

# Beyond Modality Collapse: Representations Blending for Multimodal Dataset Distillation

Xin Zhang<sup>1,2</sup>, Ziruo Zhang<sup>1,3</sup>, Jiawei Du<sup>1,2</sup>, Zuozhu Liu<sup>4</sup>, Joey Tianyi Zhou<sup>1,2</sup> 

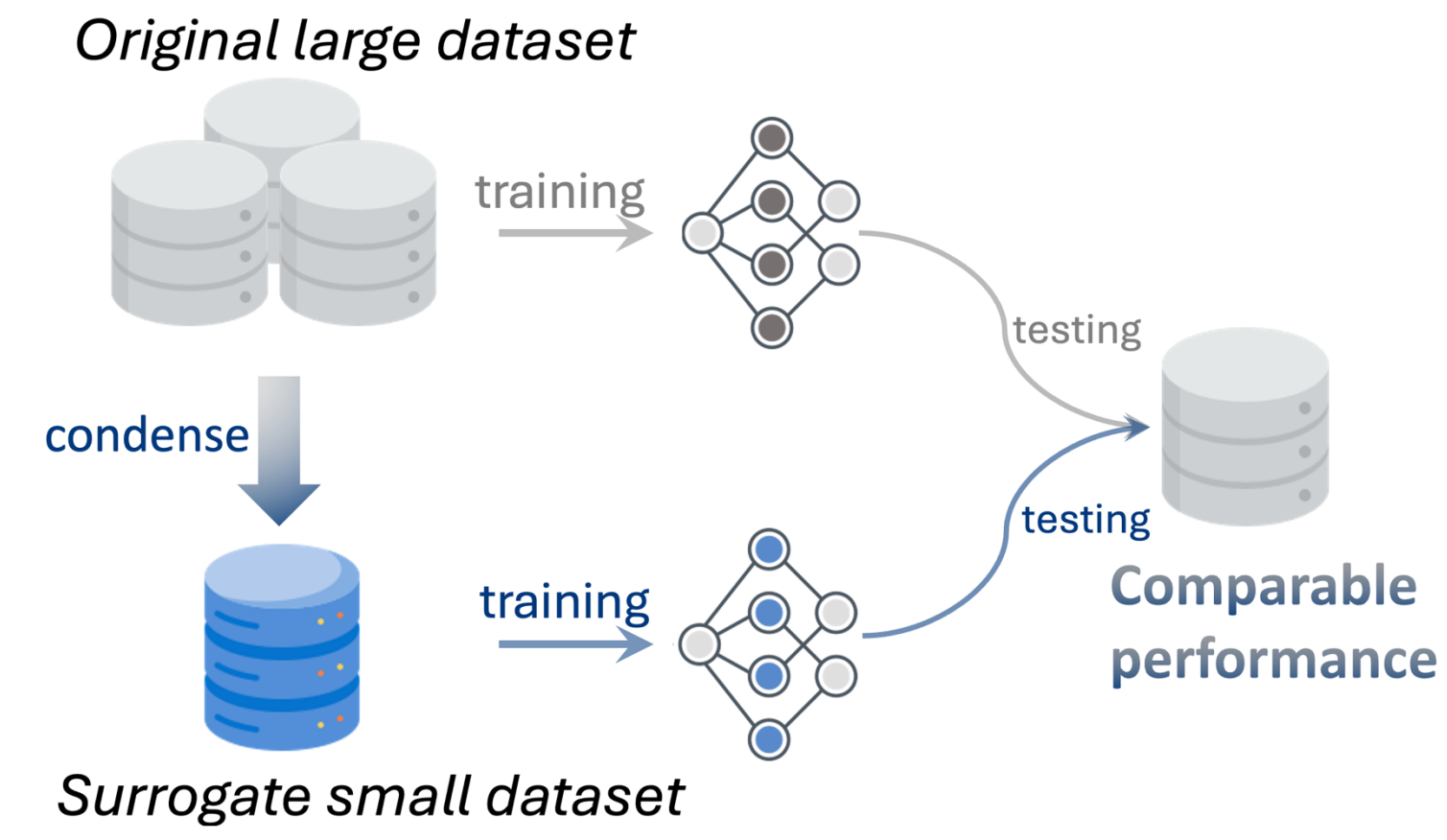
<sup>1</sup>Centre for Frontier AI Research (CFAR), A\*STAR, Singapore, <sup>2</sup>Institute of High Performance Computing (IHPC), A\*STAR, Singapore,

