

Towards Data-Centric Multimodal ML: Vision-Language Dataset Distillation



Xindi Wu

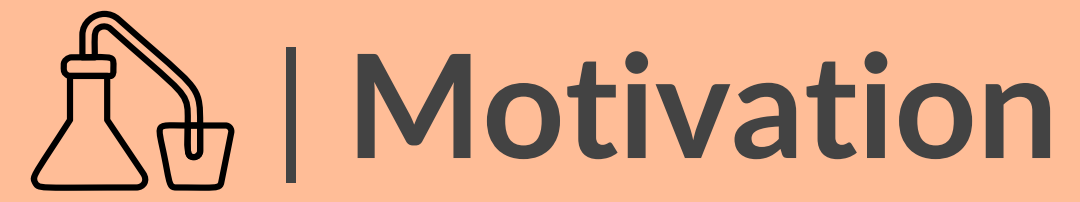
Department of Computer Science
Princeton University

<https://xindiwu.github.io/>
xindiw@princeton.edu



Overview

- 1 Motivation
- 2 Previous Work
- 3 Method
- 4 Experiments
- 5 Conclusion



| Motivation

 | **Motivation**

What is Data-centric
Multimodal ML?

Why is dataset
distillation important?



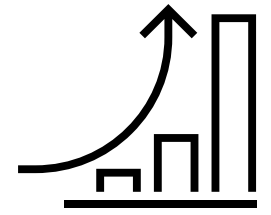
What is Multimodal?

A dictionary definition:

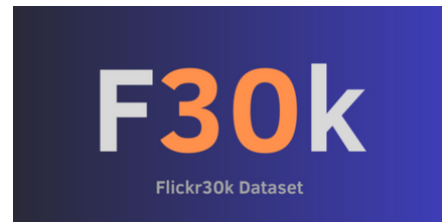
Multimodal: with multiple modalities

A research-oriented definition:

Multimodal is the science of heterogeneous and interconnected data



Data is the cornerstone of innovation in multimodal ML



Flickr30K (2014)



Visual Genome (2016)

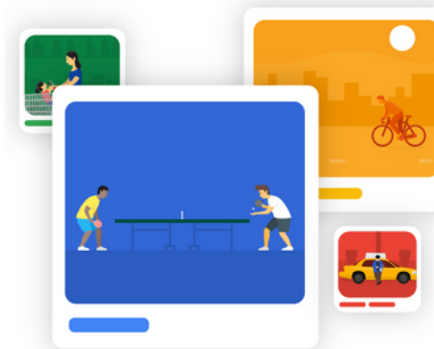


LAION 400M/5B (2021)

COCO (2014)

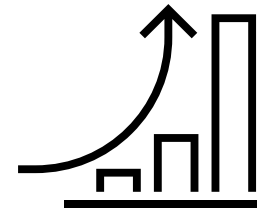


Conceptual Captions (2018)



DataComp (2023)





Data-Centric Multimodal ML



Flickr30K (2014)



Visual Genome (2016)



LAION 400M/5B (2021)

COCO (2014)



Conceptual Captions (2018)



DataComp (2023)





Data-Centric Multimodal ML

Data = Information + Noise



Data-Centric Multimodal ML

Data = Information + Noise

How can we identify the most critical information from datasets?



Data-Centric Multimodal ML

Data = **Information** + Noise

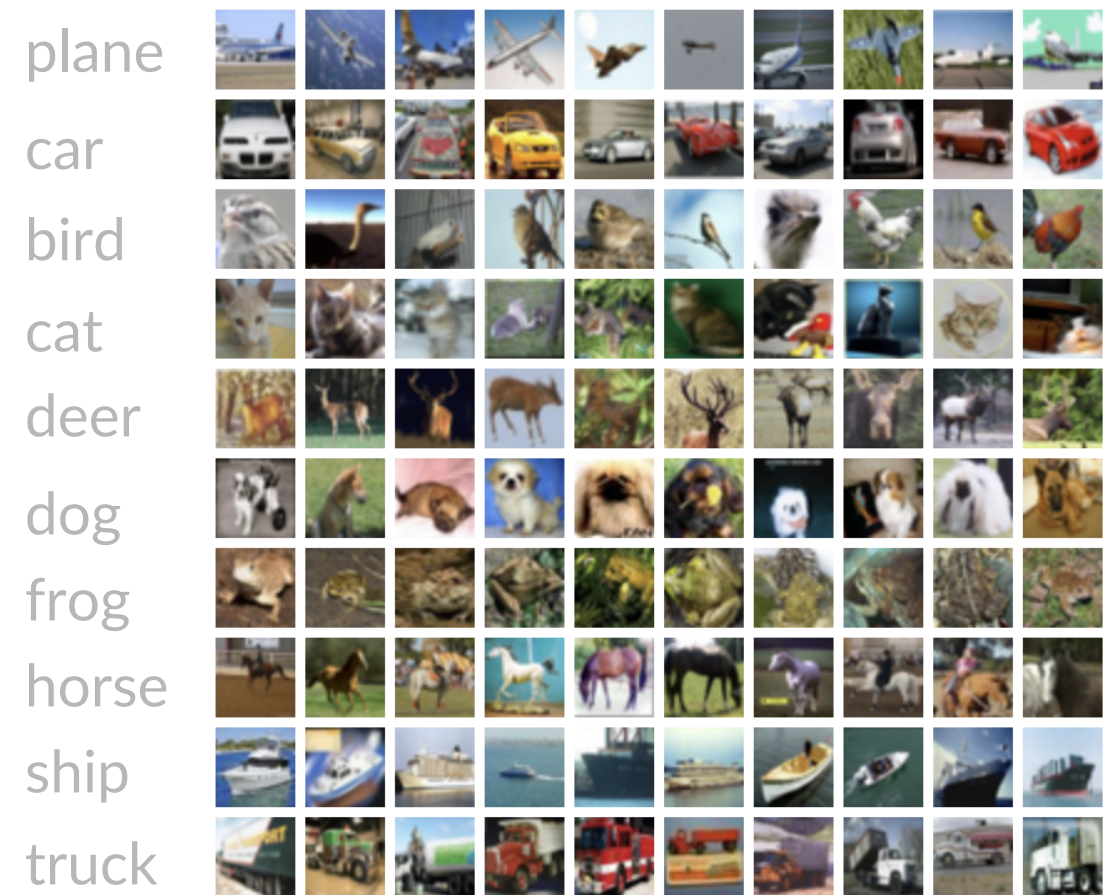
How can we identify the most critical **information** from datasets?

Dataset distillation is one promising solution!

What is Dataset Distillation?

