

Multi-Modal Learning

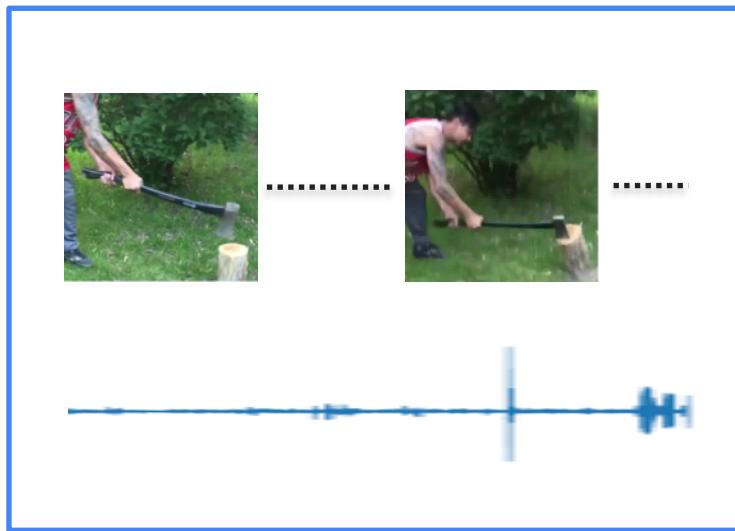
Saksham Singh Kushwaha

The University of Texas at Dallas

(DL4MIR'24 Workshop)



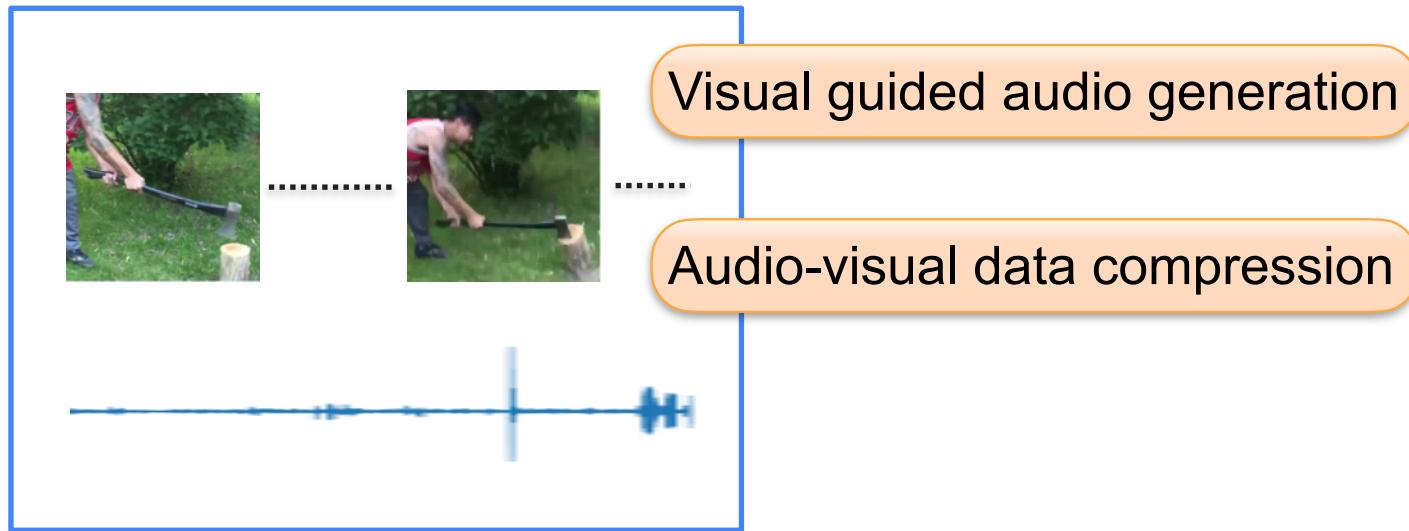
Research Interests



Audio-Visual Integration

[1] Chen et. al, "Vggsound: A Large-Scale Audio-Visual Dataset", ICASSP, 2020.

Research Interests

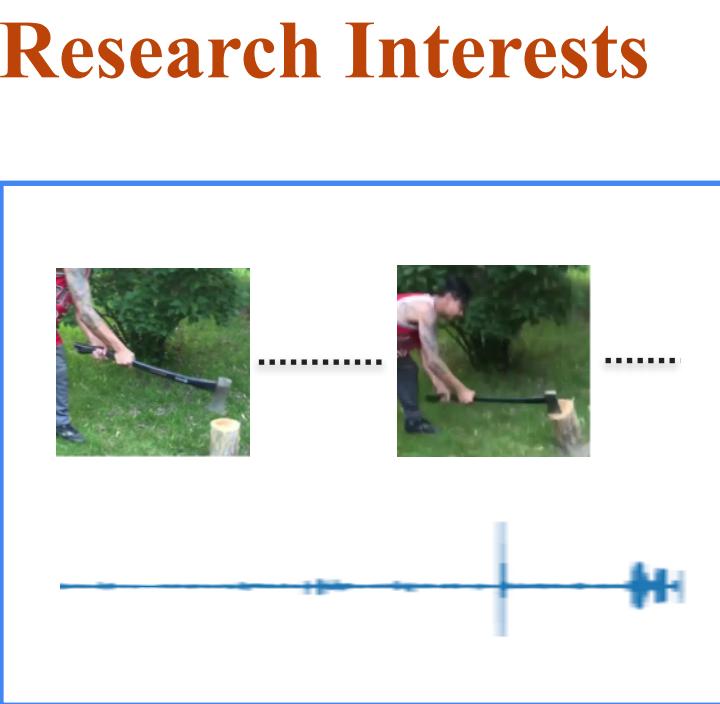


[1] Chen et. al, "Vggsound: A Large-Scale Audio-Visual Dataset", ICASSP, 2020.

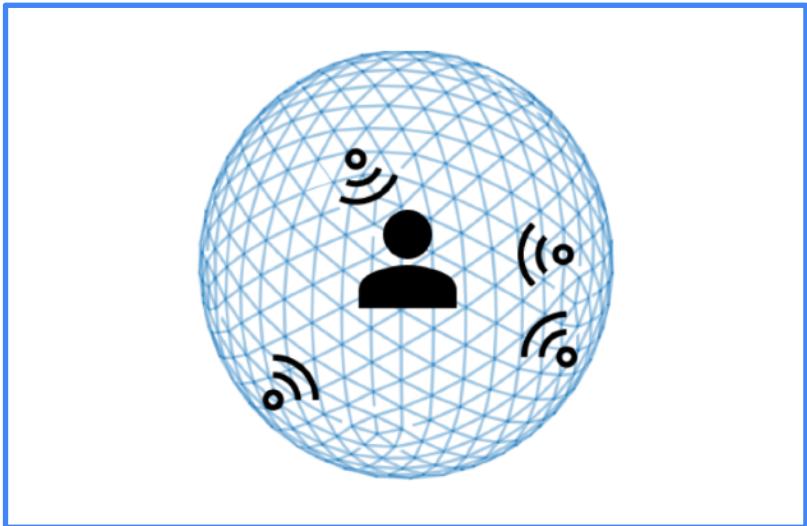
Research Interests



Audio-Visual Integration



Spatial Audio

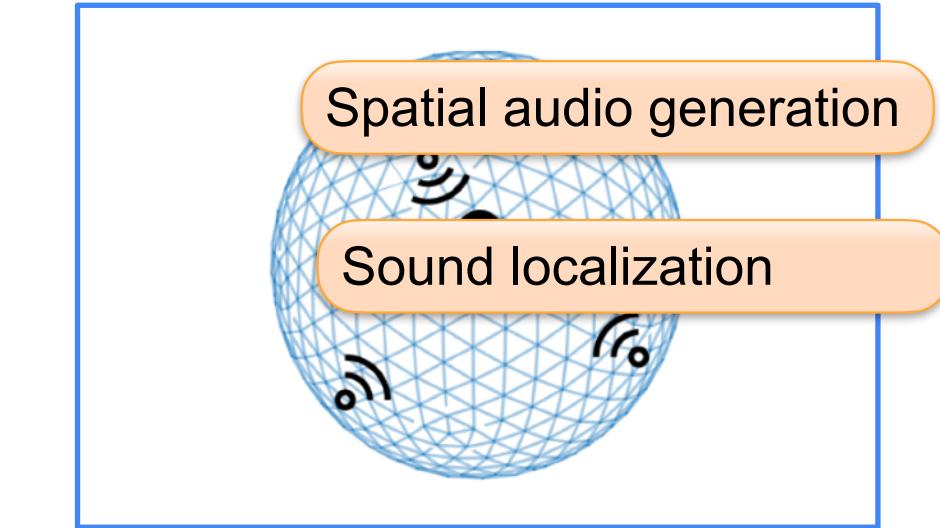


- [1] Chen et. al, "Vggsound: A Large-Scale Audio-Visual Dataset", ICASSP, 2020.
- [2] Dagli et. Al, "SEE-2-SOUND: Zero-Shot Spatial Environment-to-Spatial Sound", ArXiv, 2024

Research Interests



Audio-Visual Integration



Spatial Audio

- [1] Chen et. al, "Vggsound: A Large-Scale Audio-Visual Dataset", ICASSP, 2020.
[2] Dagli et. Al, "SEE-2-SOUND: Zero-Shot Spatial Environment-to-Spatial Sound", ArXiv, 2024

Dataset Distillation for Audio-Visual Datasets

(CVPR Sight and Sound Workshop, 2024)

Saksham Singh Kushwaha¹, Siva Sai Nagender Vasireddy¹, Kai Wang²,Yapeng Tian¹

¹University of Texas at Dallas, USA, ²National University of Singapore, Singapore

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



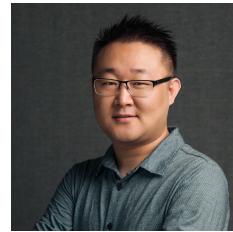
Saksham Singh Kushwaha
Ph.D. Candidate



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

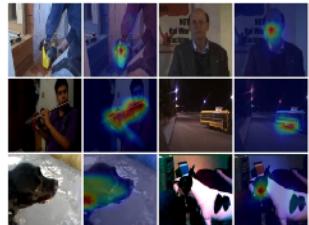


Kai Wang
Ph.D. Candidate



Yapeng Tian
Assistant Professor

Increase in dataset sizes



AVE (2018)
(~4k)