<sup>3</sup>National University of Singapore, Singapore, <sup>4</sup>Zhejiang University, China



## Background

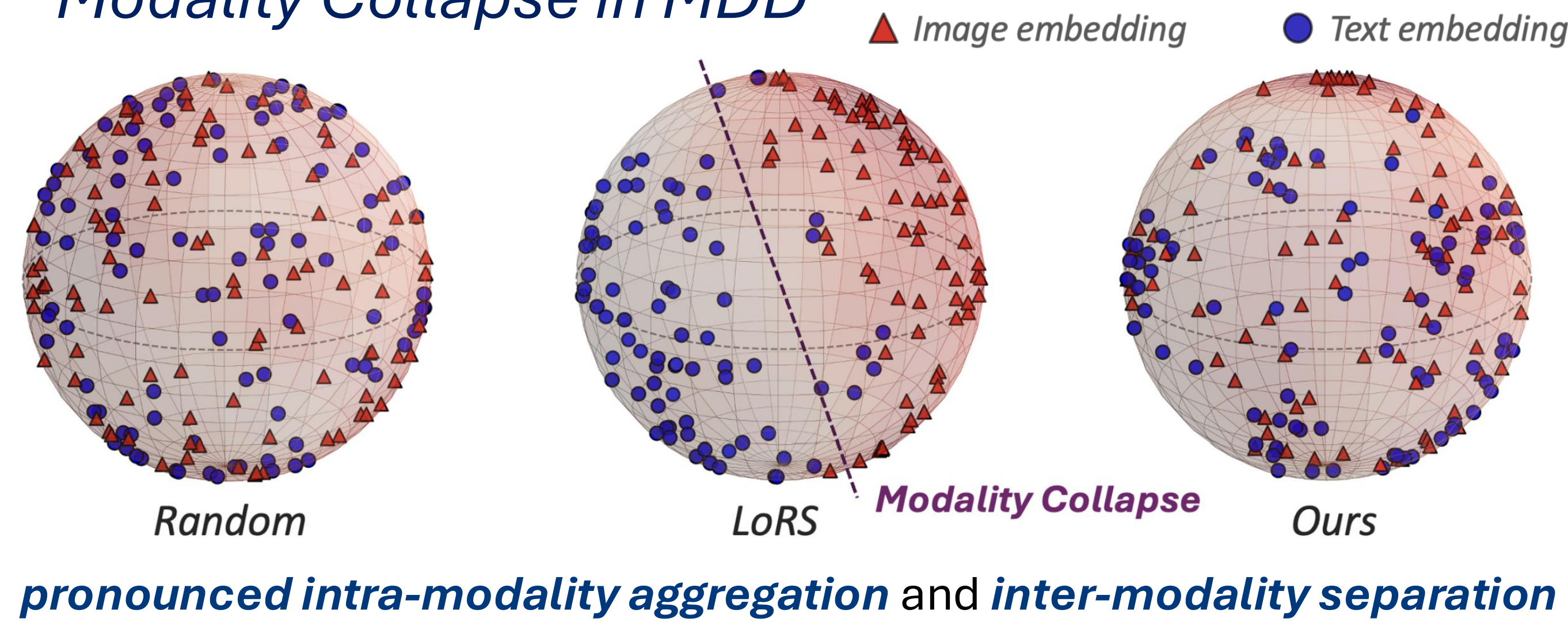
### What is Dataset Distillation (DD)?