CIFAR10 examples

N images



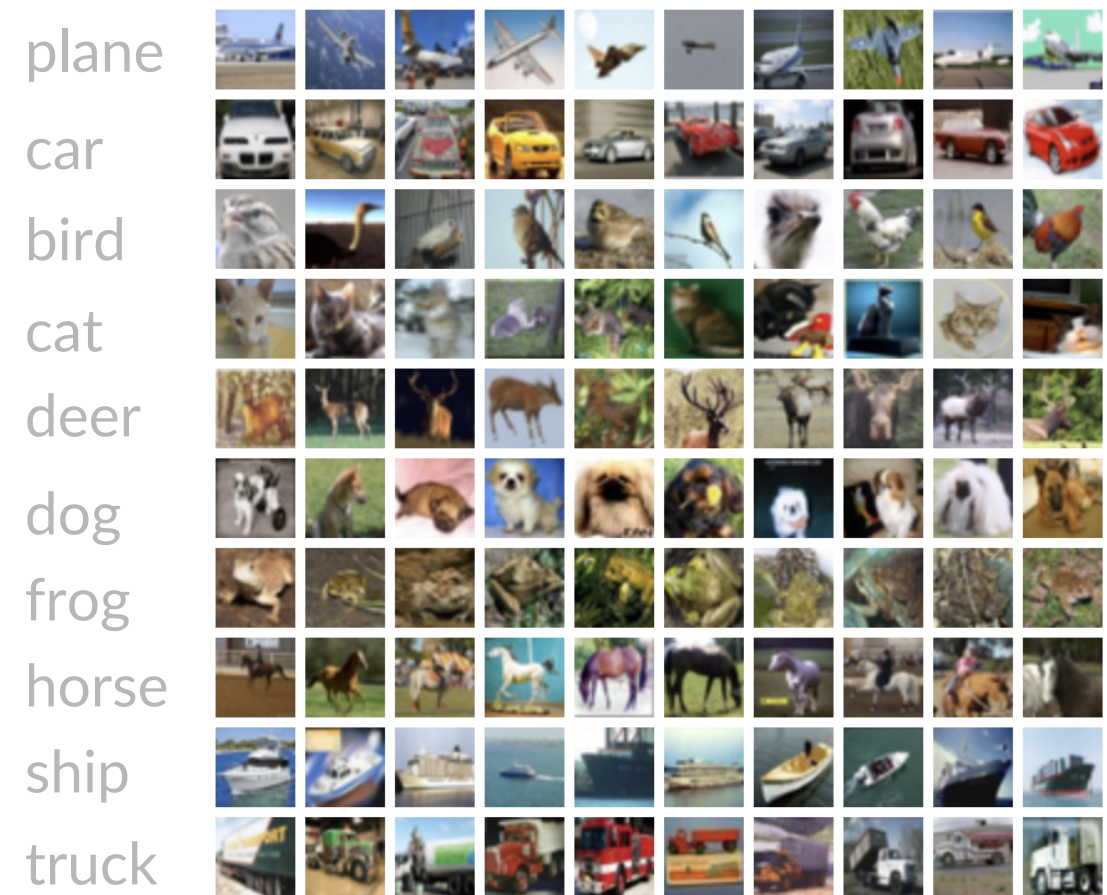
CIFAR10

(5000 images/class)

What is Dataset Distillation?

CIFAR10 examples

N images



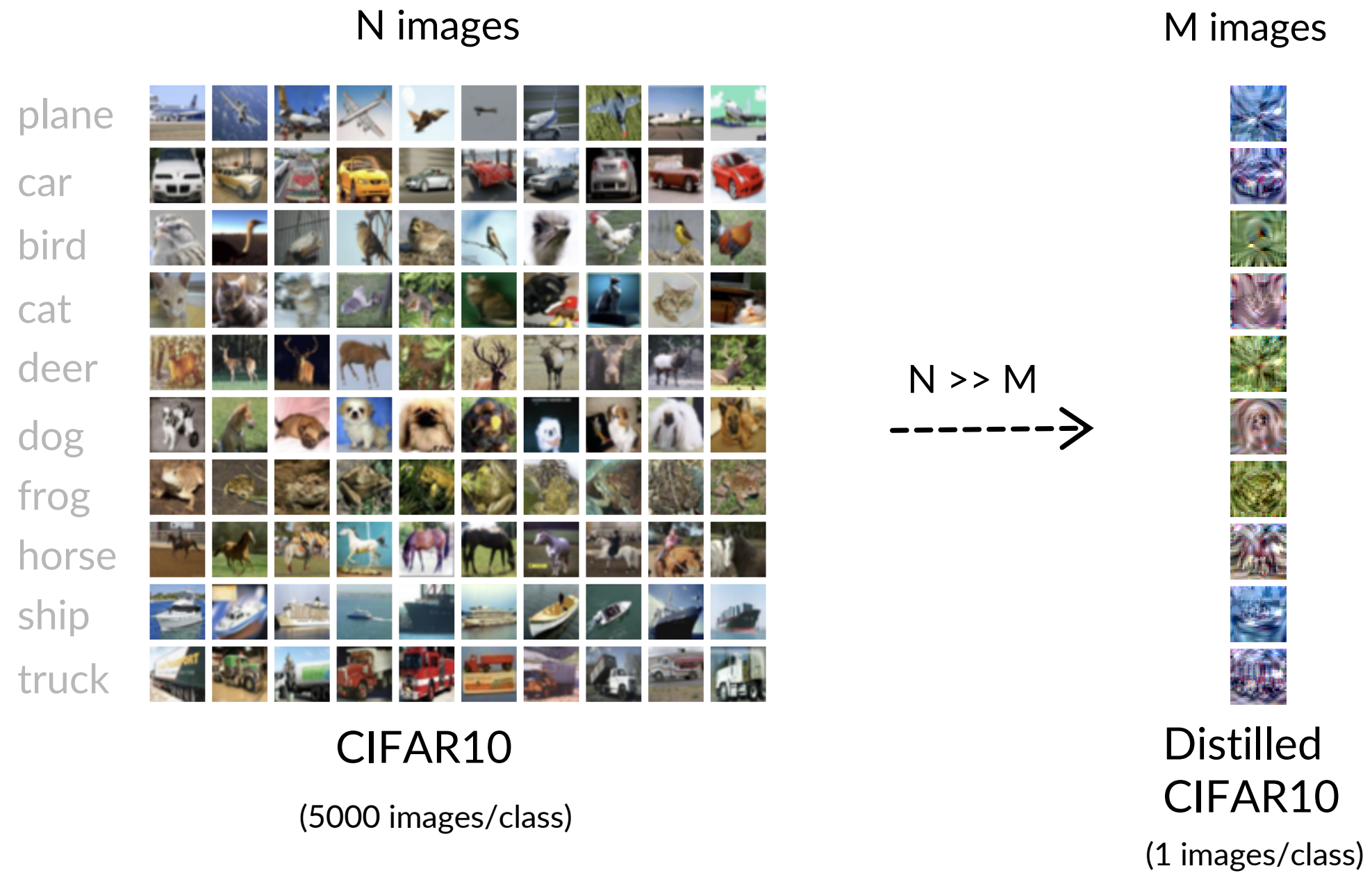
$N \gg M$
----->

CIFAR10

(5000 images/class)

What is Dataset Distillation?

