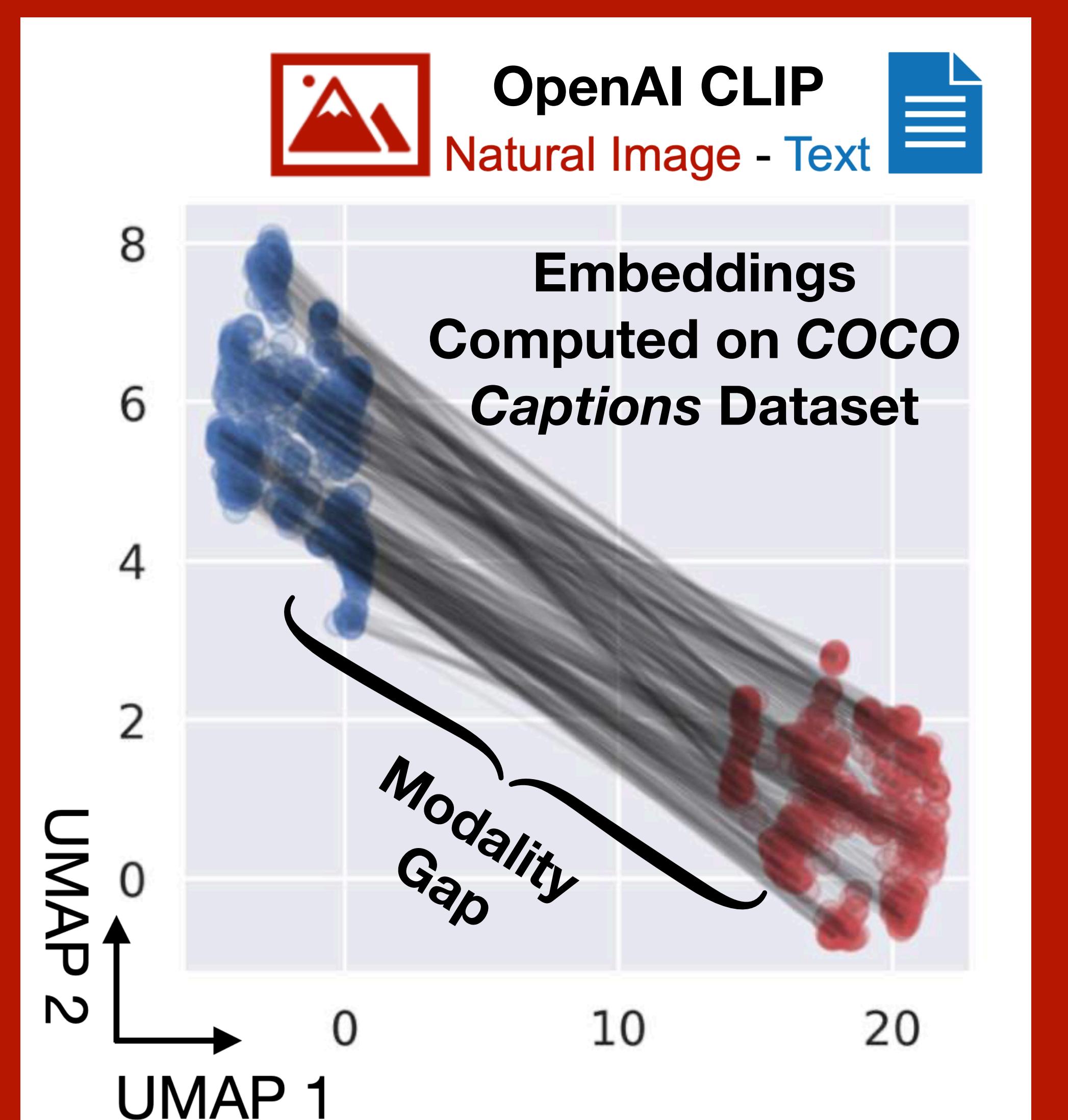