Pandas-70M (2024)

 **Panda-70M**

COCO (2014)
(~120k)

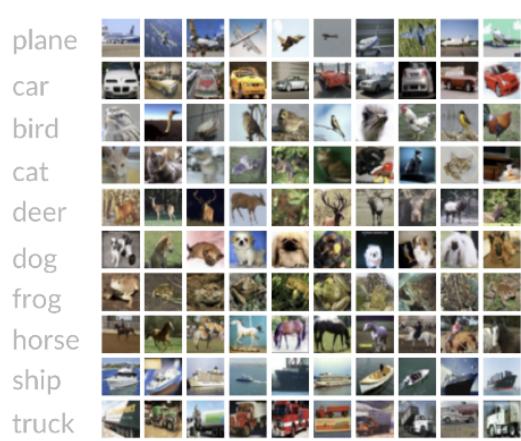


VGGSound (2020)
(~200k)

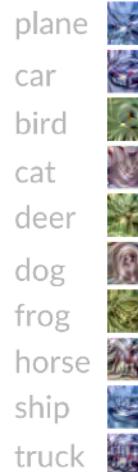


- [1] Chen et. al, "VggSound: A Large-Scale Audio-Visual Dataset", ICASSP, 2020.
- [2] Tian et. al, "Audio-visual event localization in unconstrained videos", ECCV, 2018
- [3] Yin et. al, "Microsoft **coco**: Common objects in context", ECCV, 2014
- [4] Chen et. al, "**Panda-70m**: Captioning **70m** videos with multiple cross-modality teachers", CVPR, 2024

Goal: Dataset distillation/compression



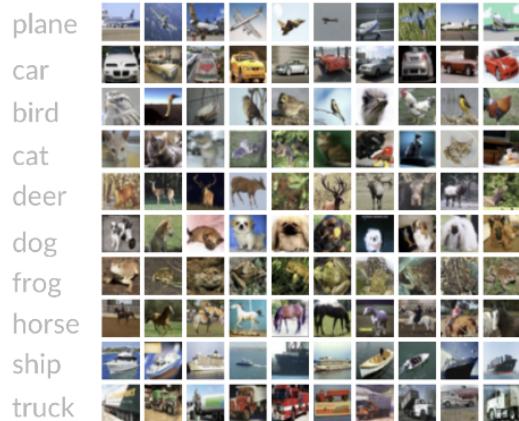
Large training data



Representative data

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

Goal: Dataset distillation/compression



Large training data

Reduce storage ✓

lane
car

Faster hyper-parameter search ✓

deer

Reduce training time ✓

truck

Effective Continual learning ✓

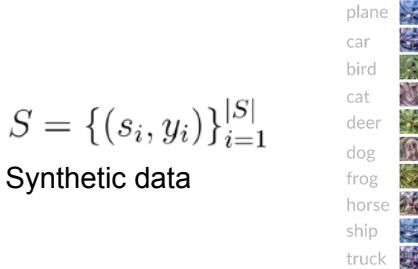
truck

Representative data

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

Representativeness

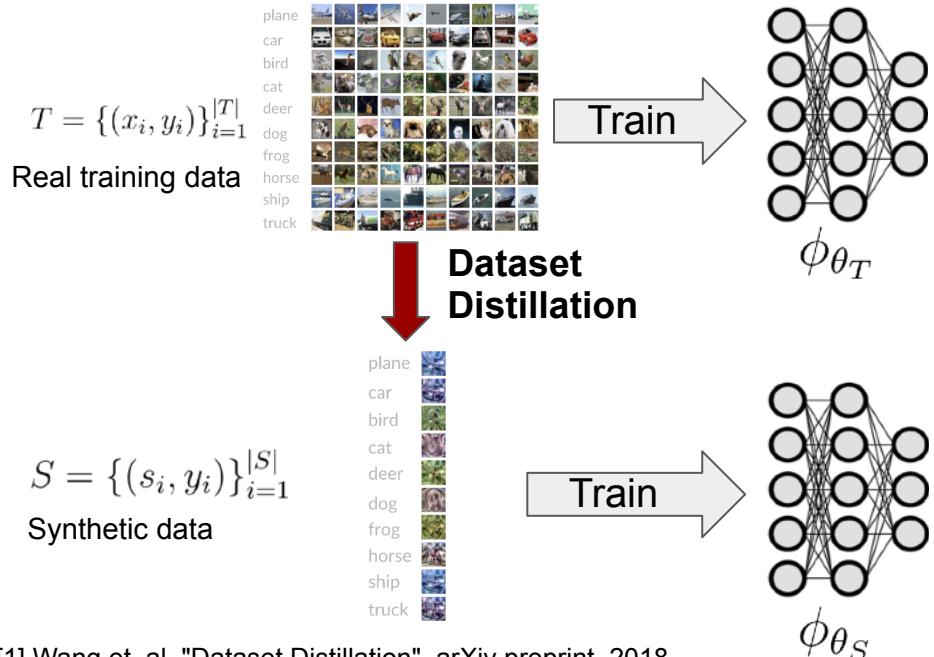
- Comparable performance on unseen real test data



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

Representativeness

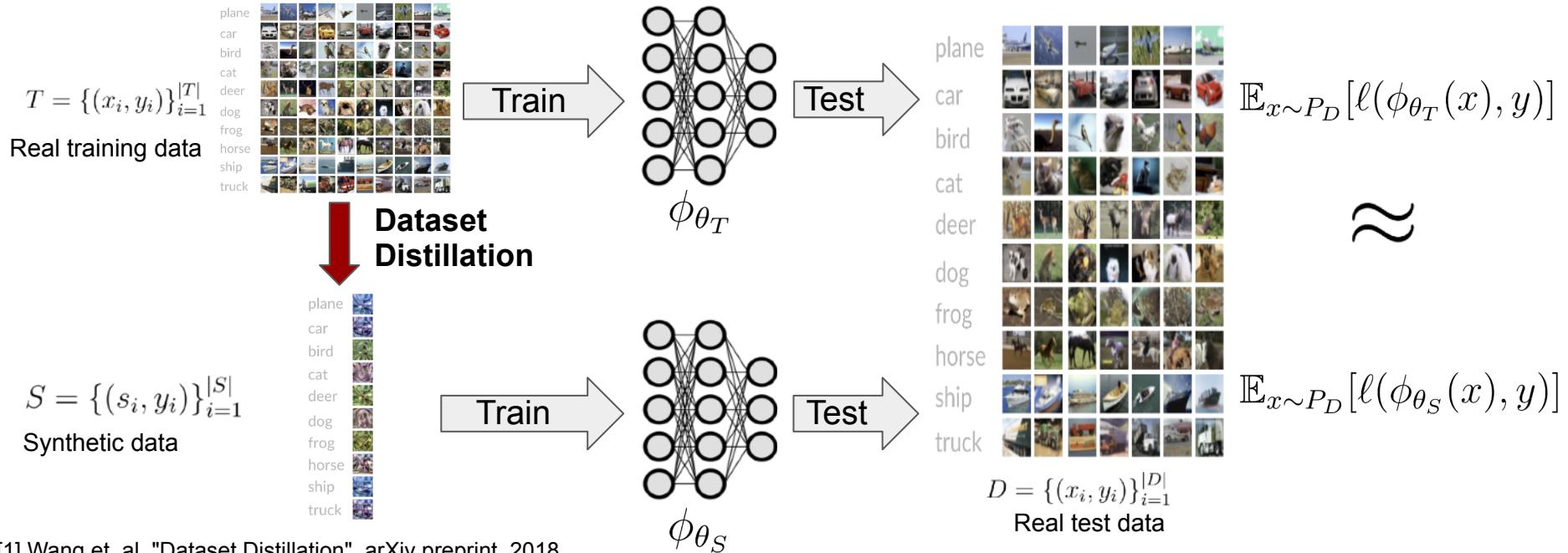
- Comparable performance on unseen real test data