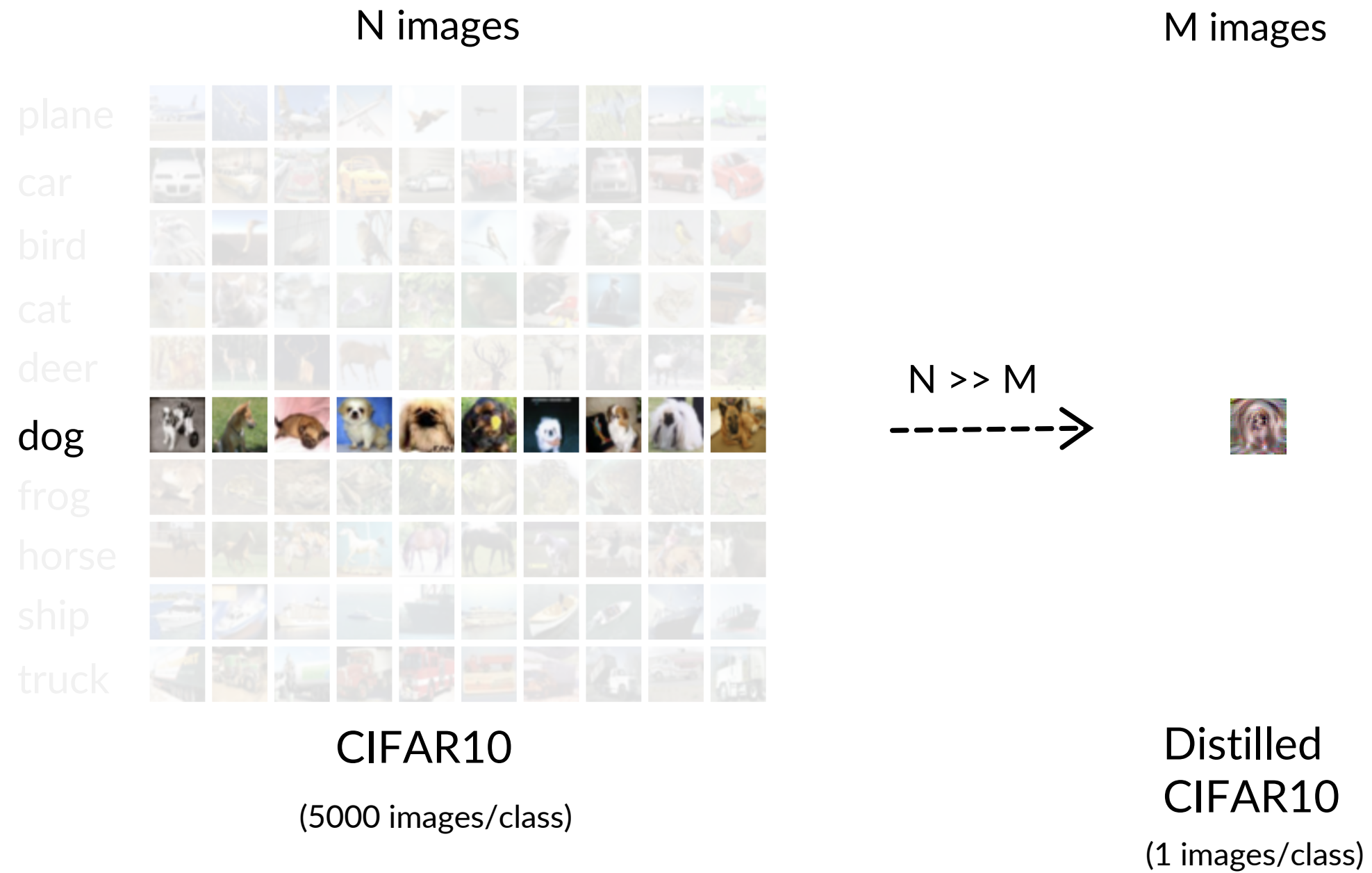
CIFAR10 examples





What is Dataset Distillation?

CIFAR10 examples

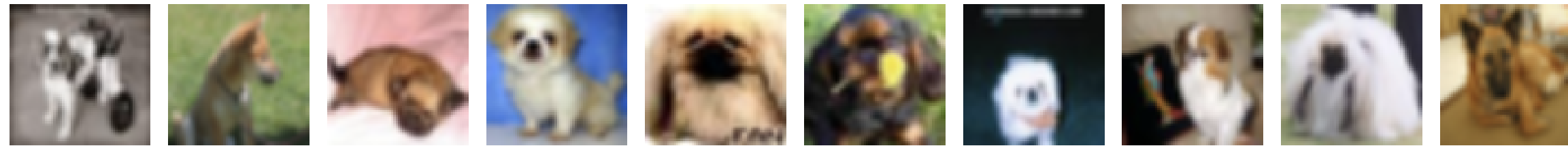


What is Dataset Distillation?

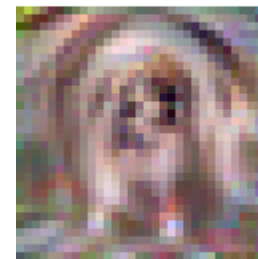
CIFAR10 examples

dog

CIFAR10
(5000 images/class)