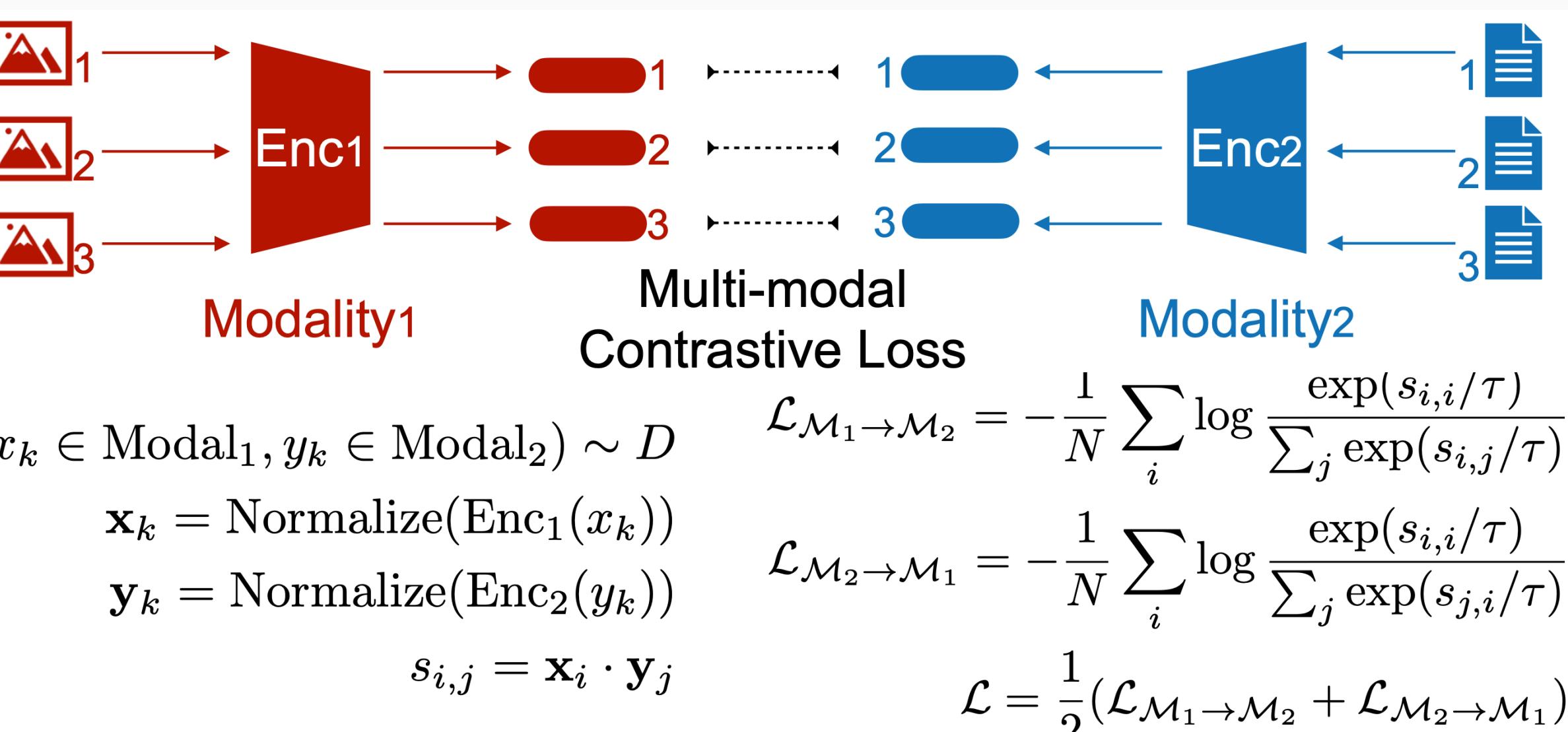