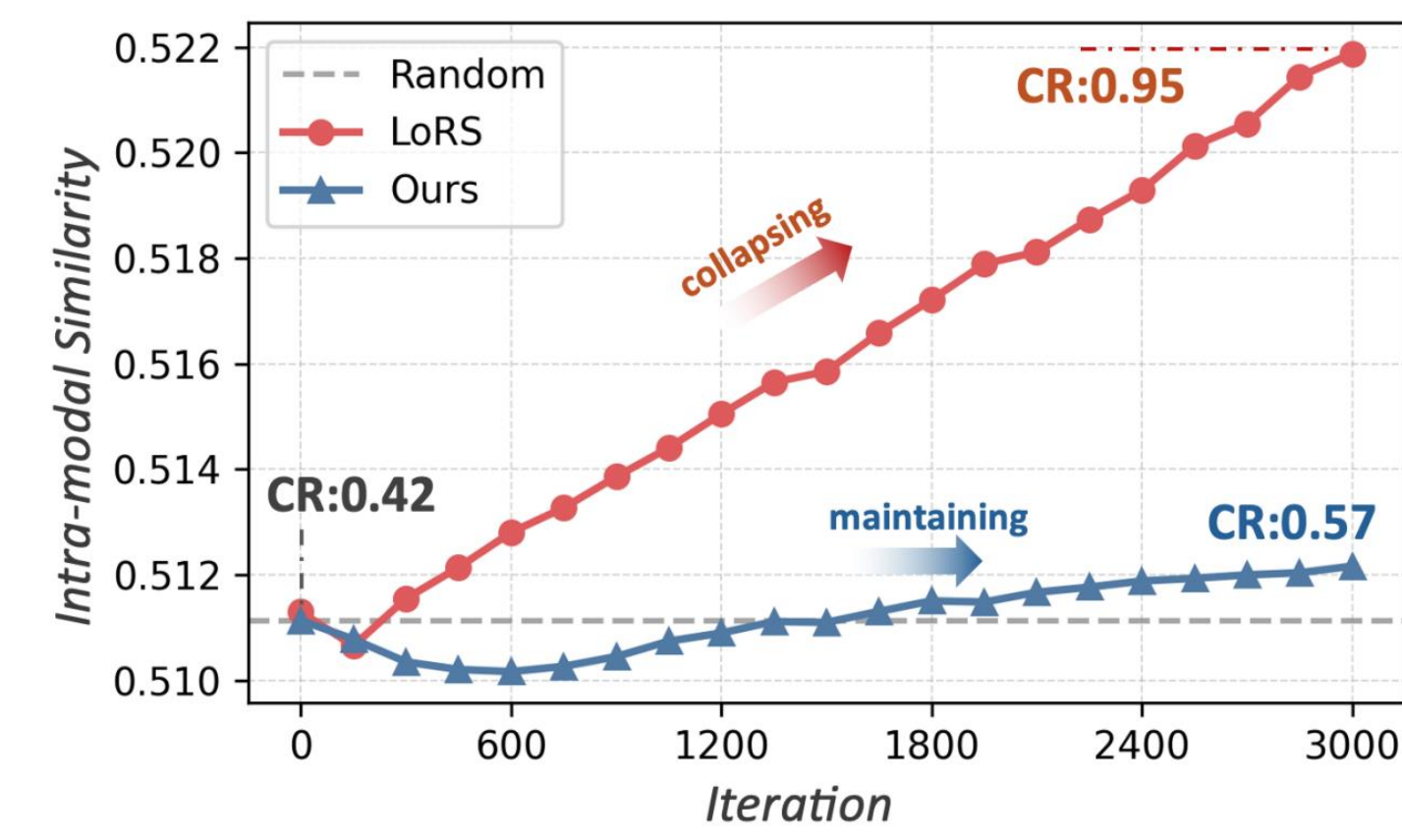
**Dataset distillation** is to **synthesize a tiny and compact dataset** from a given **real and large dataset**, such that the former can yield a **comparable performance** as the latter.