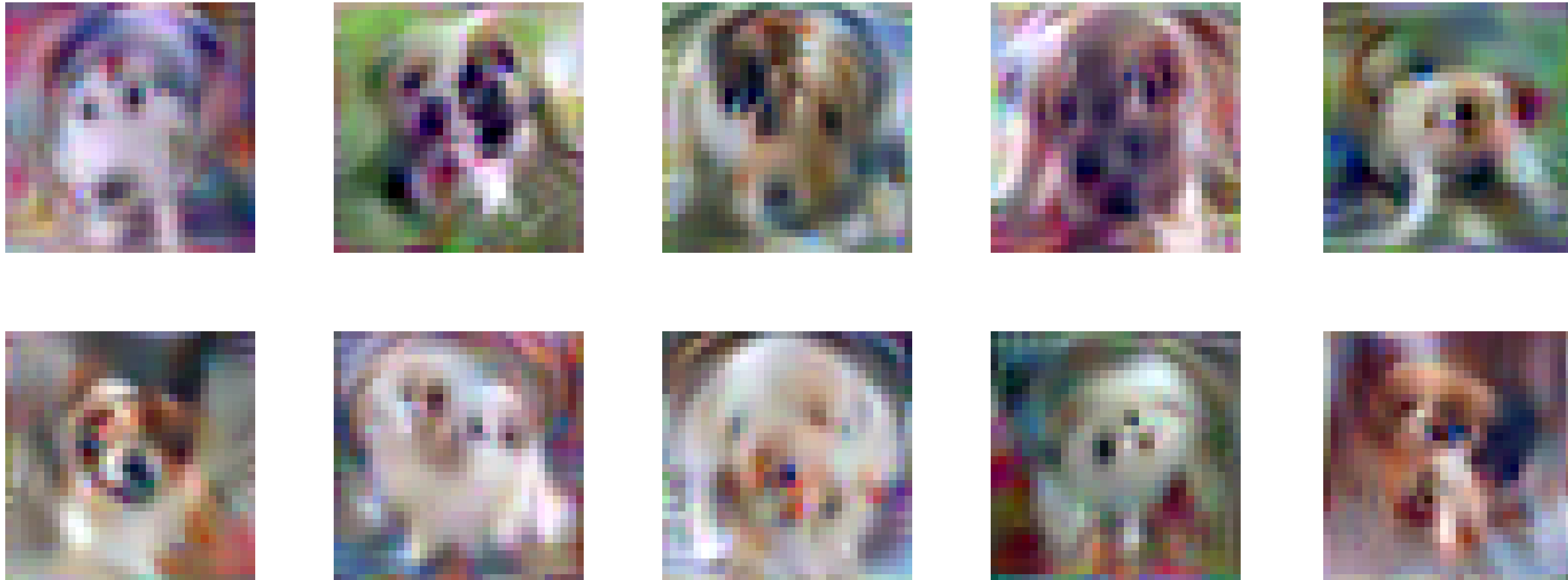


Distilled
CIFAR10
(1 images/class)



What is Dataset Distillation?