We present **modality gap**, an intriguing geometric phenomenon of multi-modal representation space, caused by a combination of model **initialization** and **contrastive learning optimization**.

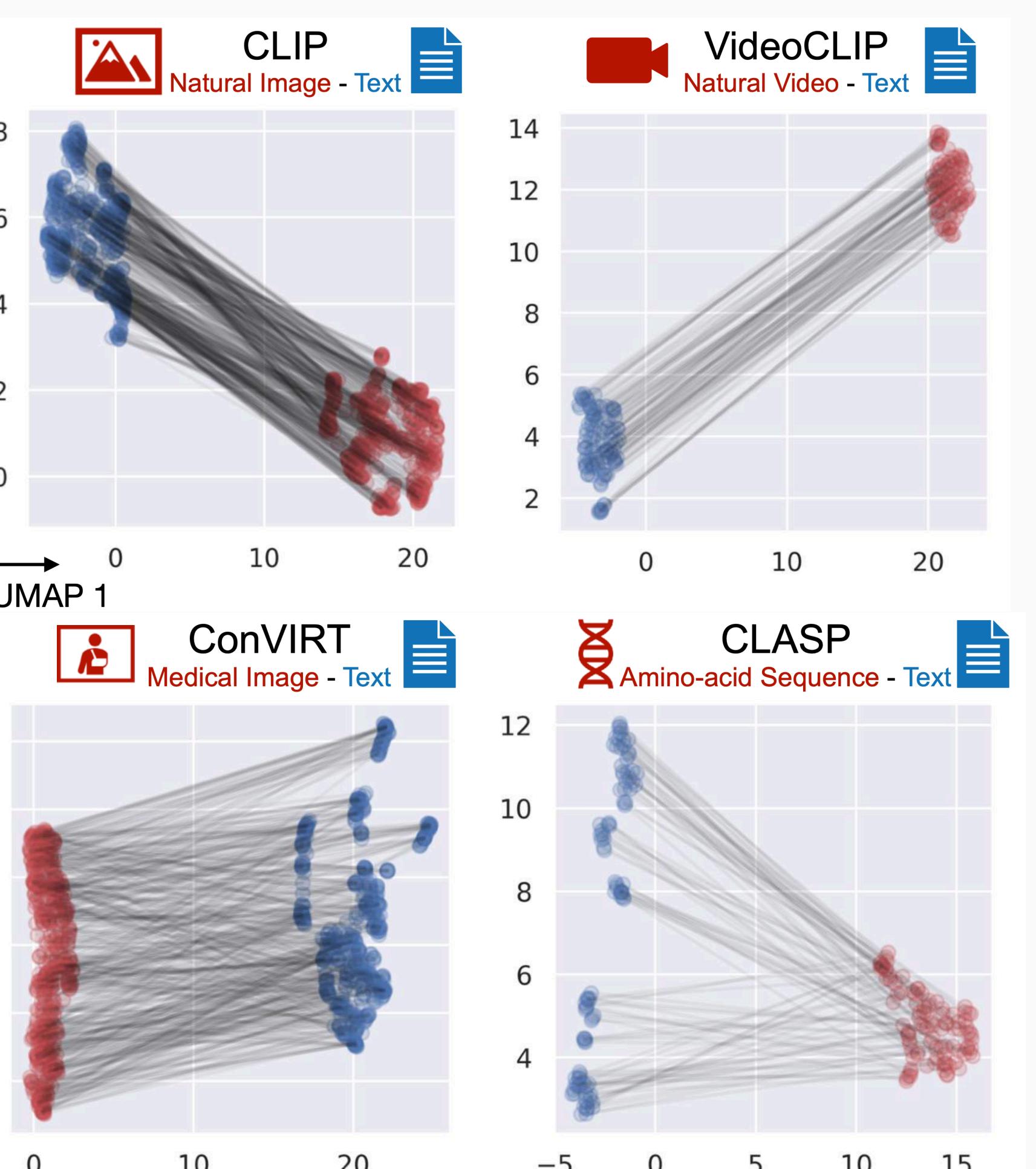


Paper

Modality Gap in Multi-Modal Contrastive Learning



Multi-modal Contrastive Learning: Paired inputs from different data modalities are mapped into a shared representation space (e.g., OpenAI CLIP).