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

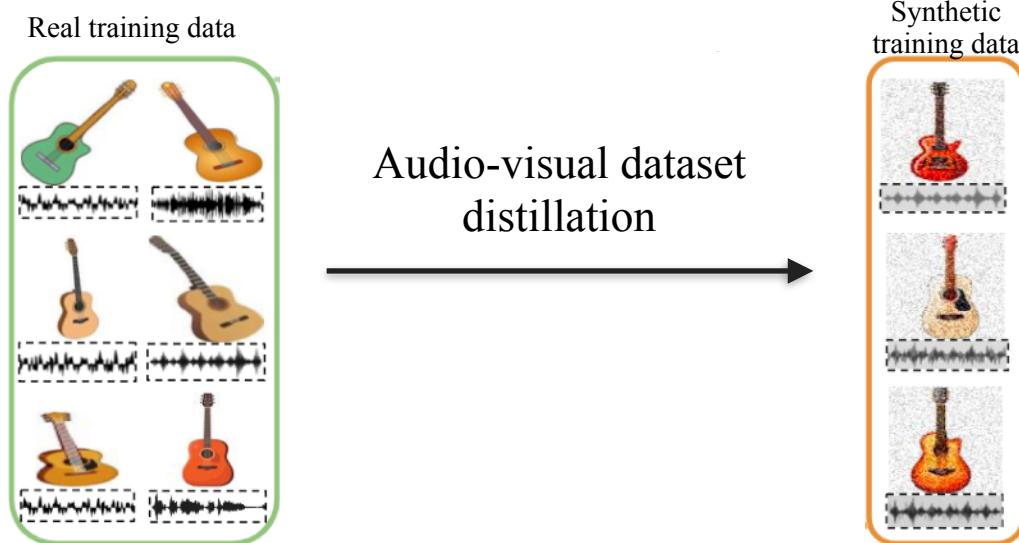
Representativeness

- Comparable performance on unseen real test data



Audio-Visual dataset distillation

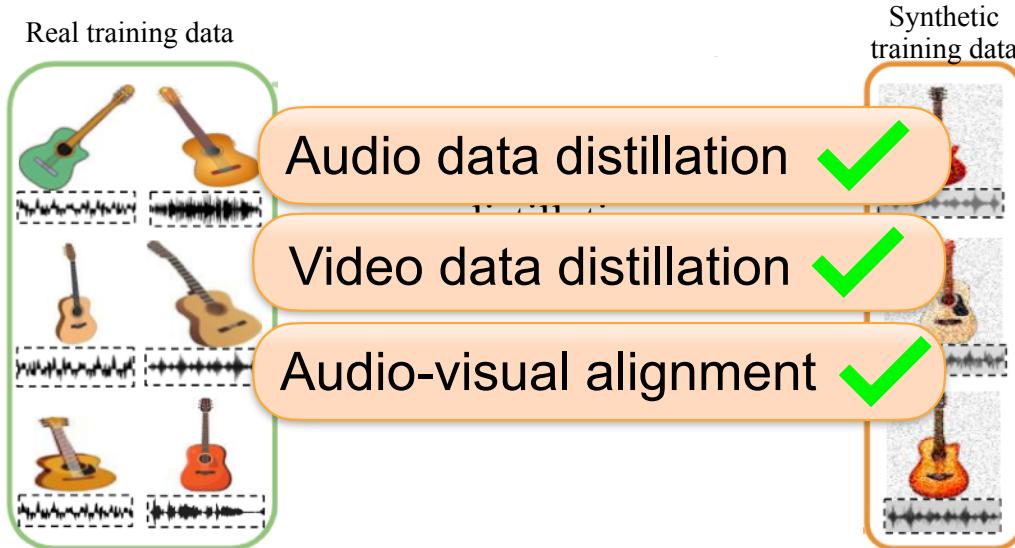
- Previous dataset distillation research focuses on image-only
- Explosion in large multimodal datasets



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

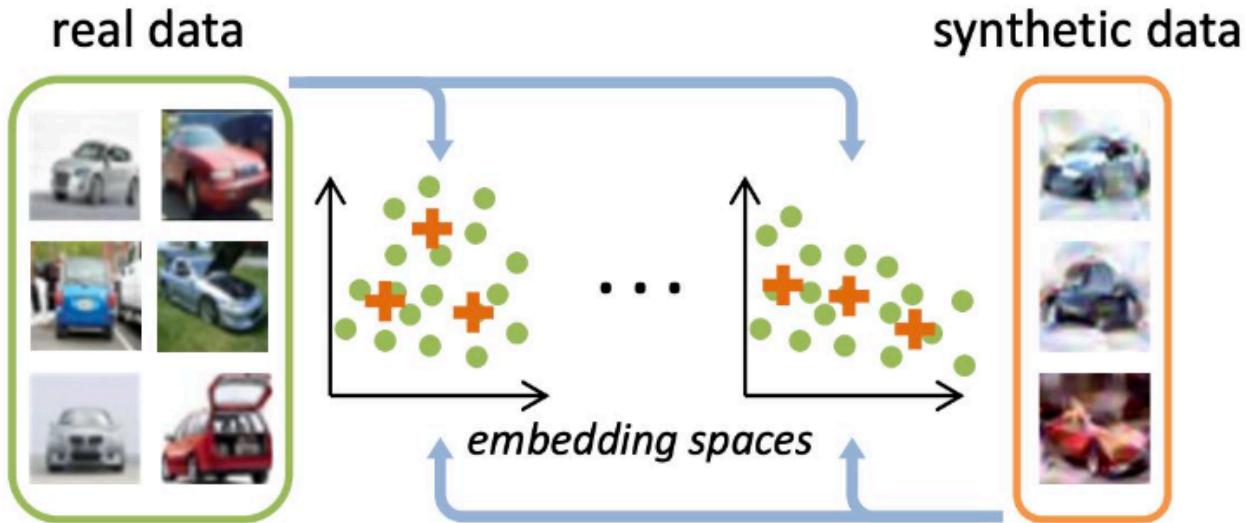
Audio-Visual dataset distillation

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



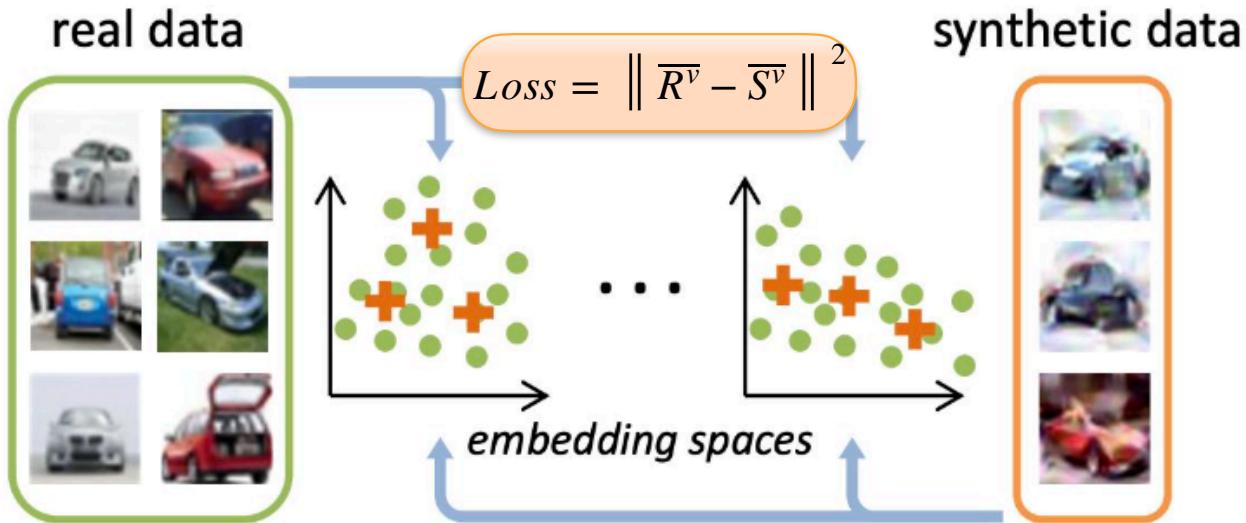
[1] Wang et. al, "Dataset Distillation", arXiv preprint, 2018.

Background: Distribution Matching



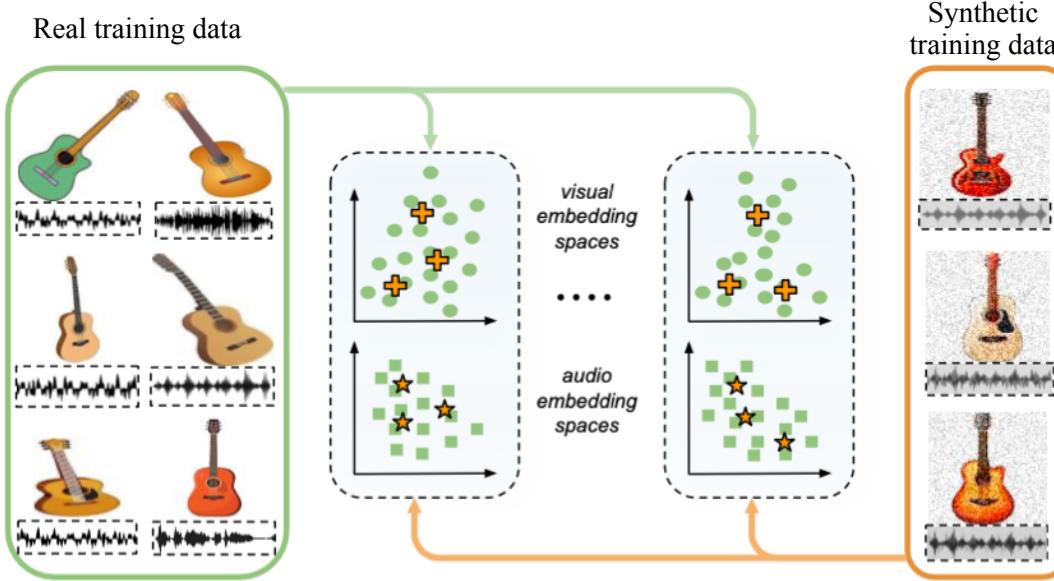
- [1] B. Zhao "Data Condensation using distribution matching" WACV 2023
- [2] G. Zhao "Improved Distribution Matching for Dataset condensation" CVPR 2023
- [3] H. Zhang "M3D: Dataset condensation by minimizing maximum mean discrepancy" AAAI 2024

Background: Distribution Matching



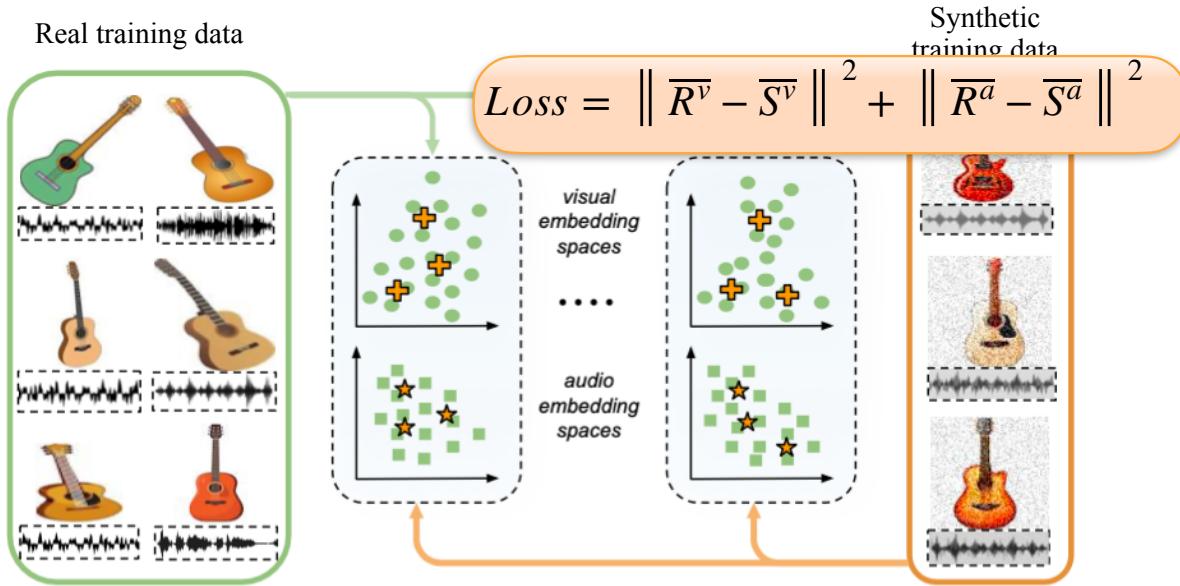
- [1] B. Zhao "Data Condensation using distribution matching" WACV 2023
- [2] G. Zhao "Improved Distribution Matching for Dataset condensation" CVPR 2023
- [3] H. Zhang "M3D: Dataset condensation by minimizing maximum mean discrepancy" AAAI 2024