### How do existing methods achieve MDD?

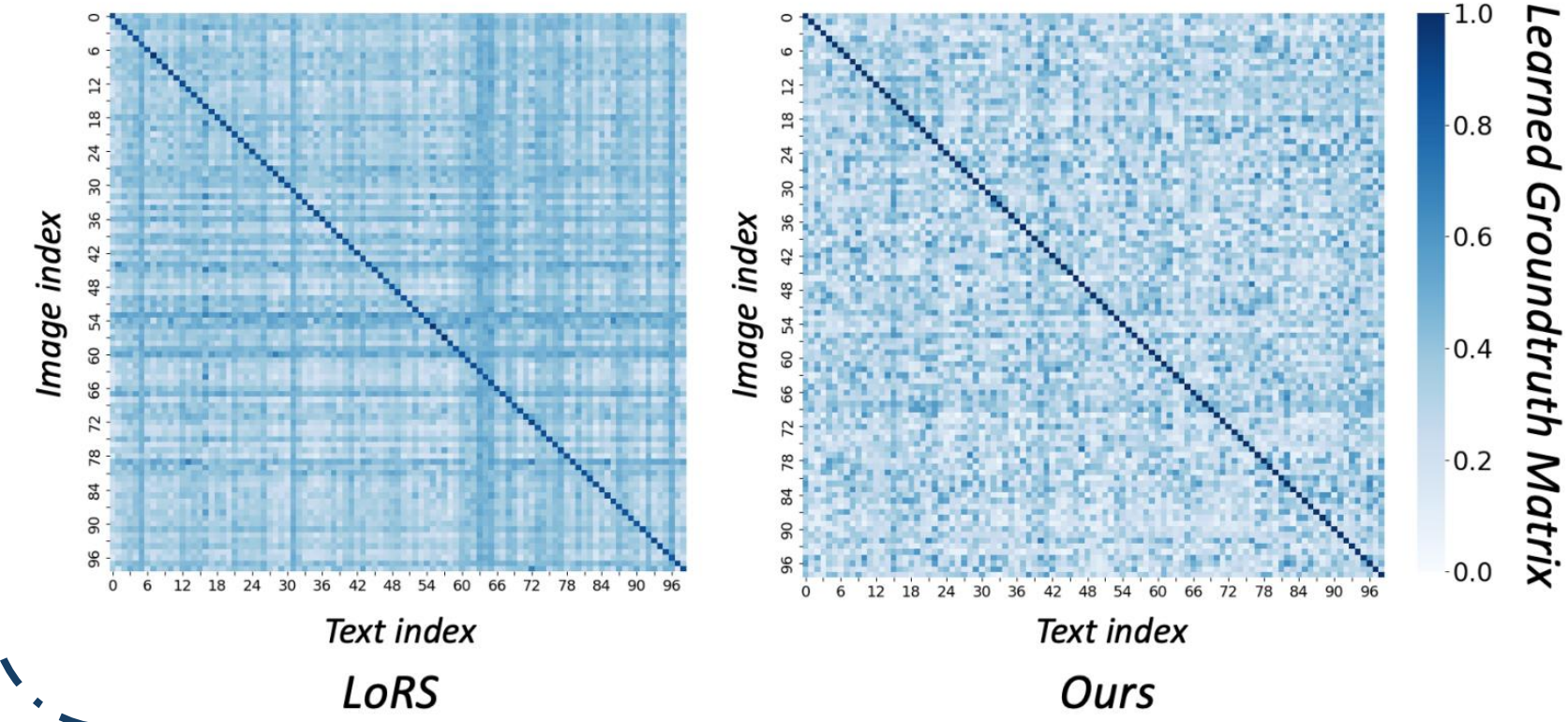
#### ❖ Modality Collapse in MDD



observable effects:



➤ The intra-modality similarity consistently increases throughout the distillation process.



➤ Feature centralization widens the modality gap, making non-matching pairs indistinguishable.