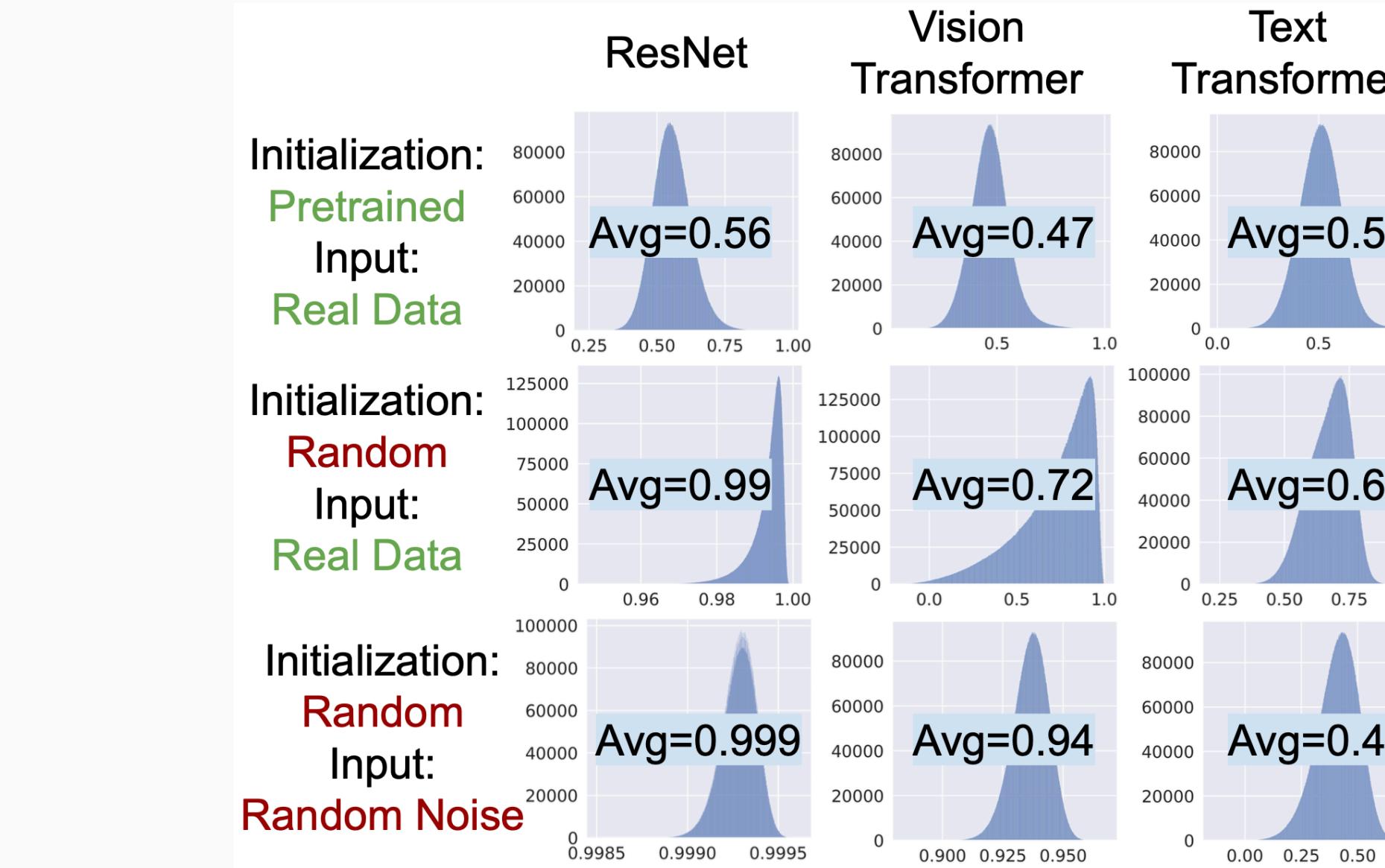
Pervasive Modality Gap: Clear gap is observed when visualizing embeddings from different modalities in 2D using UMAP (lines indicate pairs).