Vanilla Audio-Visual Dataset 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

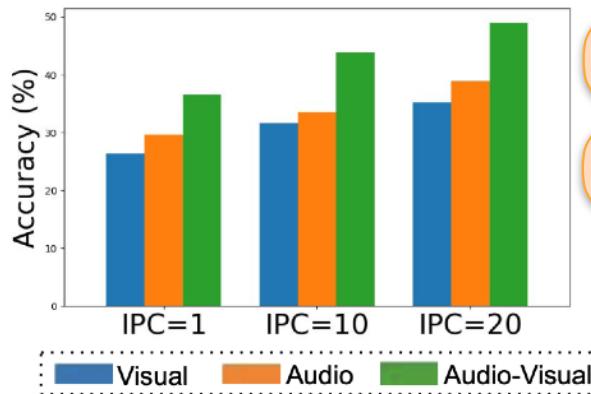
Vanilla Audio-Visual Dataset distillation



- [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

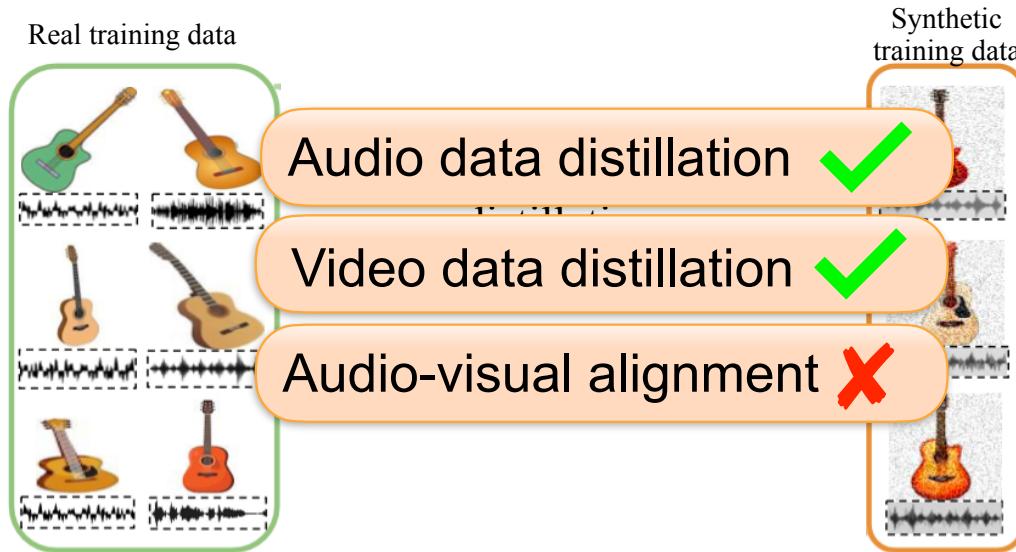
Does integration hold for AV distilled data?

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



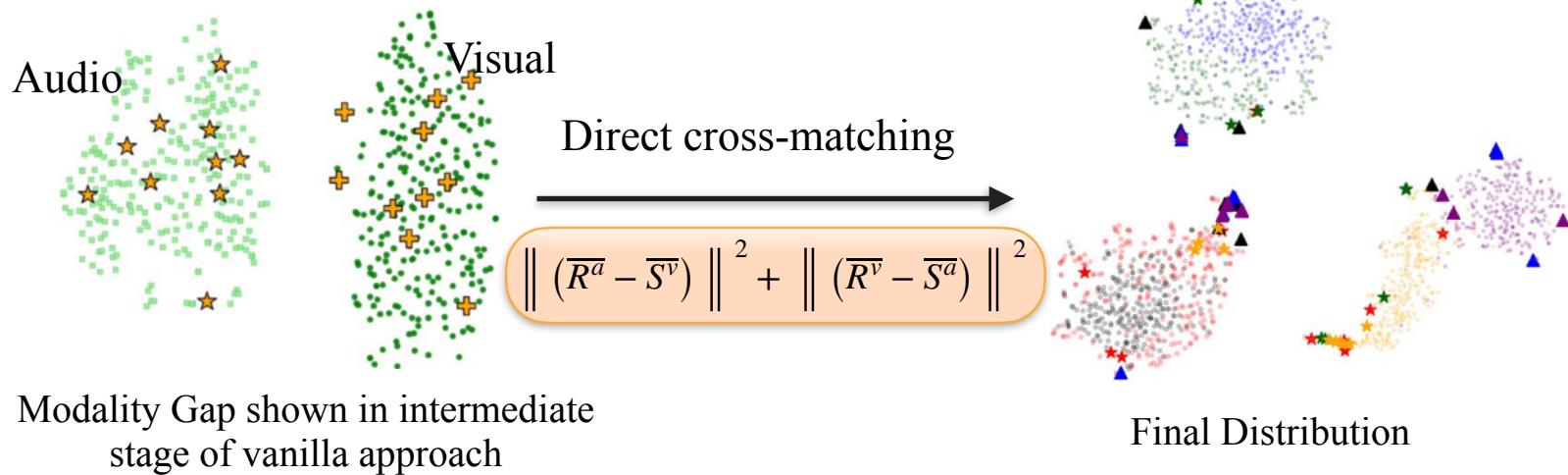
Audio-visual integration helps
Ensemble method outperforms

Vanilla Audio-Visual Dataset distillation



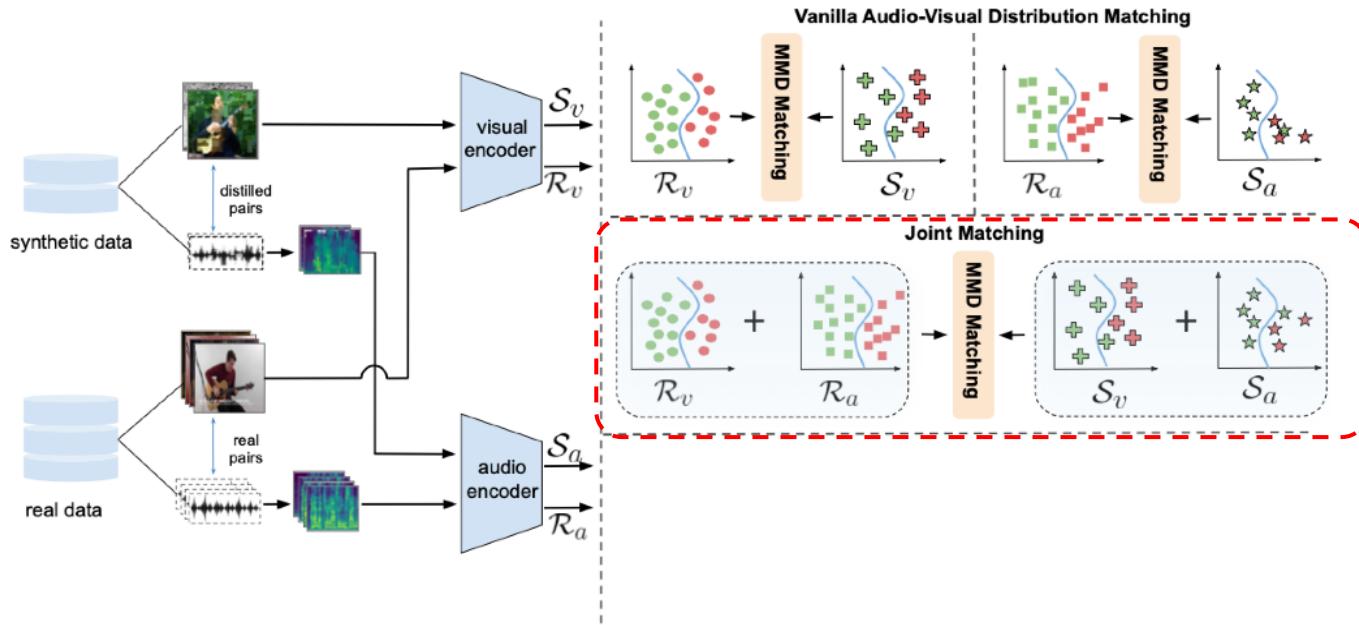
- [1] B. Zhao "Data Condensation using distribution matching" WACV 2023
- [2] G. Zhao "Improved Distribution Matching for Dataset condensation" CVPR 2023
- [3] H. Zhang "M3D: Dataset condensation by minimizing maximum mean discrepancy" AAAI 2024

