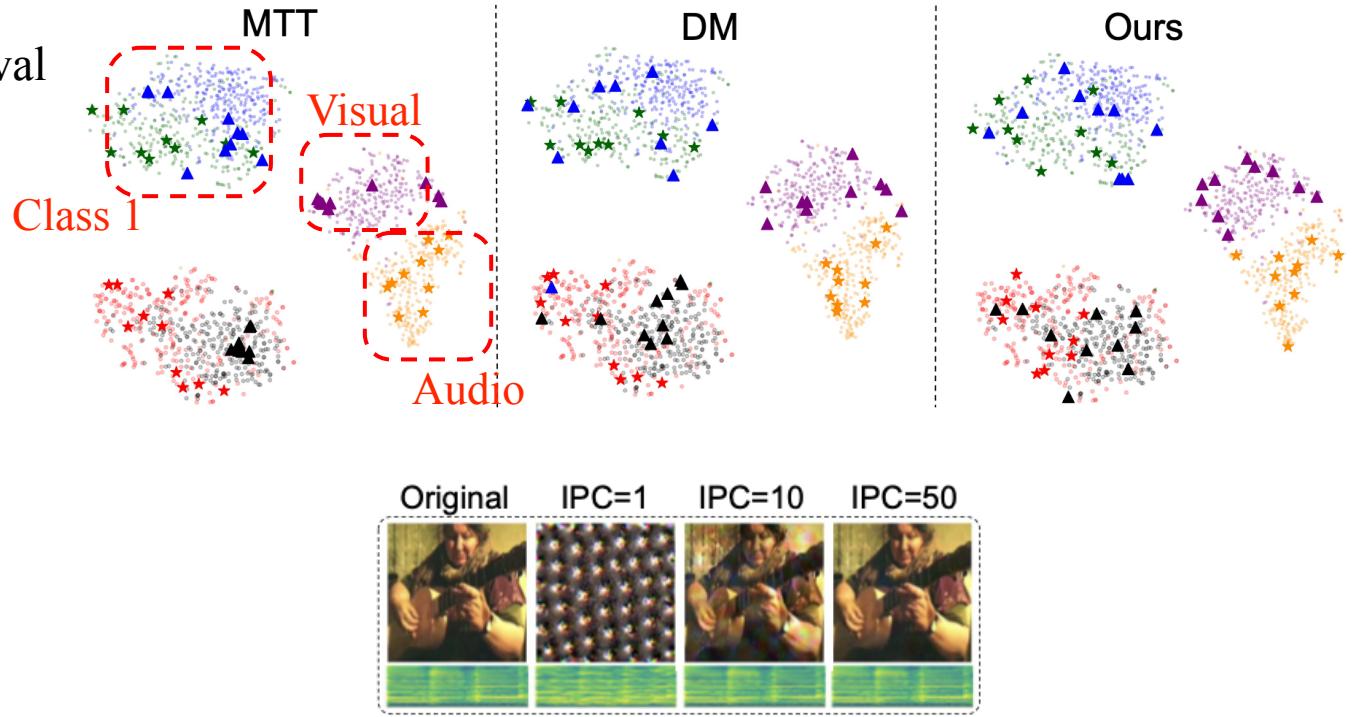
- Ablation study
- Cross-modal retrieval

Method	VGGS-10k test subset			AVE test subset		
	R@1↑	R@5↑	MedR↓	R@1↑	R@5↑	MedR↓
A→V	Random	13.33±5.03	52.00±14.00	5.83±1.75	7.62±3.21	30.23±4.06
	DM[8]	8.66±1.15	47.33±5.77	6.66±1.52	6.90±2.29	32.14±0.71
	<b>Ours</b>	<b>19.33±2.30</b>	<b>59.33±1.15</b>	<b>3.66±0.57</b>	<b>13.09±2.88</b>	<b>35.00±1.88</b>
V→A	Whole data	44.00±2.00	74.00±5.03	2.00±0.00	27.61±5.35	51.66±4.06
	Random	10.66±2.30	49.33±5.77	6.00±0.86	9.04±1.48	26.66±2.29
	DM[8]	11.33±3.05	44.00±4.00	6.66±1.15	<b>10.95±3.59</b>	29.52±3.52
	<b>Ours</b>	<b>27.33±2.30</b>	<b>59.33±7.02</b>	<b>3.83±0.57</b>	<b>6.43±3.11</b>	<b>34.52±3.30</b>
Whole data	45.33±5.03	76.00±2.00	1.83±0.28	17.14±0.71	44.76±1.79	7.16±0.288



# Further experiments

- Ablation study
- Cross-modal retrieval
- Visualisations
  - Distribution
  - Distilled data





# Conclusion

- New problem: Audio-Visual dataset distillation
- Audio-visual integration still hold for synthetic data
- Need carefully designed cross-modal alignment losses
- Extensive experiments on Audio-visual recognition and retrieval tasks
- Future Work and Limitation:
  - Extend to longer videos
  - Reduce gap with whole data
  - Extend to Instance-wise distillation