Modality gap is caused by ① + ②

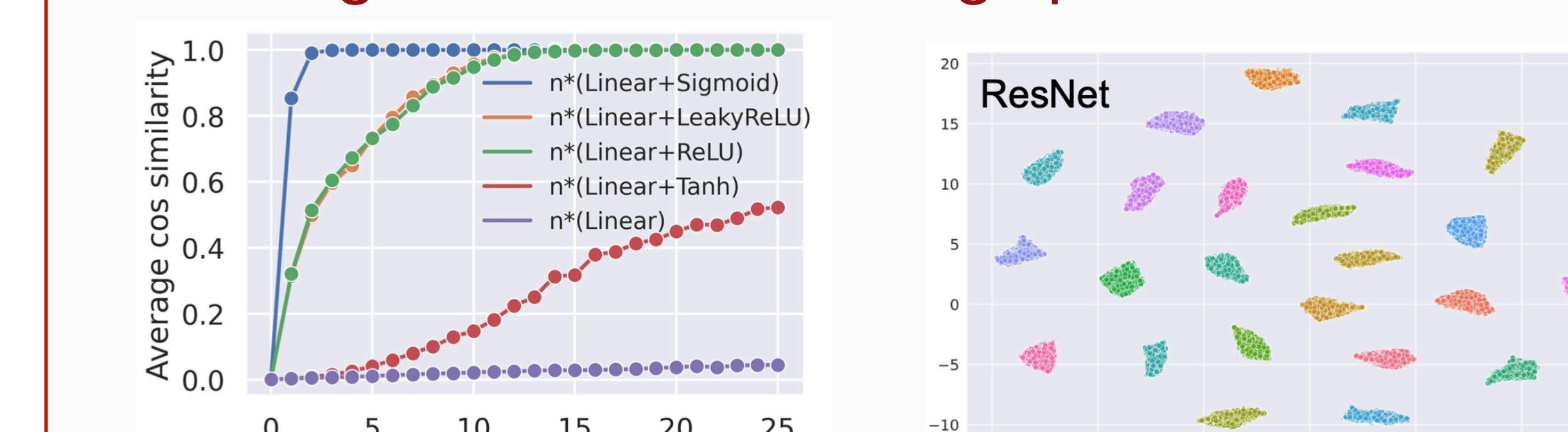
Mind the Gap: Understanding the Modality Gap in Multi-modal Contrastive Representation Learning