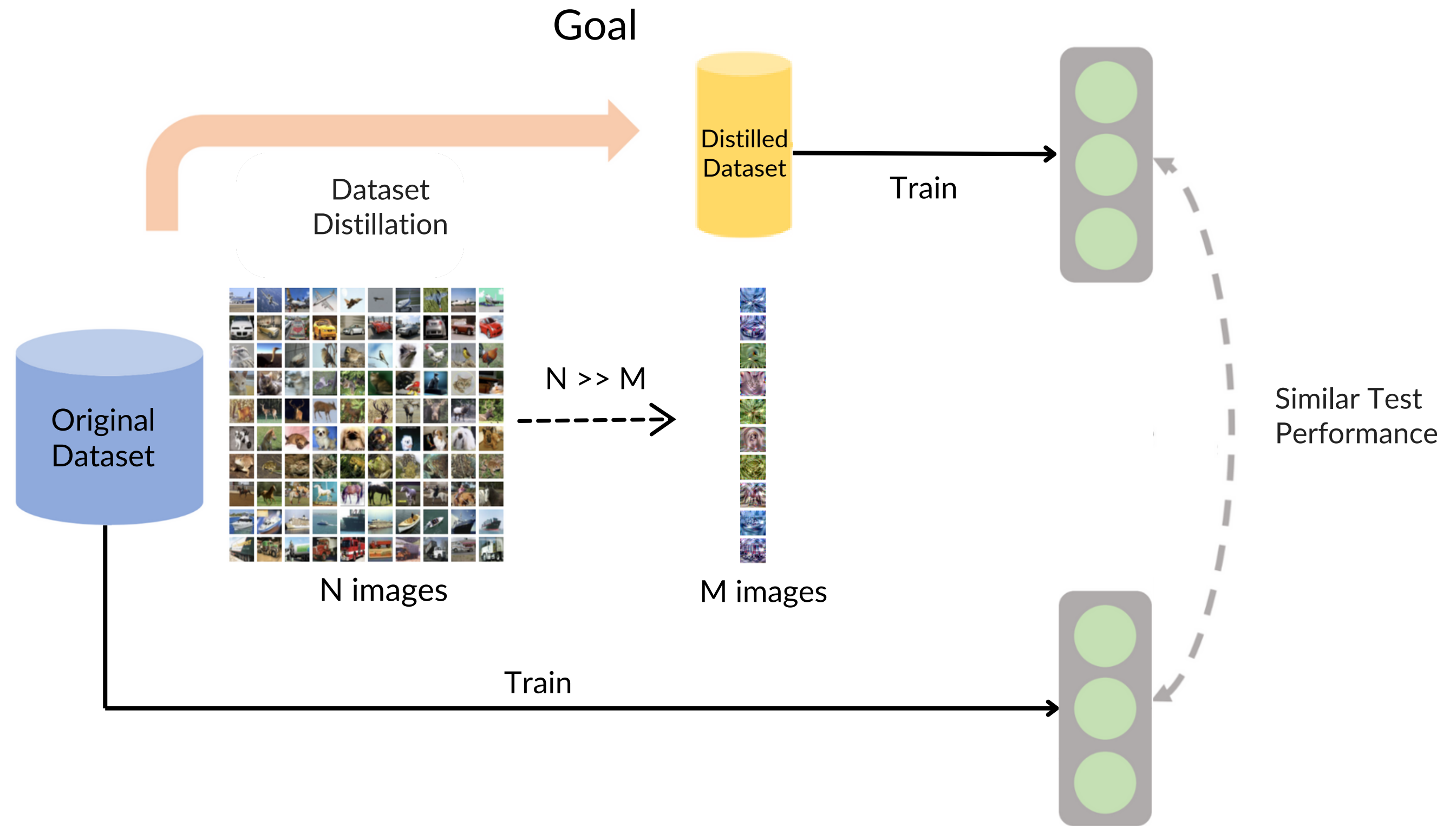
CIFAR10 examples



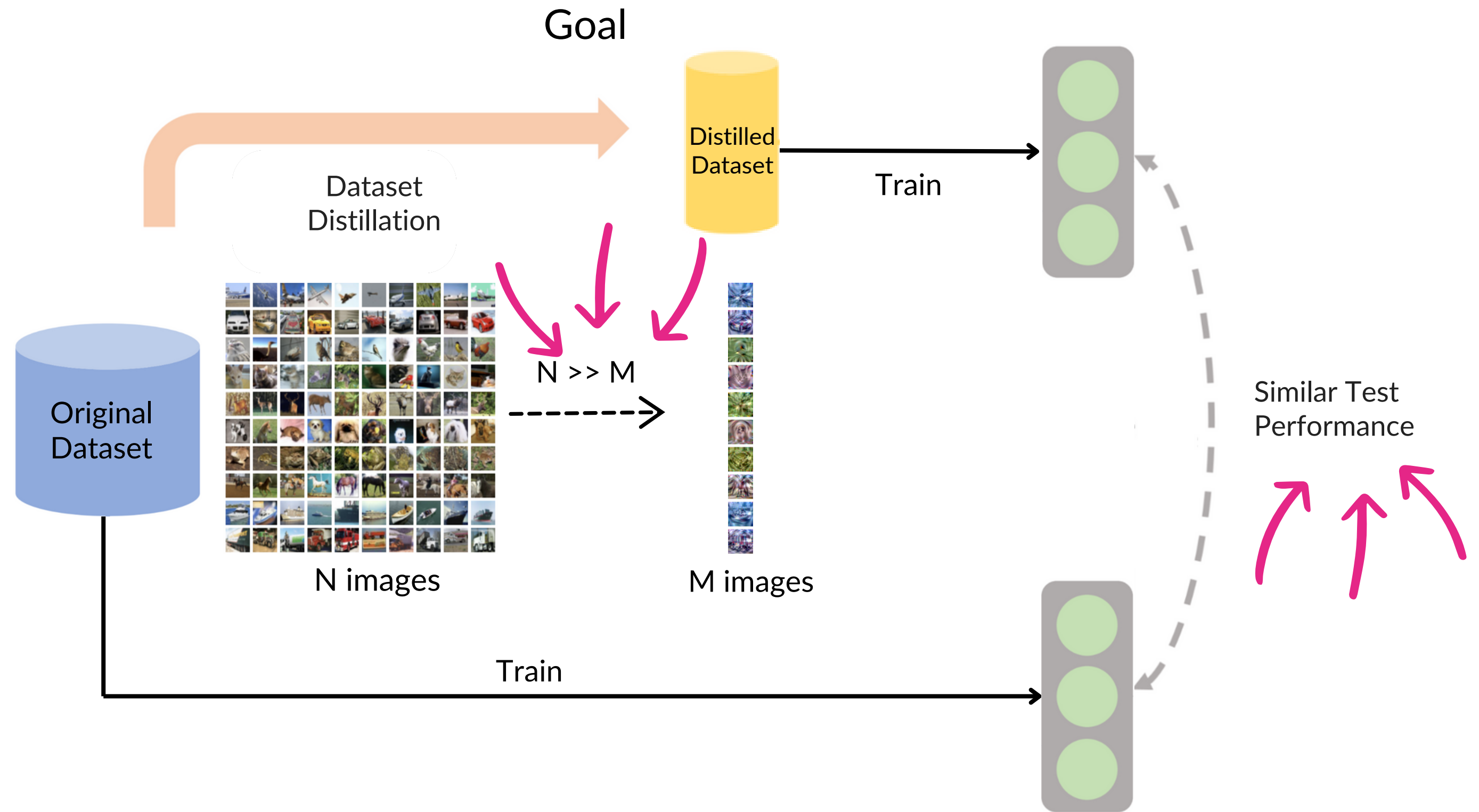
Distilled CIFAR10 (dog)

(10 images/class)

What is Dataset Distillation?



What is Dataset Distillation?





What is Dataset Distillation?

Defination

(Loose) Definition[1]:

Approaches that aim to synthesize tiny and high-fidelity data summaries which distill the most important knowledge from a given target dataset.



What is Dataset Distillation?

Defination

(Loose) Definition[1]:

Approaches that aim to synthesize tiny and high-fidelity data summaries which distill the most important knowledge from a given target dataset.

Such distilled summaries are optimized to serve as **effective** drop-in **replacements** of the original dataset for **efficient and accurate** data-usage applications.



What is Dataset Distillation?

Application

Neural
Architecture
Search

Federated
Learning

Continual
Learning

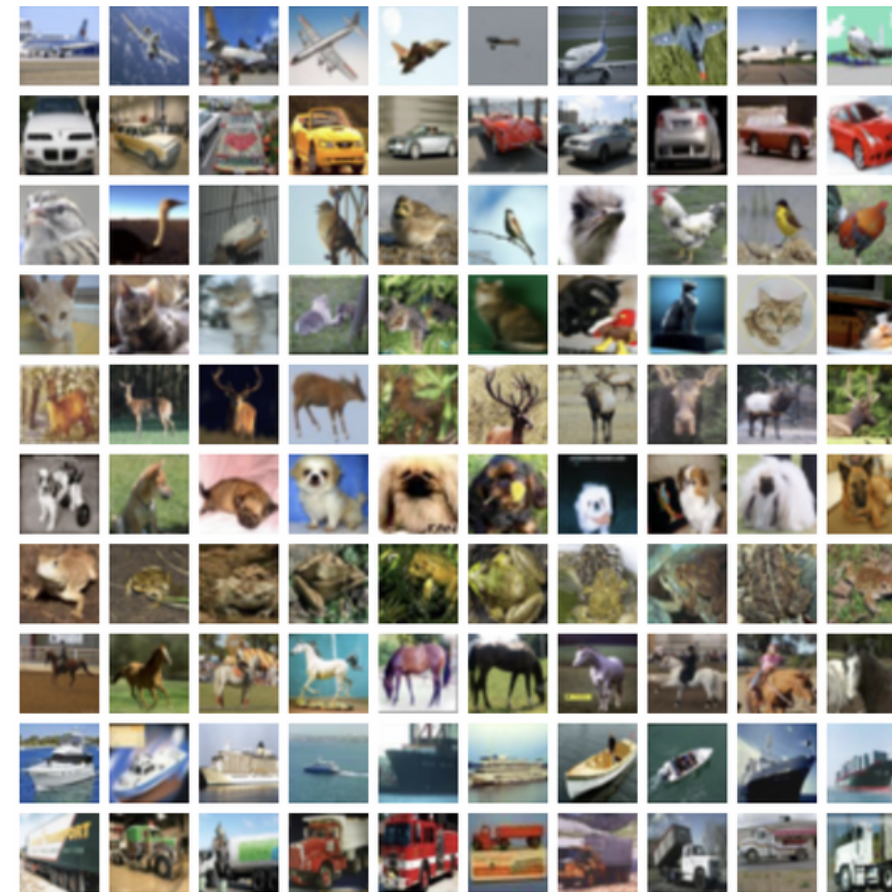
Differential
Privacy

 | Previous Work

Previous Work

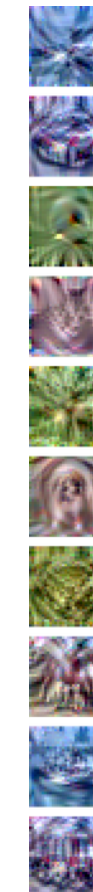
Image Classification

plane
car
bird
cat
deer
dog
frog
horse
ship
truck



CIFAR10
(5000 images/class)