Simple cross-modal alignment fails



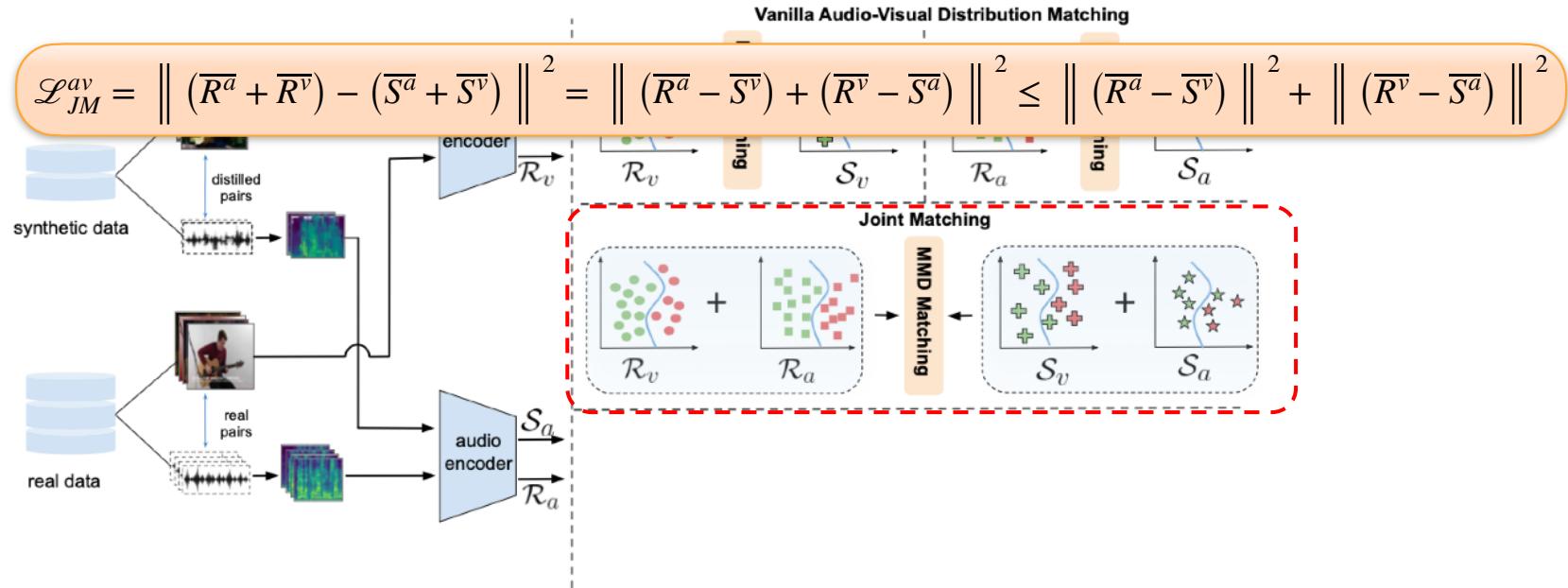
Joint Matching Loss

- Implicitly distills cross-modal alignment



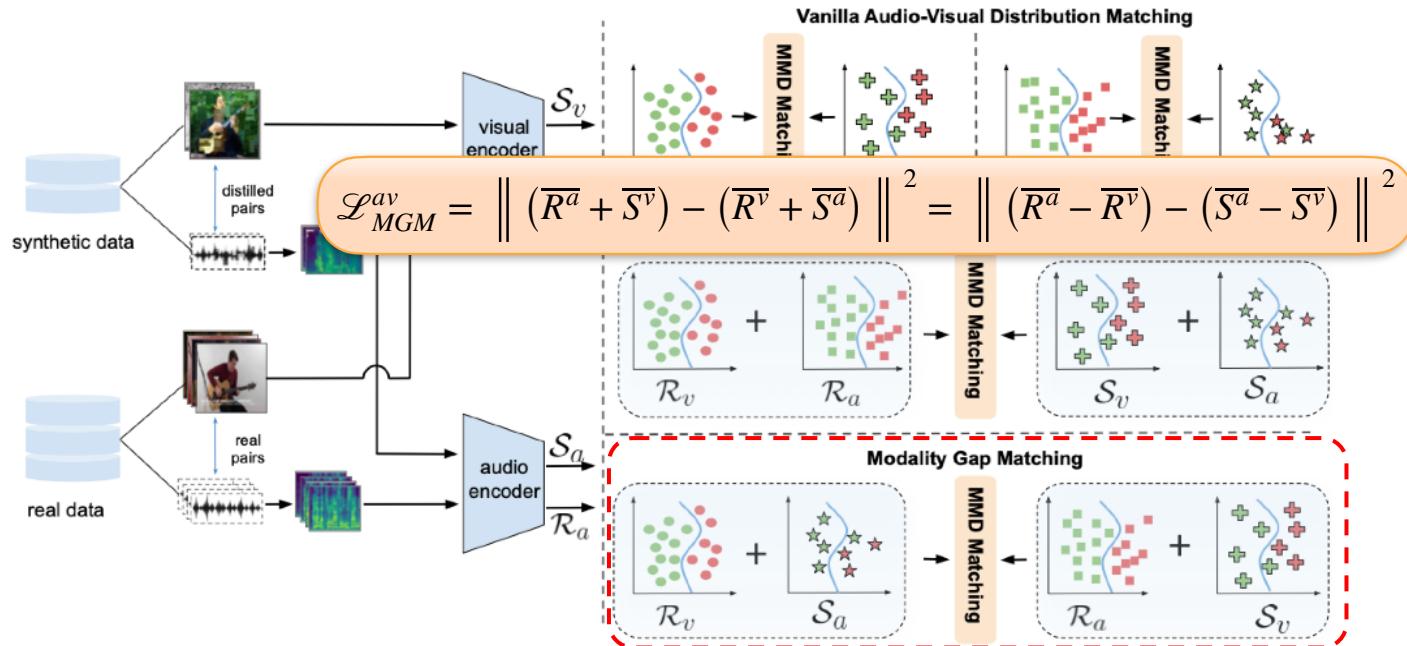
Joint Matching Loss

- Implicitly distills cross-modal alignment

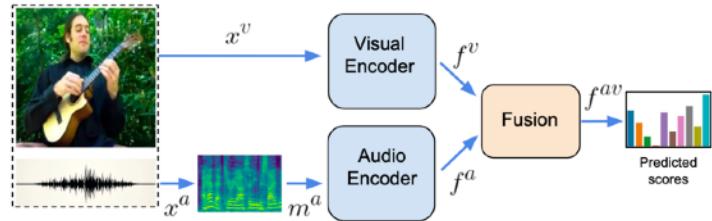


Modality Gap Matching Loss

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

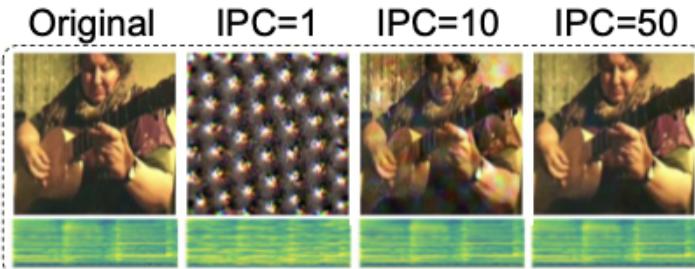
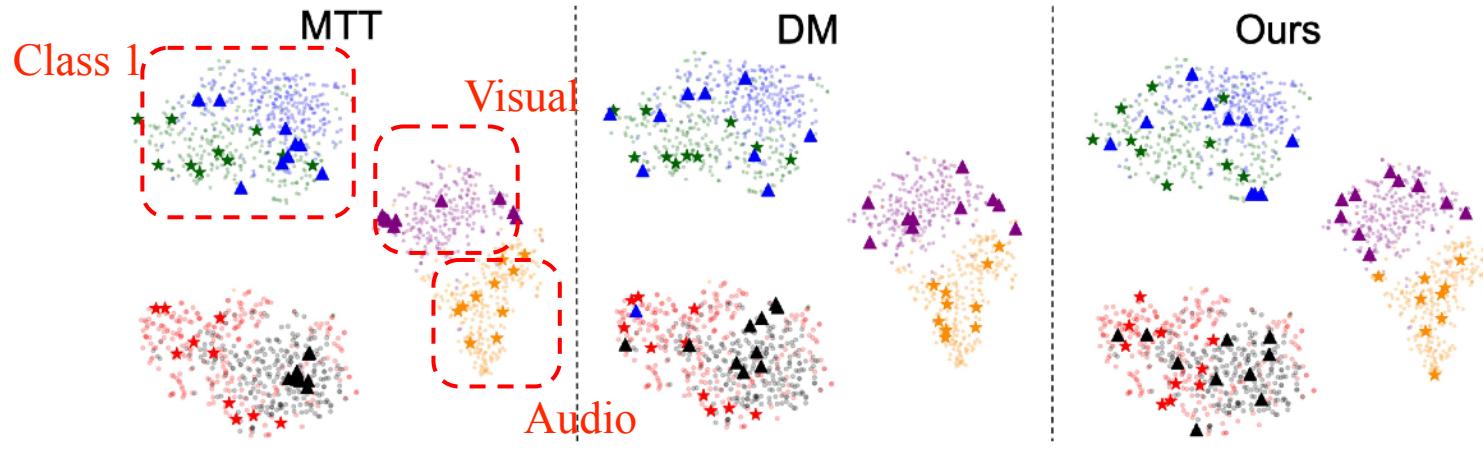


Audio-visual event recognition



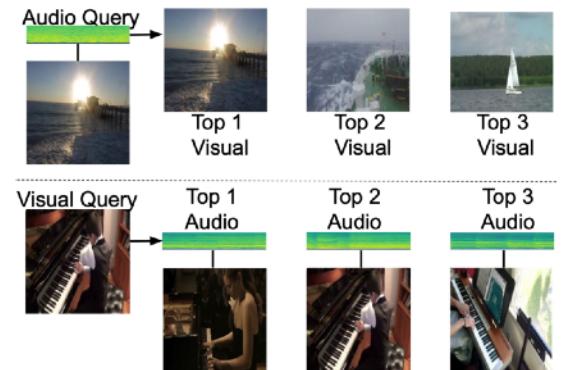
IPC	Ratio%	Coreset Selection		Training Set Synthesis				Ours	Upper Bound (Whole data)
		Random	Herding	DC	DSA	MTT	DM		
VGGS-10K	1	0.11	15.44 ± 1.87	20.77 ± 2.11	18.28 ± 1.36	19.32 ± 1.35	34.13 ± 3.62	36.54 ± 2.52	40.41 ± 1.81
	10	1.13	32.01 ± 1.64	39.89 ± 1.64	32.10 ± 0.84	36.61 ± 1.04	36.79 ± 1.97	43.85 ± 1.75	54.99 ± 1.73
	20	2.26	45.1 ± 2.31	50.2 ± 0.74	-	-	51.87 ± 1.26	49.01 ± 2.44	58.04 ± 1.68
VGGSound	1	0.18	1.38 ± 0.12	2.14 ± 0.17	-	-	1.55 ± 0.15	3.08 ± 0.21	4.97 ± 0.30
	10	1.87	5.55 ± 0.19	7.09 ± 0.09	-	-	-	6.40 ± 0.28	8.23 ± 0.24
	20	3.74	8.00 ± 0.14	9.64 ± 0.14	-	-	-	8.64 ± 0.14	9.85 ± 0.34
MUSIC-21	1	0.014	24.12 ± 4.33	26.15 ± 2.01	22.60 ± 1.13	22.98 ± 1.16	28.71 ± 1.23	38.26 ± 1.32	44.02 ± 2.21
	10	0.14	45.77 ± 1.74	51.89 ± 1.39	-	-	42.25 ± 1.07	54.78 ± 1.39	68.07 ± 0.98
	20	0.28	54.86 ± 1.85	59.98 ± 0.85	-	-	-	61.06 ± 1.31	70.30 ± 0.69
AVE	1	0.10	10.07 ± 1.16	11.84 ± 0.4	10.45 ± 0.39	10.76 ± 0.62	12.13 ± 0.41	21.70 ± 1.46	23.00 ± 1.37
	10	1.0	20.0 ± 1.45	26.86 ± 0.52	-	-	23.15 ± 0.95	28.14 ± 1.80	36.82 ± 0.88
	20	2.0	26.32 ± 1.01	33.04 ± 0.38	-	-	-	32.57 ± 0.97	40.13 ± 1.00