## Proposed Method

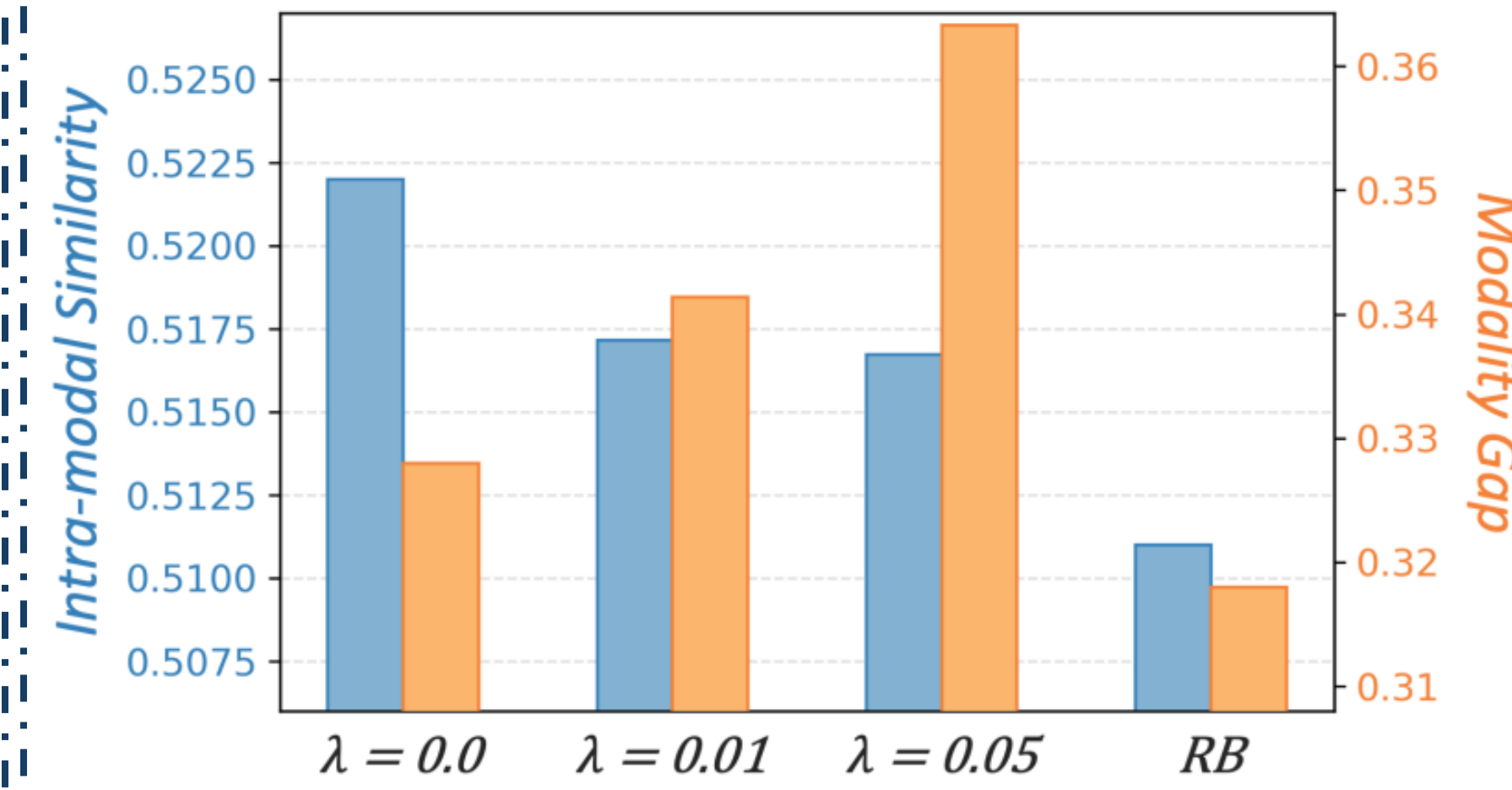
### Beyond Modality Collapse

❖ underlying theoretical cause

$$S^* = \arg \min_S \mathbb{E}_{(x, \tau) \sim \mathcal{P}} [\mathcal{L}(f_{\theta_S}(x, \tau), y)] \quad \text{s.t.} \quad \theta_S = \arg \min_{\theta} \mathbb{E}_{(\tilde{x}, \tilde{\tau}) \sim S} [\mathcal{L}(f_{\theta}(\tilde{x}, \tilde{\tau}), \tilde{y})],$$