$N \gg M$
----->

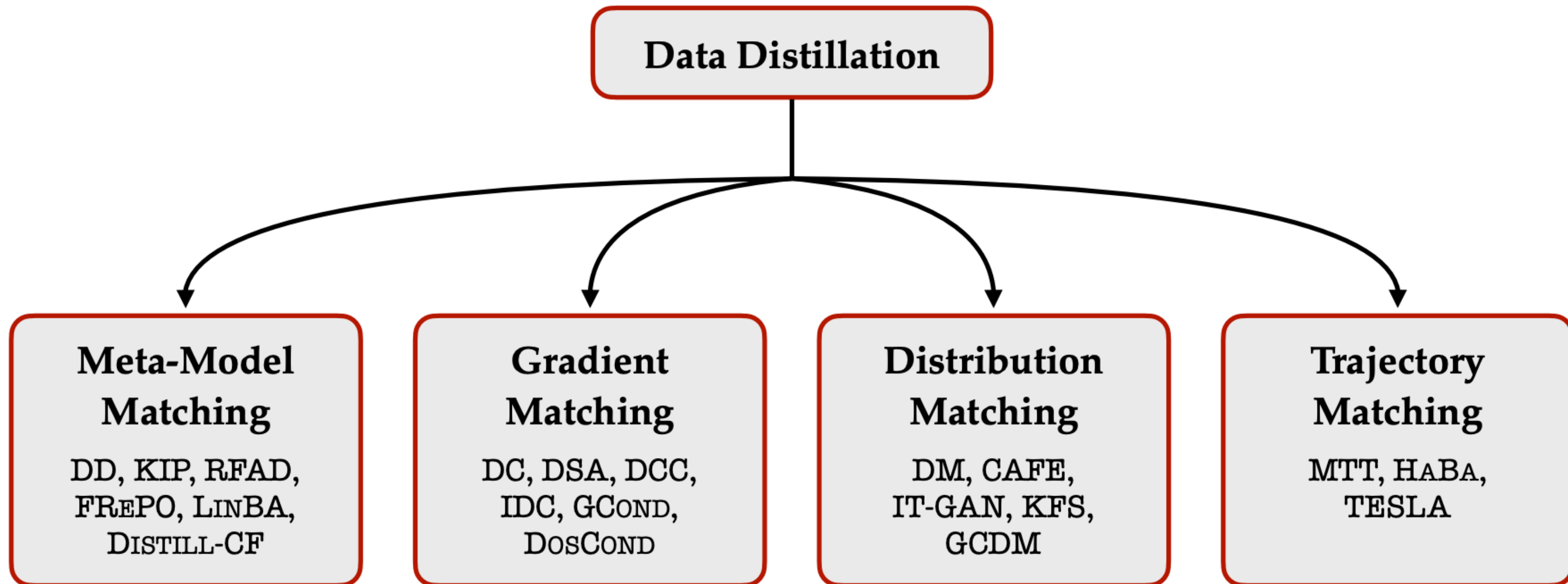


Distilled
CIFAR10
(1 images/class)



Previous Work

Taxonomy





Previous Work

Meta-Model Matching Framework

Dataset distillation can be formulated as a bi-level meta-learning problem

$$\begin{array}{ll} \text{Large Scale Dataset} & \mathbf{D} = \{(x_i, y_i)\}_{i=1}^N \\ \text{Small Synthetic Dataset} & \hat{\mathbf{D}} = \{(\hat{x}_j, \hat{y}_j)\}_{j=1}^M \end{array} \quad M \ll N$$



Previous Work

Meta-Model Matching Framework

Dataset distillation can be formulated as a bi-level meta-learning problem

<i>Large Scale Dataset</i>	$\mathbf{D} = \{(x_i, y_i)\}_{i=1}^N$	
<i>Small Synthetic Dataset</i>	$\hat{\mathbf{D}} = \{(\hat{x}_j, \hat{y}_j)\}_{j=1}^M$	$M \ll N$
<i>Model</i>	$f(\cdot; \theta)$	
<i>Generalization Loss</i>	$\ell(f(\hat{\mathbf{D}}; \theta), \mathbf{D})$	



Previous Work

Meta-Model Matching Framework

Dataset distillation can be formulated as a bi-level meta-learning problem

Inner level $f(\hat{\mathbf{D}}; \theta)$



Previous Work

Meta-Model Matching Framework

Dataset distillation can be formulated as a bi-level meta-learning problem

$$\begin{array}{ll} \text{Inner level} & f(\hat{\mathbf{D}}; \theta) \\ \text{Outer level} & \hat{\mathbf{D}}' = \arg \min_{\hat{\mathbf{D}}} F(\hat{\mathbf{D}}) \end{array}$$