Visualization



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 \pm 5.03	52.00 \pm 14.00	5.83 \pm 1.75	7.62 \pm 3.21	30.23 \pm 4.06
	DM	8.66 \pm 1.15	47.33 \pm 5.77	6.66 \pm 1.52	6.90 \pm 2.29	32.14 \pm 0.71
	Ours	19.33 \pm 2.30	59.33 \pm 1.15	3.66 \pm 0.57	13.09 \pm 2.88	35.00 \pm 1.88
	Whole data	44.00 \pm 2.00	74.00 \pm 5.03	2.00 \pm 0.00	27.61 \pm 5.35	51.66 \pm 4.06
V→A	Random	10.66 \pm 2.30	49.33 \pm 5.77	6.00 \pm 0.86	9.04 \pm 1.48	26.66 \pm 2.29
	DM	11.33 \pm 3.05	44.00 \pm 4.00	6.66 \pm 1.15	10.95 \pm 3.59	29.52 \pm 3.52
	Ours	27.33 \pm 2.30	59.33 \pm 7.02	3.83 \pm 0.57	6.43 \pm 3.11	34.52 \pm 3.30
	Whole data	45.33 \pm 5.03	76.00 \pm 2.00	1.83 \pm 0.28	17.14 \pm 0.71	44.76 \pm 1.79





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



Multimodal Prototypical approach for Unsupervised sound classification

INTERSPEECH 2023

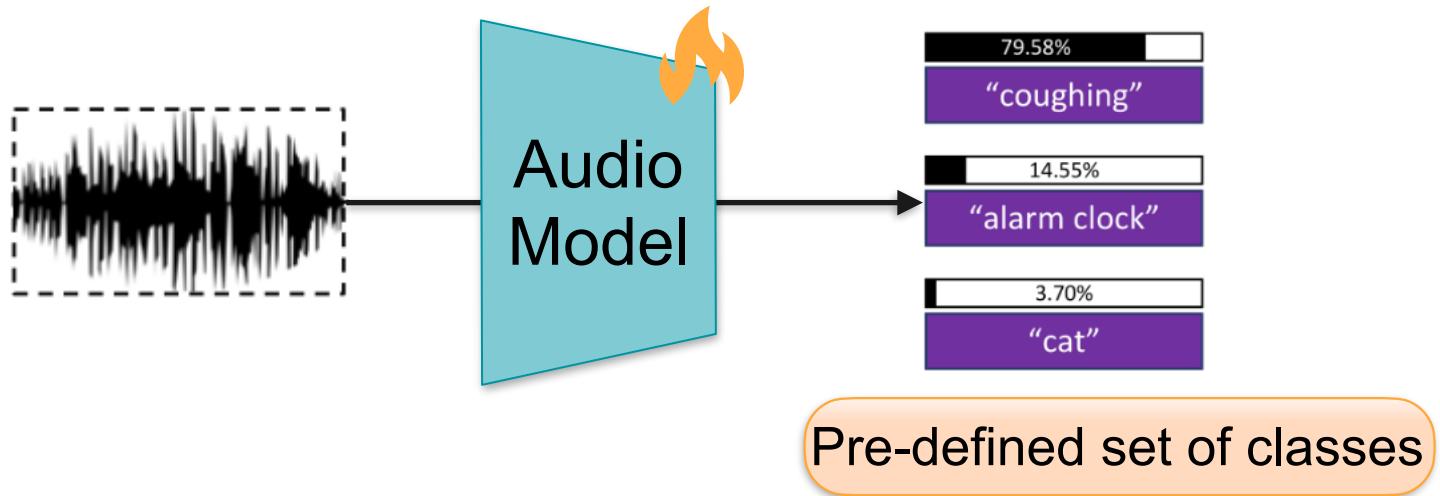


**Saksham Singh
Kushwaha**
MSCS, NYU



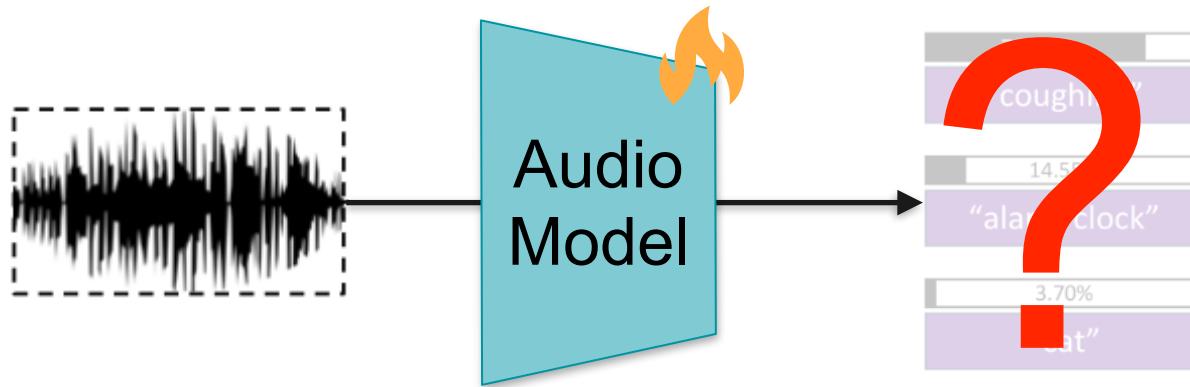
Magdalena Fuentues
Assistant Prof., NYU Tandon

Audio classification



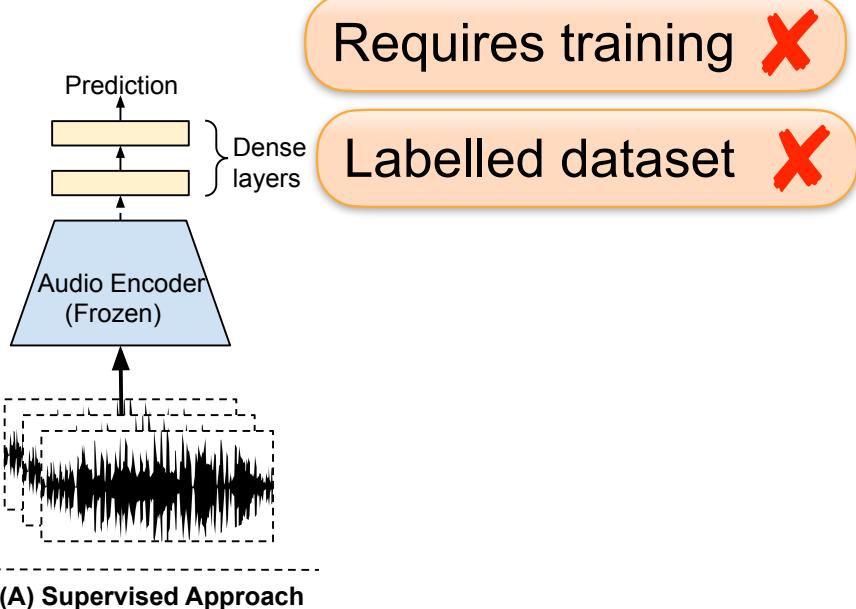


Audio classification



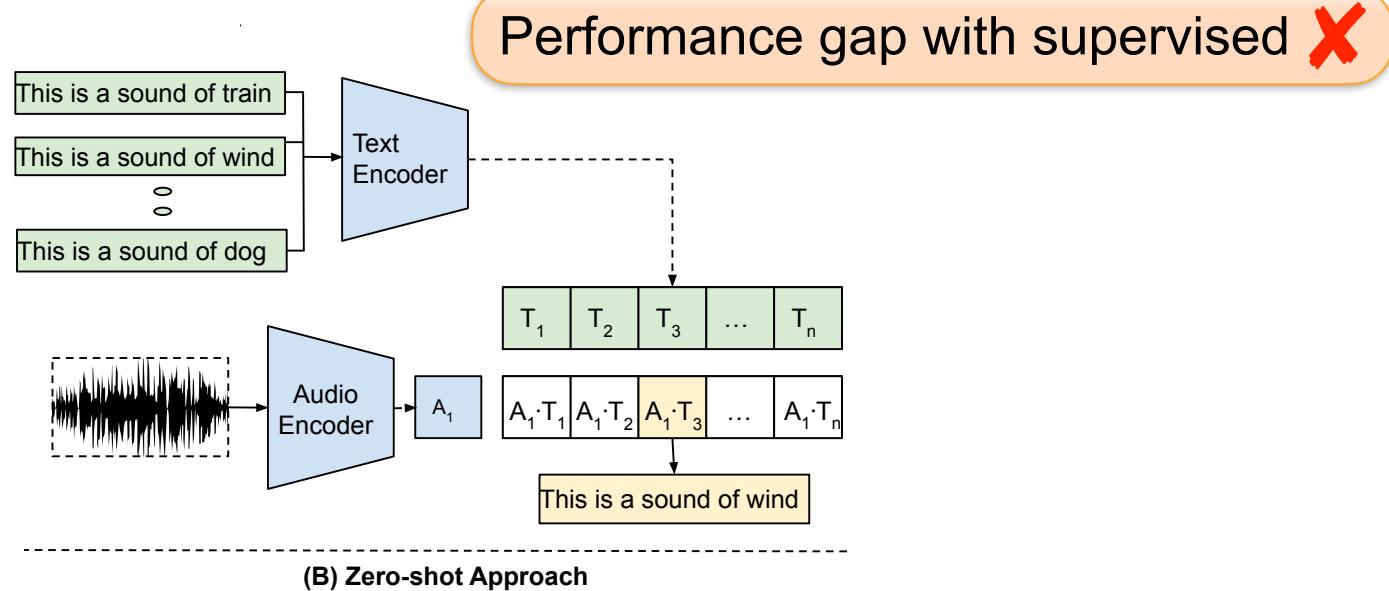
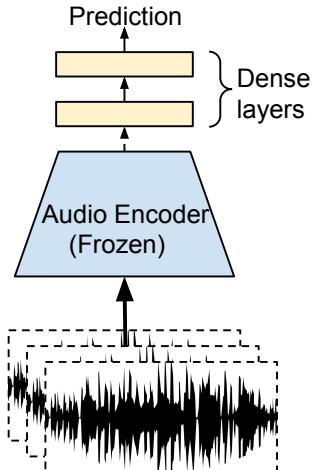