$$\frac{\partial \mathcal{L}}{\partial \tilde{x}'_n} \frac{\partial \mathcal{L}}{\partial \tilde{x}'_m} \leftarrow \frac{w_{nm} w_{mn}}{\gamma^2} [\sigma(\hat{y}_{nm})/t - \tilde{y}_{nm}] [\sigma(\hat{y}_{mn})/t - \tilde{y}_{mn}] \tilde{\tau}'_m{}^T \tilde{\tau}'_n$$

✓ Mitigating Modality Collapse via Representation Blending

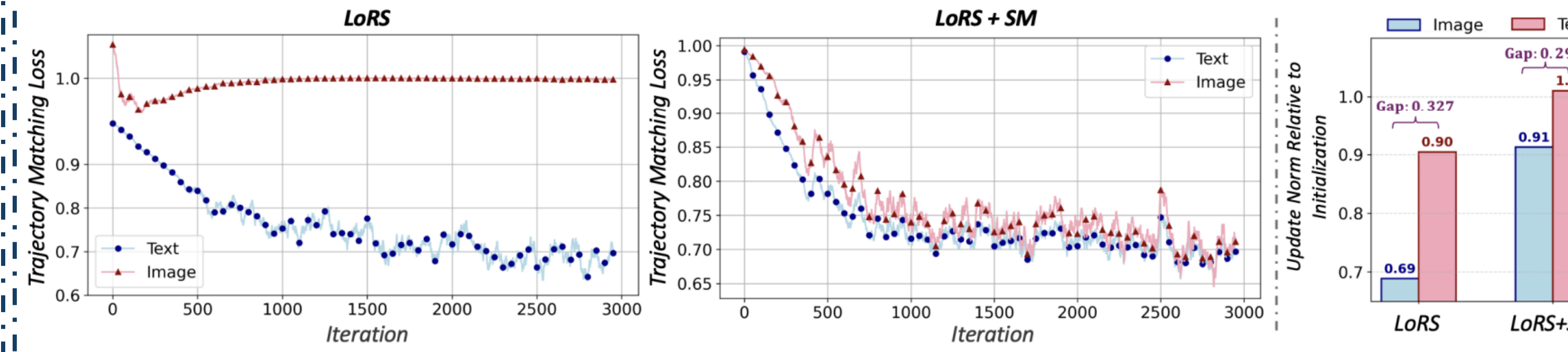


Reduced intra-modality similarity  
&  
Reduced inter-modality gap

$$\tilde{\tau}_m^{+/noise} = \text{Normalize} \left( f^{\text{textP}}((1 - \lambda)\tilde{\tau}_m + \lambda\tilde{\tau}_n) \right) \quad \tilde{\tau}_m^{\text{blend}} = \text{Normalize} \left( f^{\text{textP}}((1 - \lambda)\tilde{\tau}_m + \lambda\tilde{\tau}_i) \right)$$

$$\tilde{\tau}_n^{+/noise} = \text{Normalize} \left( f^{\text{textP}}((1 - \lambda)\tilde{\tau}_n + \lambda\tilde{\tau}_m) \right) \quad \tilde{\tau}_n^{\text{blend}} = \text{Normalize} \left( f^{\text{textP}}((1 - \lambda)\tilde{\tau}_n + \lambda\tilde{\tau}_j) \right)$$