Yuhui Zhang*, Weixin Liang*, Yongchan Kwon*, Serena Yeung, James Zou

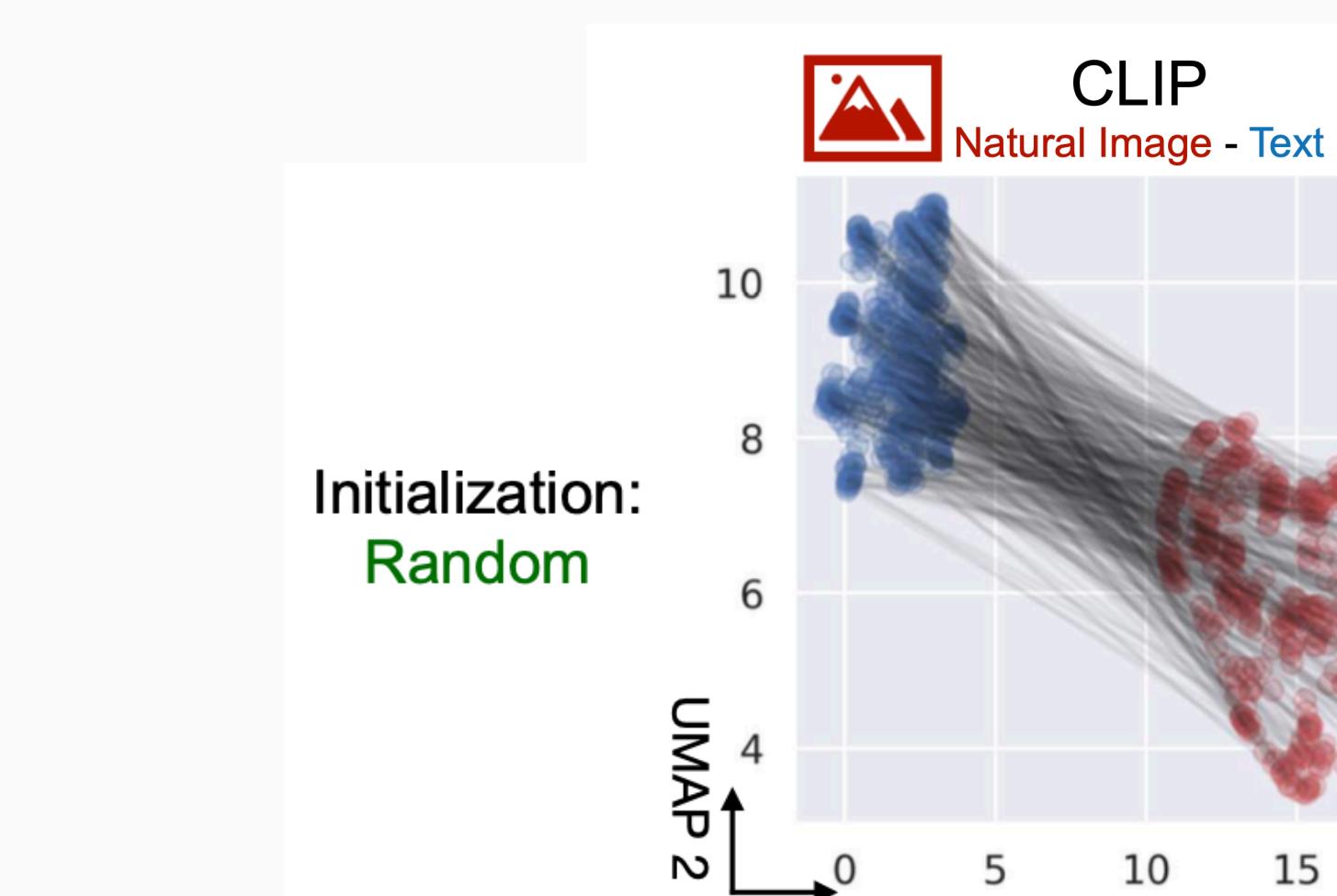
① Initialization: The Cone Effect Induces A Modality Gap



Cone Effect: The average cosine similarity between all pairs of embeddings is substantially larger than 0, indicating that the embedding space is a narrow cone.