Previous approaches





Previous approaches



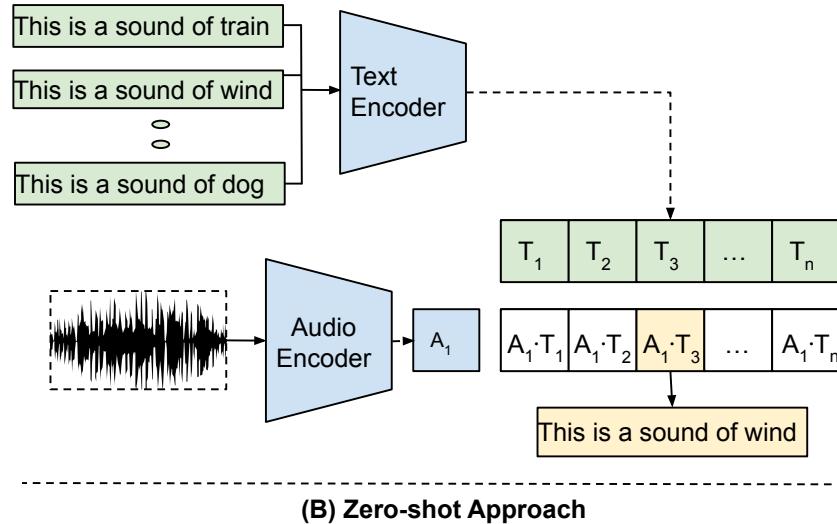
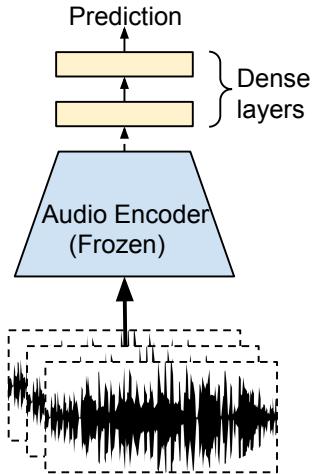
(A) Supervised Approach

(B) Zero-shot Approach





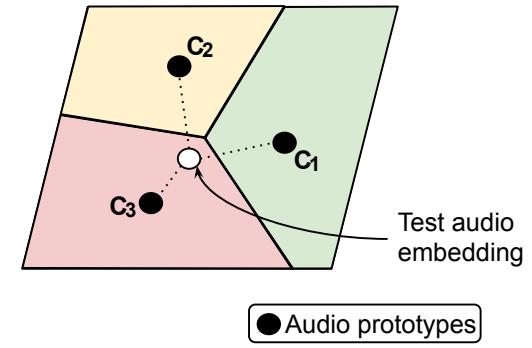
Previous approaches



(A) Supervised Approach

(B) Zero-shot Approach

Labelled data X

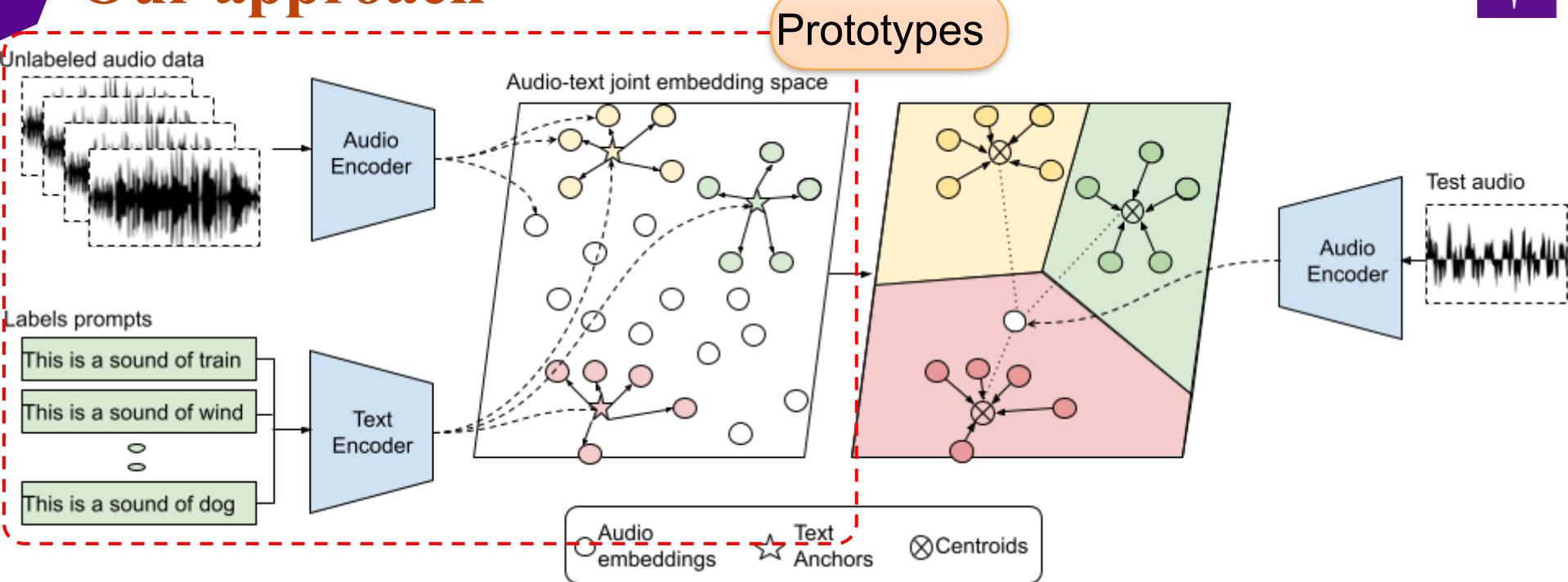


(C) Prototypical Approach



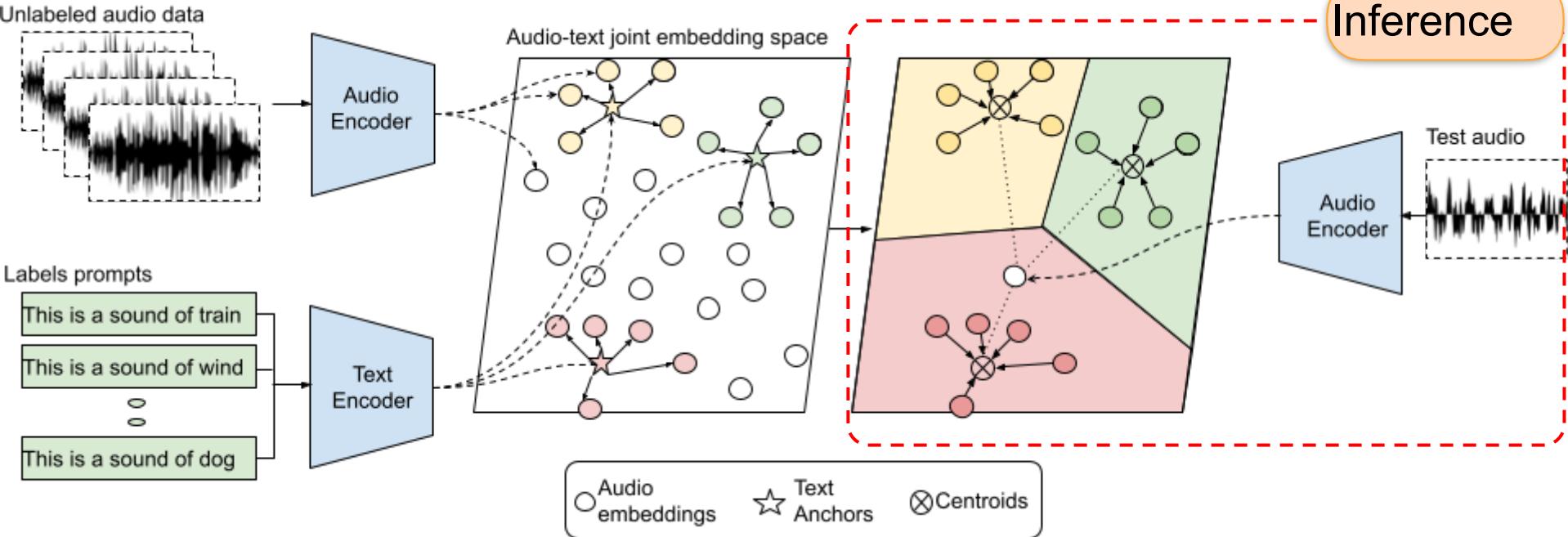


Our approach





Our approach





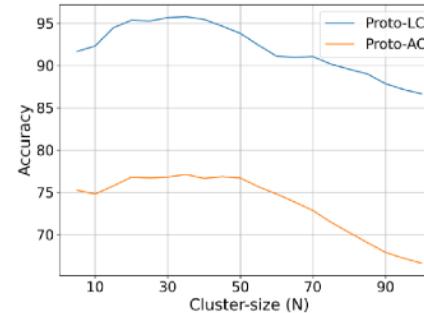
Results

	ESC-50 (acc)		US8k (acc)		FSD50K (mAP)	
	Zero Shot	Supervised	Zero Shot	Supervised	Zero Shot	Supervised
Wav2Clip	0.41	0.86	0.40	0.81	0.03	0.43
AudioClip	0.68	0.88	0.62	0.86	0.20	0.50
CLAP	0.83	0.97	0.73	0.88	0.30	0.59
LAION-CLAP	0.91	0.96	0.72	0.89	0.22	0.61
Proto-AC(Ours)	0.78	0.82	0.71	0.77	0.40	0.48
Proto-LC(Ours)	0.96	0.97	0.73	0.83	0.52	0.65