Previous Work

Meta-Model Matching Framework

Dataset distillation can be formulated as a bi-level meta-learning problem

Inner level

$$f(\hat{\mathbf{D}}; \theta)$$

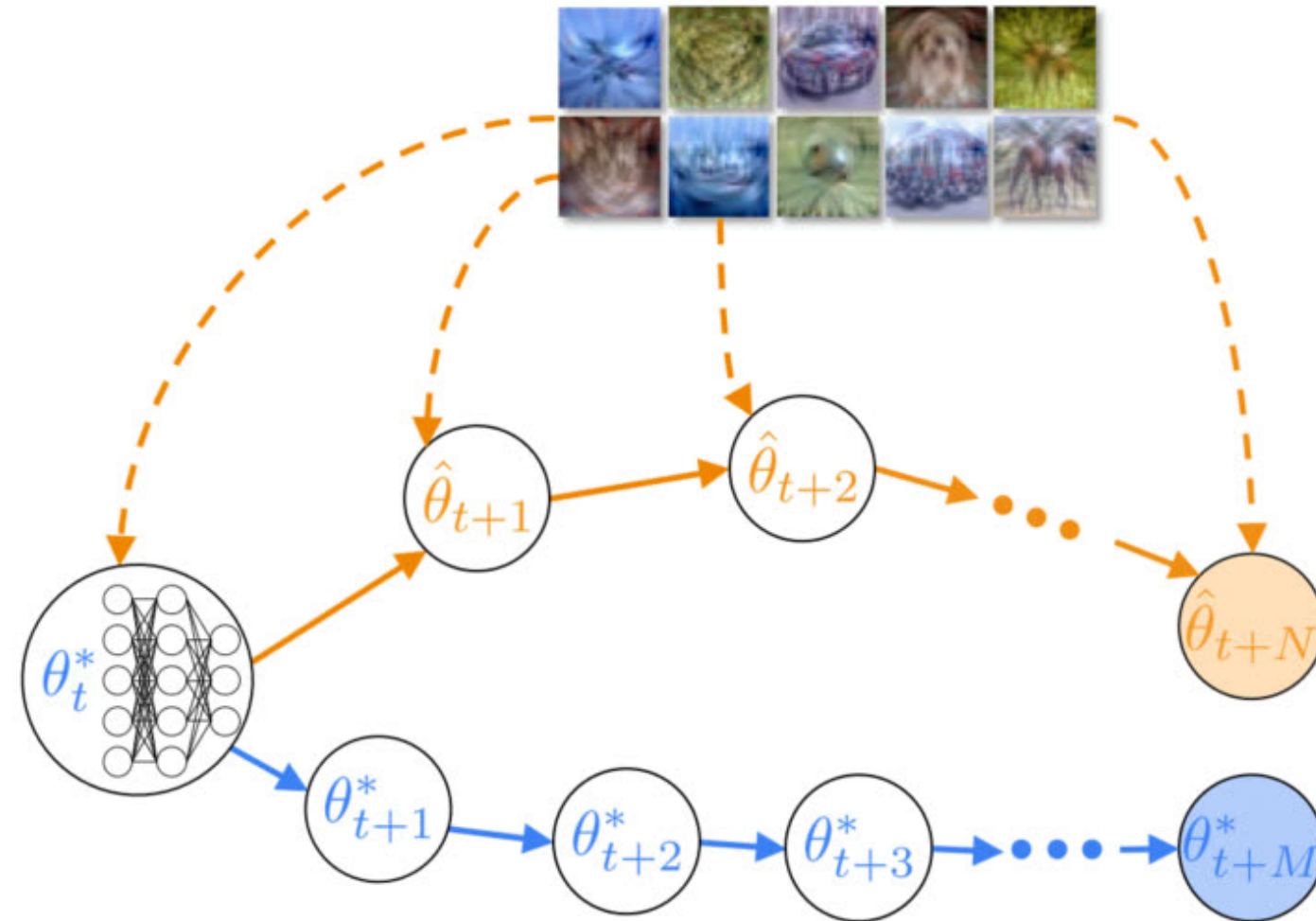
Outer level

$$\hat{\mathbf{D}}' = \arg \min_{\hat{\mathbf{D}}} F(\hat{\mathbf{D}})$$

$$F(\hat{\mathbf{D}}) = \mathbb{E}_{\theta \sim P(\theta)} \ell(f(\hat{\mathbf{D}}; \theta), \mathbf{D})$$