✓ Enhancing Cross-modal Alignment via Symmetric Projection Trajectory Matching



$$\tilde{x}^*, \tilde{\tau}^*, \tilde{y}^* = \arg \min_{\tilde{x}, \tilde{\tau}, \tilde{y}} \left( \left\| \theta_{S_{\text{img}}}^{t+T} - \theta_{D_{\text{img}}}^{t+M} \right\|_2^2 + \left\| \theta_{S_{\text{text}}}^{t+T} - \theta_{D_{\text{text}}}^{t+M} \right\|_2^2 \right) / \left( \left\| \theta_{D_{\text{img}}}^t - \theta_{D_{\text{img}}}^{t+M} \right\|_2^2 + \left\| \theta_{D_{\text{text}}}^t - \theta_{D_{\text{text}}}^{t+M} \right\|_2^2 \right)$$

$$\tilde{x}^*, \tilde{\tau}^*, \tilde{y}^* = \arg \min_{\tilde{x}, \tilde{\tau}, \tilde{y}} \left( \left\| \theta_{S_{\text{img}}}^{t+T} - \theta_{D_{\text{img}}}^{t+M} \right\|_2^2 + \left\| \theta_{S_{\text{text}}}^{t+T} - \theta_{D_{\text{text}}}^{t+M} \right\|_2^2 \right) / \left( \left\| \theta_{D_{\text{img}}}^t - \theta_{D_{\text{img}}}^{t+M} \right\|_2^2 + \left\| \theta_{D_{\text{text}}}^t - \theta_{D_{\text{text}}}^{t+M} \right\|_2^2 \right)$$