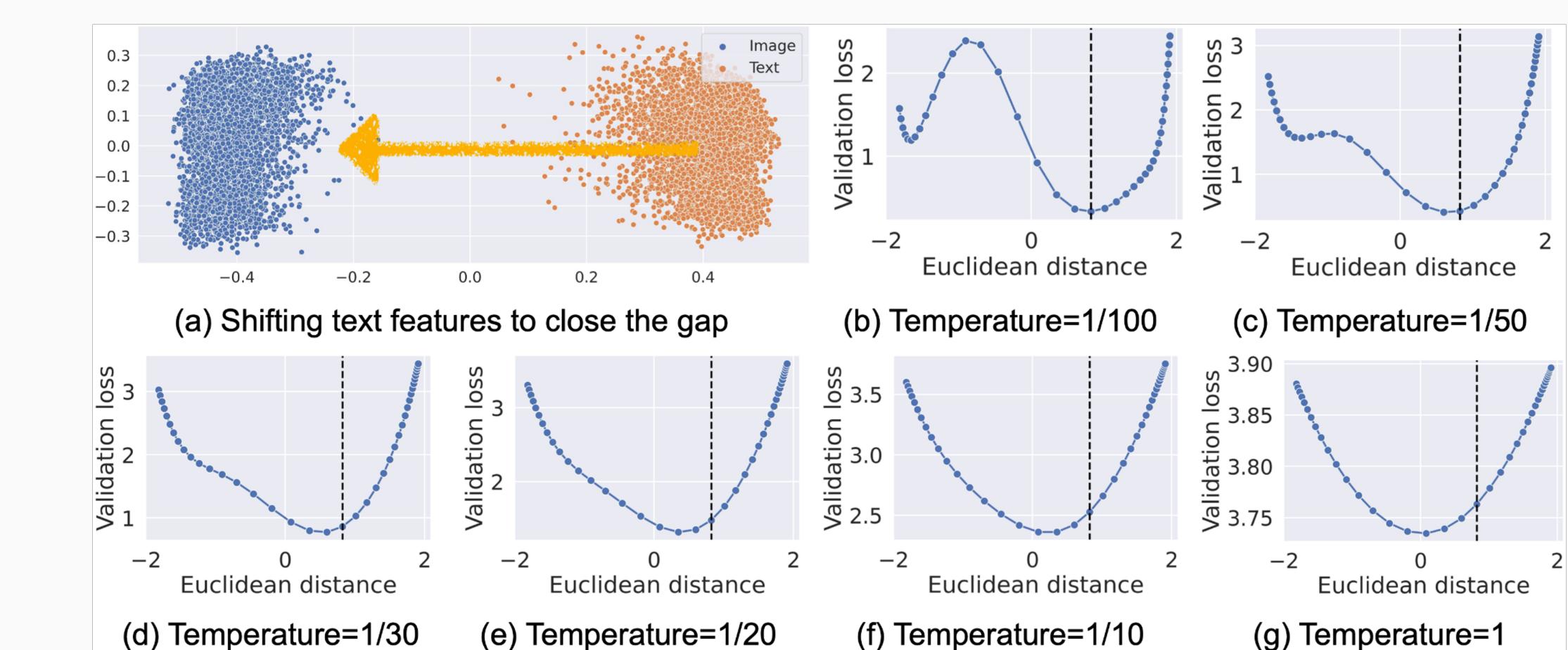


Analysis of Cone Effect: Deeper networks create narrower cones; different random initializations create distinctively different cones.