Ablation Studies

- Cluster-size vs. accuracy
- Prompt template vs. accuracy



	AudioClip	LAION-CLAP	Proto-AC	Proto-LC
Prompt 1	0.67	0.83	0.77	0.94
Prompt 2	0.68	0.86	0.72	0.92
Prompt 3	0.69	0.90	0.72	0.96
Prompt 4	0.68	0.88	0.78	0.95
Prompt 5	0.67	0.92	0.70	0.96

Table 2: Accuracy on ESC-50 with different prompts. Prompt 1: '{Class label}', Prompt 2: 'I can hear {class label}', Prompt 3: 'This is an audio of {class label}', Prompt 4: 'This is {class label}', Prompt 5: 'This is a sound of {class label}'.





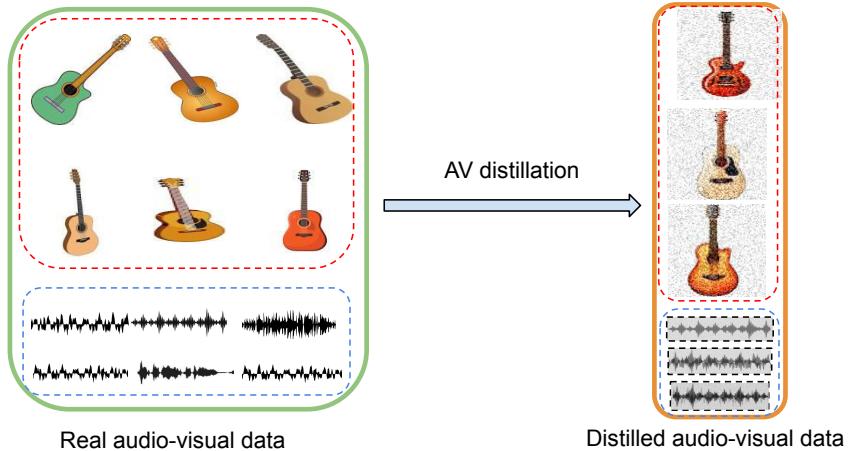
Conclusion

- Personalised audio-classification without human intervention
- Outperform previous zero-shot approach
- Test on different pretrained models and single/multi-label datasets



Future research directions

- Multimodal data compression/distillation
 - Current: Class-wise distillation



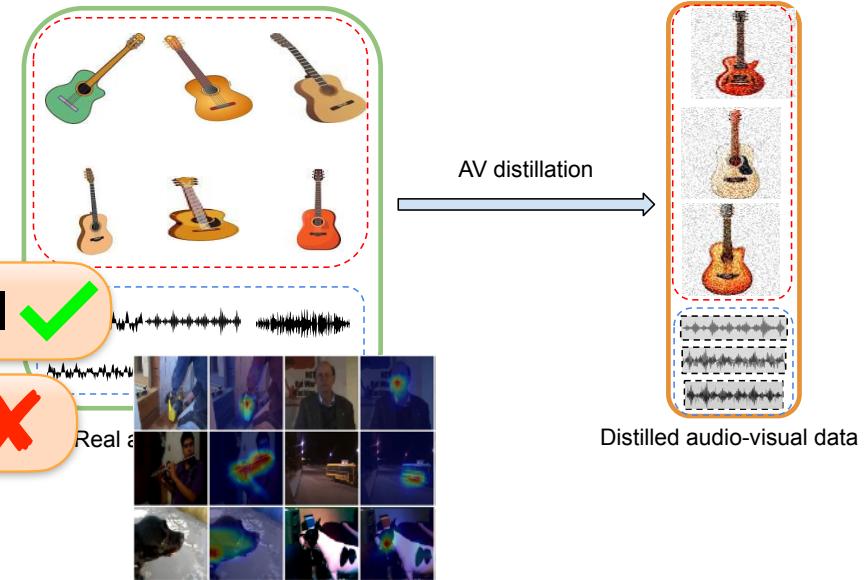
Future research directions

- Multimodal data compression/distillation
 - Current: Class-wise distillation

Classification ✓

Cross-modal retrieval ✓

Source localization ✗



Future research directions

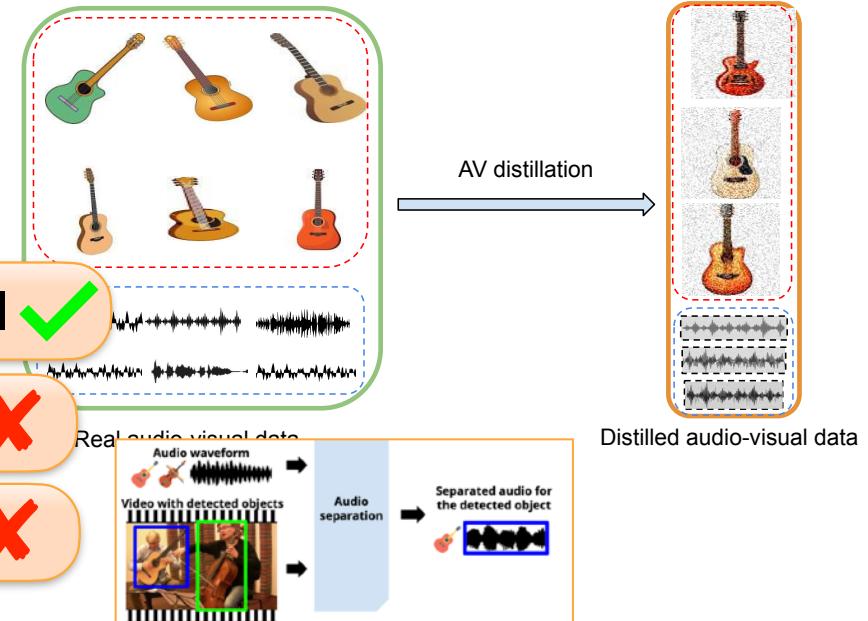
- Multimodal data compression/distillation
 - Current: Class-wise distillation

Classification ✓

Cross-modal retrieval ✓

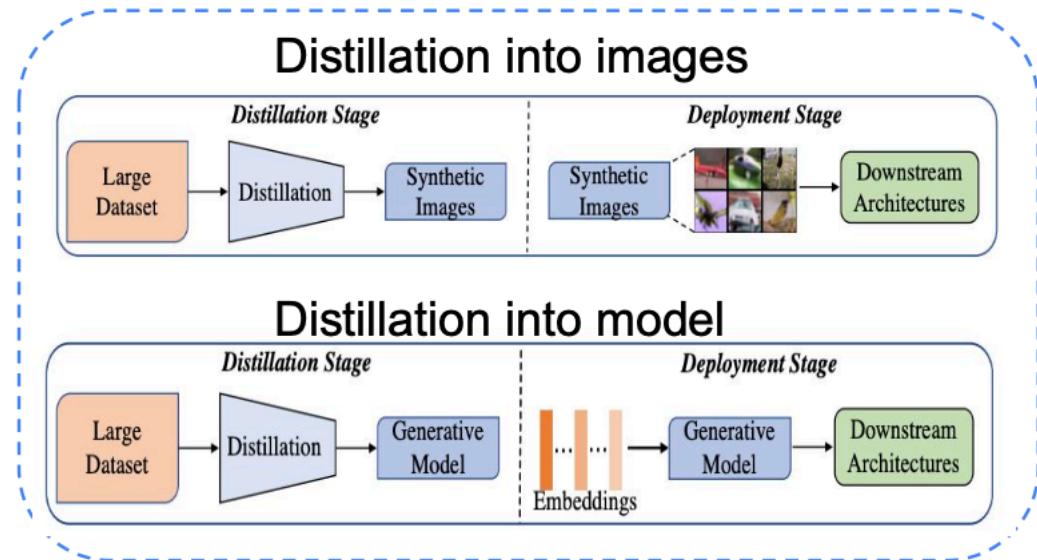
Source localization ✗

Source Separation ✗



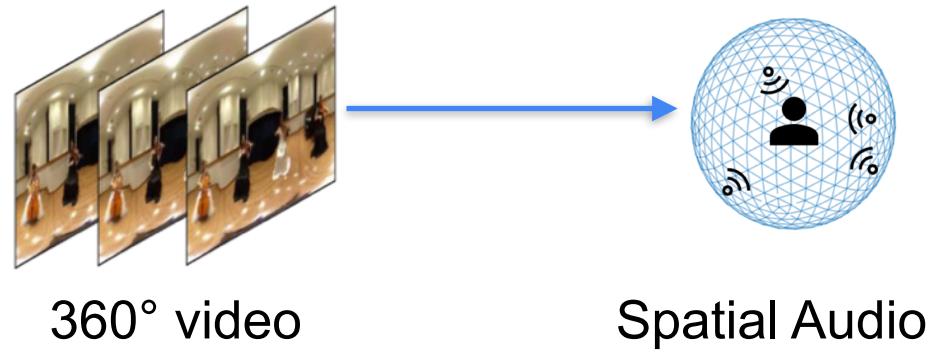
Future research directions

- Multimodal data compression/distillation
 - Current: Class-wise distillation
 - Distillation into model



Future research directions

- Spatial audio generation
 - Text/360° videos to spatial audio generation
 - Limited data
 - Computational requirements





Thank You

- Questions?