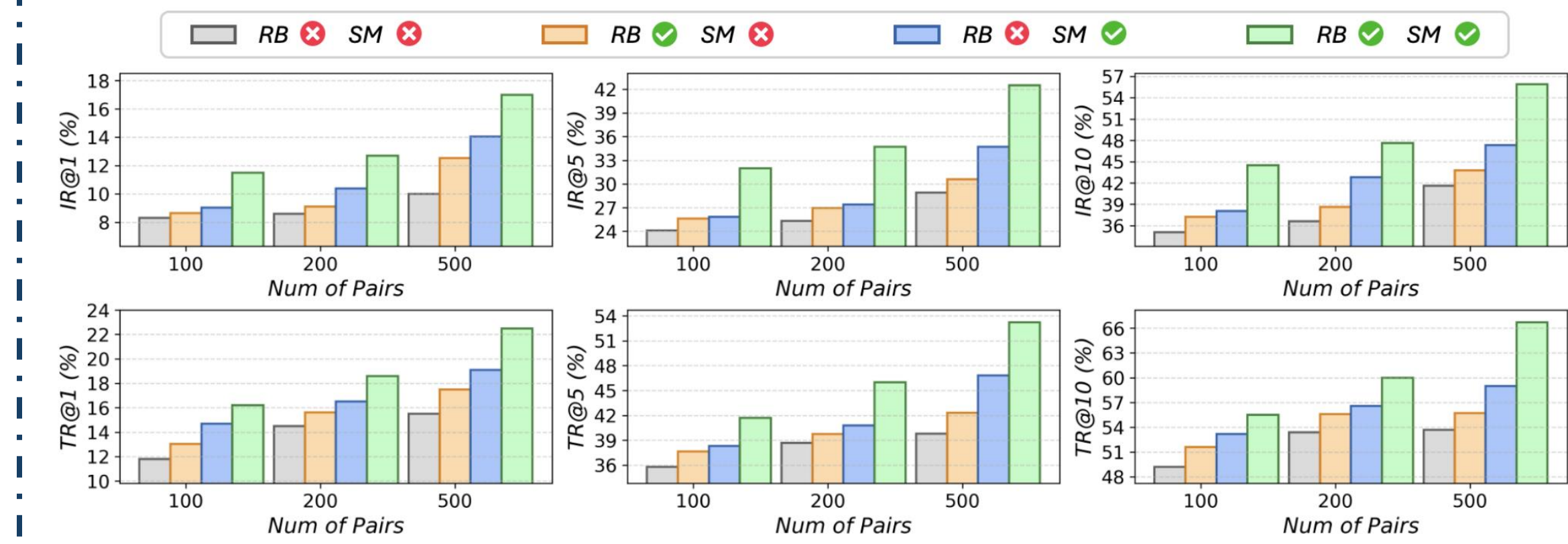
Enhance the image-side distillation

## Experiments

### Results:

Pairs	Ratio	Metric	Coreset Selection				Dataset Distillation			
			Rand	Herd [50]	K-Cent [16]	Forget [45]	MTT-VL [51]	TESLA-VL [52]	LoRS [52]	Ours
Flickr-30k										
500	1.7%	IR@1	2.4	3.0	3.5	1.8	6.6 $\pm$ 0.3	1.1 $\pm$ 0.2	10.0 $\pm$ 0.2	17.0 $\pm$ 0.6
		IR@5	10.5	10.0	10.4	9.0	20.2 $\pm$ 1.2	7.3 $\pm$ 0.4	28.9 $\pm$ 0.7	42.5 $\pm$ 0.5
		IR@10	17.4	17.0	17.3	15.9	30.0 $\pm$ 2.1	12.6 $\pm$ 0.5	41.6 $\pm$ 0.6	55.9 $\pm$ 0.6
		TR@1	5.2	5.1	4.9	3.6	13.3 $\pm$ 0.6	5.1 $\pm$ 0.2	15.5 $\pm$ 0.7	22.5 $\pm$ 0.4
		TR@5	18.3	16.4	16.4	12.3	32.8 $\pm$ 1.8	15.3 $\pm$ 0.5	39.8 $\pm$ 0.4	53.2 $\pm$ 0.3
		TR@10	25.7	24.3	23.3	19.3	46.8 $\pm$ 0.8	23.8 $\pm$ 0.3	53.7 $\pm$ 0.3	66.7 $\pm$ 0.3
COCO										
500	4.4%	IR@1	1.1	1.7	1.1	0.8	2.5 $\pm$ 0.5	0.8 $\pm$ 0.2	2.8 $\pm$ 0.2	6.2 $\pm$ 0.1
		IR@5	5.0	5.3	6.3	5.8	8.9 $\pm$ 0.7	3.6 $\pm$ 0.6	9.9 $\pm$ 0.5	19.9 $\pm$ 0.3
		IR@10	8.7	9.9	10.5	8.2	15.8 $\pm$ 1.5	6.7 $\pm$ 0.9	16.5 $\pm$ 0.7	30.6 $\pm$ 0.1
		TR@1	1.9	1.9	2.5	2.1	5.0 $\pm$ 0.4	1.7 $\pm$ 0.4	5.3 $\pm$ 0.5	7.0 $\pm$ 0.2
		TR@5	7.5	7.8	8.7	8.2	17.2 $\pm$ 1.3	5.9 $\pm$ 0.8	18.3 $\pm$ 1.5	22.0 $\pm$ 0.3
		TR@10	12.5	13.7	14.3	13.0	26.0 $\pm$ 1.9	10.2 $\pm$ 1.0	27.9 $\pm$ 1.4	32.9 $\pm$ 0.6

Comparison with SOTAs on Flickr-30k & COCO



The ablation study of RepBlend

Methods	ImageNet-10 Classification		TextCaps Retrieval					
	ACC@1	ACC@5	IR@1	IR@5	IR@10	TR@1	TR@5	TR@10
LoRS [55]	21.4	74.4	1.7	5.1	8.4	0.4	1.7	3.1
Ours	27.6	76.2	3.1	9.4	14.5	1.9	6.2	10.3

Zero-Shot Generalization

Methods	LoRS [52]	Ours
(IR@1, TR@1) (%)	(8.3, 11.8)	(11.5, 16.2)
<b>Buffer</b>		
Speed (min/traj)	70	40
Memory (GB/traj)	1.63	0.73
<b>Distillation</b>		
Speed (s/iter)	11.5	1.71
Peak GPU VRAM (GB)	21.78	10.17

**Better performance & higher speed**



Visualization

### References:

Beyond Modality Collapse: Representations Blending for Multimodal Dataset Distillation. **Zhang, Xin** and Zhang, Ziruo, and Du, Jiawei and Liu, Zuozhu and Zhou, Joey Tianyi. NeurIPS 2025.

Low-Rank Similarity Mining for Multimodal Dataset Distillation. Xu, Yue and Lin, Zhilin and Qiu, Yusong and Lu, Cewu and Li, Yong-Lu. ICML 2024