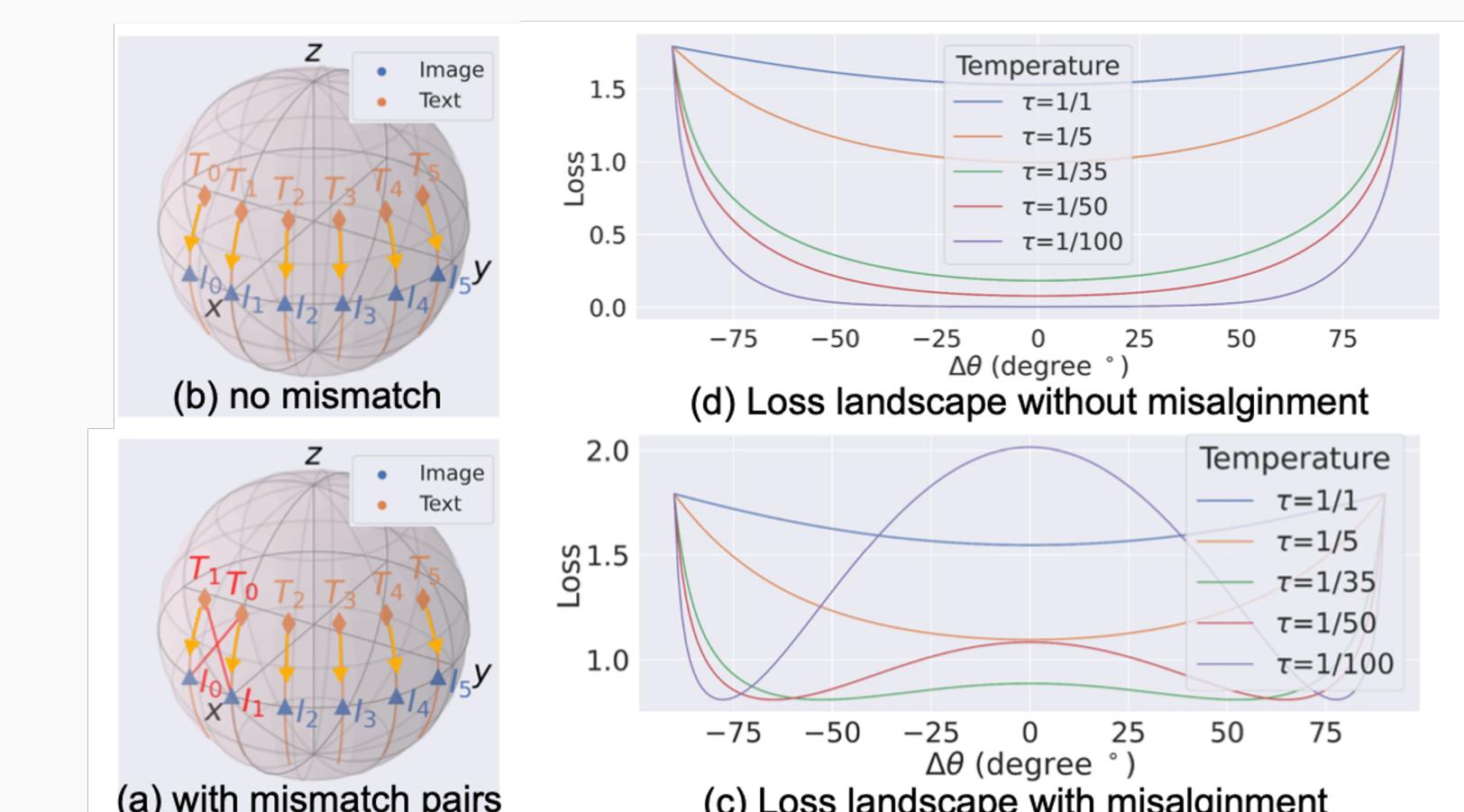


Cone Effect Induces Modality Gap: For multi-modal models with two encoders, representations of two modalities are embedded into two different cones at the initialization stage, thus creating the modality gap.

② Optimization: Contrastive Learning Preserves Modality Gap



Embedding Shift Analysis: The landscape of contrastive loss when manually shifting the embeddings from two modalities. The gap is optimal at default temperature.



Simulation Analysis: Six embedding pairs on a 3D sphere; text embeddings are shifted towards closing the modality gap. Mismatched pairs encourages the gap.

Modality Gap Implications

Dataset	Original gap	Modified gap	Direction	Denigration Biases	Original gap	Modified gap				
				Crime	Non related human	Sum	Crime	Non related human	Sum	
Coarse-grained Classification										
CIFAR10	0.9013	0.9081	↑	Black	1.0%	0.1%	1.1%	0.8%	0.1%	1.0%
CIFAR100	0.6658	0.6737	↓	White	15.5%	0.2%	15.7%	13.2%	0.4%	1.1%
Fine-grained Classification				Indian	1.2%	0.0%	1.2%	1.1%	0.0%	1.1%
EuroSAT	0.5410	0.5645	↓	Latino	2.8%	0.1%	2.8%	1.9%	0.1%	2.0%
Optical Character Recognition				Middle Eastern	6.3%	0.0%	6.3%	5.2%	0.0%	5.2%
SVHN	0.5389	0.5396	↑	Southeast Asian	0.5%	0.0%	0.5%	0.3%	0.0%	0.3%
HatefulMemes	0.5800	0.5811	↑	East Asian	0.7%	0.0%	0.7%	0.6%	0.0%	0.6%

Modifying the modality gap can improve zero-shot performances for downstream tasks and reduces model biases.