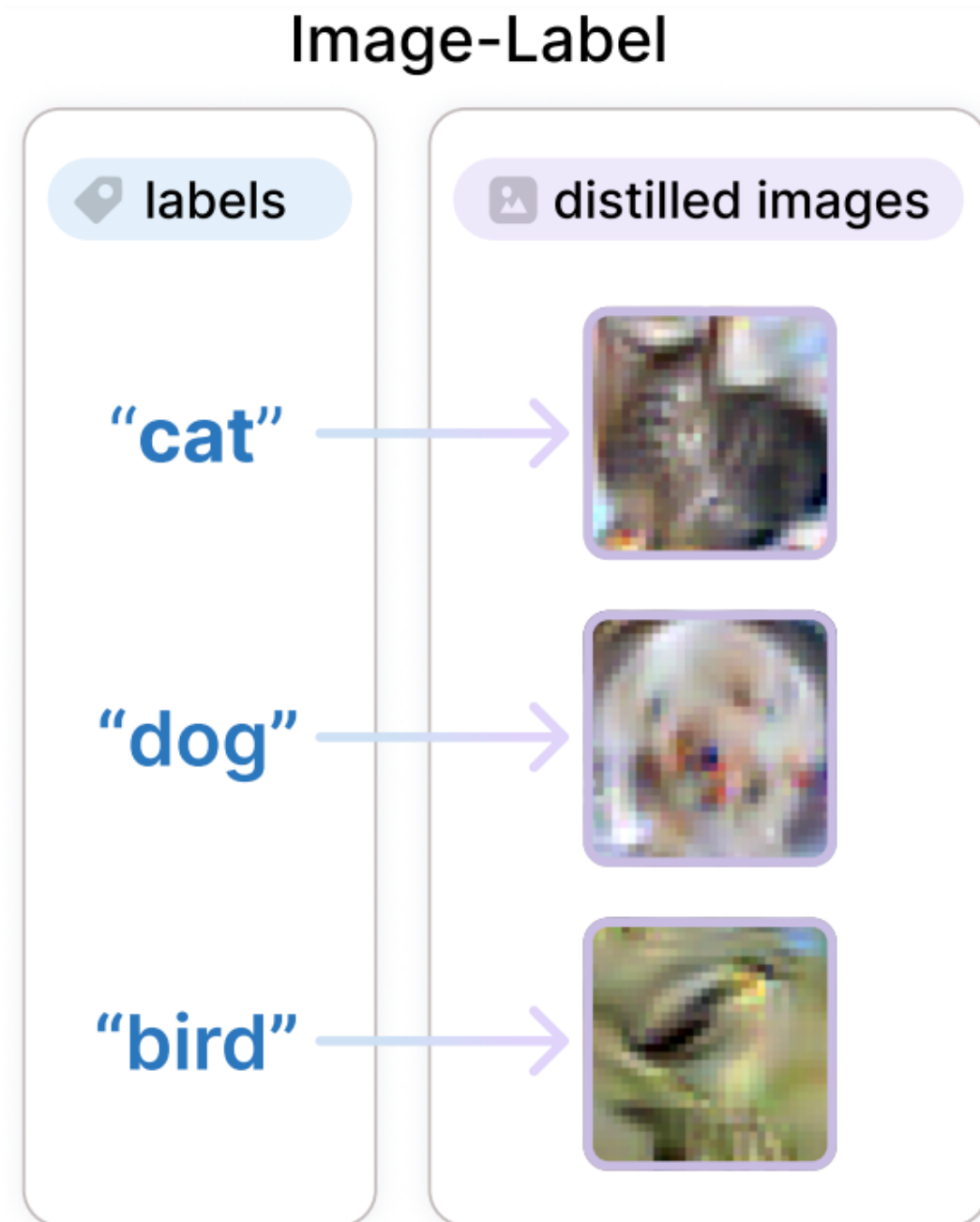
Previous Work

Trajectory Matching Framework

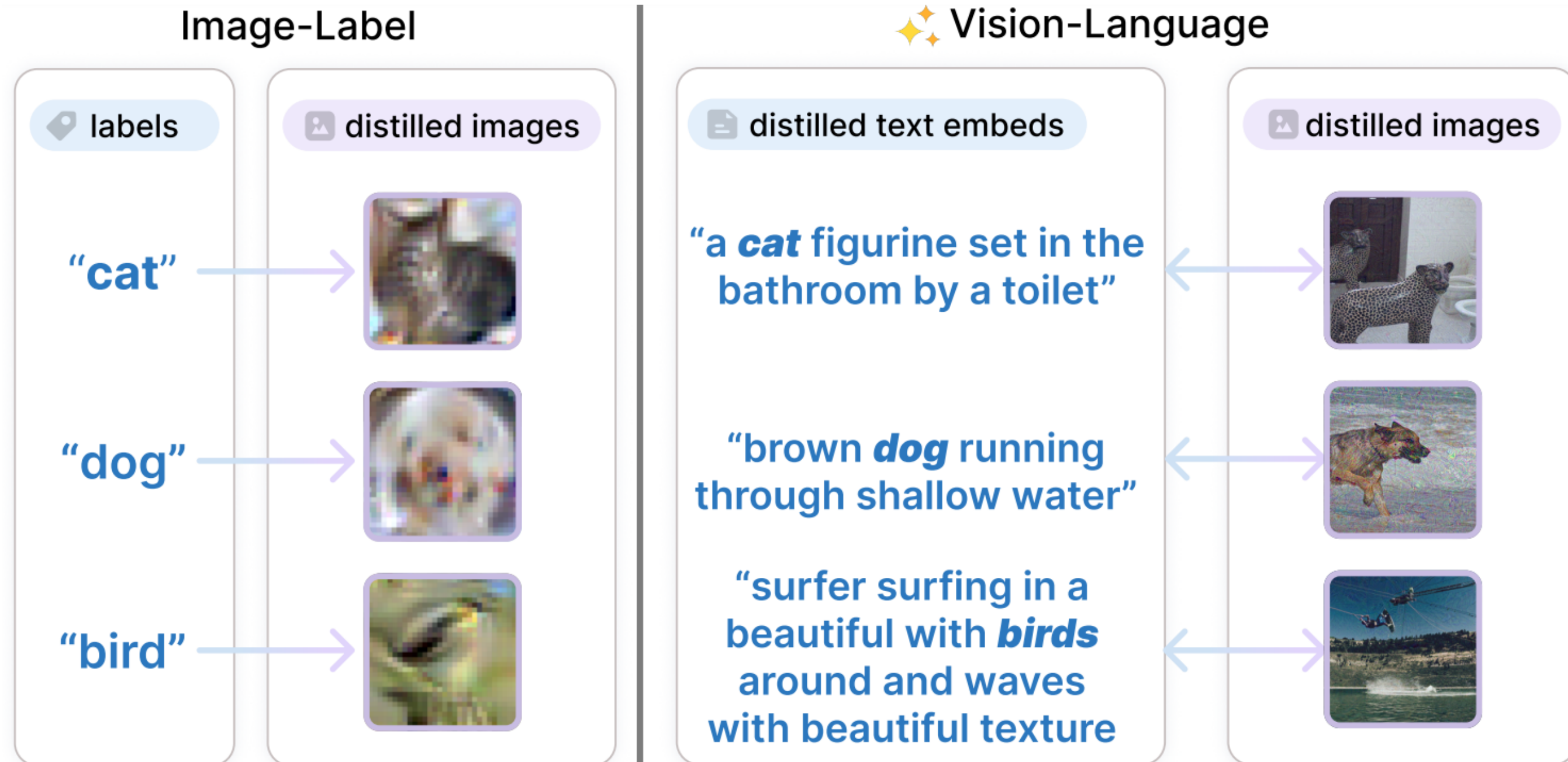


Student Trajectories are trained on Synthetic Data

Comparison



Comparison



Comparison

What's unique about Vision-Language Dataset?

image



text

A dog wearing a green sweater and fanny pack walks on a snow-covered field.

A small dog wearing a green sweater and a backpack walks through snow.

A dog wearing a green sweater and a backpack walking on snow.

A dog wearing a green sweater and backpack running in the snow.

A small dog wearing a sweater walking in the snow.

Comparison

What's unique about Vision-Language Dataset?

image



text

A **dog** wearing a green sweater and fanny pack walks on a snow-covered field.

A small **dog** wearing a green sweater and a backpack walks through snow.

A **dog** wearing a green sweater and a backpack walking on snow.

A **dog** wearing a green sweater and backpack running in the snow.

A small **dog** wearing a sweater walking in the snow.

Comparison

What's unique about Vision-Language Dataset?

image



text

A **dog** wearing a green sweater and fanny pack **walks** on a snow-covered field.

A small **dog** wearing a green sweater and a backpack **walks** through snow.

A **dog** wearing a green sweater and a backpack **walking** on snow.

A **dog** wearing a green sweater and backpack **running** in the snow.

A small **dog** wearing a sweater **walking** in the snow.

Comparison

What's unique about Vision-Language Dataset?

image



text

A **dog** wearing a **green sweater** and fanny pack walks on a snow-covered field.

A small dog wearing a **green sweater** and a backpack walks through snow.

A **dog** wearing a **green sweater** and a backpack walking on snow.

A **dog** wearing a **green sweater** and backpack running in the snow.

A small **dog** wearing a **sweater** walking in the snow.

Comparison

What's unique about Vision-Language Dataset?

image



text

A **dog** wearing a **green sweater** and **fanny pack** walks on a snow-covered field.

A small dog wearing a **green sweater** and a **backpack** walks through snow.

A **dog** wearing a **green sweater** and a **backpack** walking on snow.

A **dog** wearing a **green sweater** and **backpack** running in the snow.

A small **dog** wearing a **sweater** walking in the snow.

Comparison

What's unique about Vision-Language Dataset?

image



text

A **dog** wearing a **green sweater** and **fanny pack** walks on a **snow-covered field**.

A small **dog** wearing a **green sweater** and a **backpack** walks through **snow**.

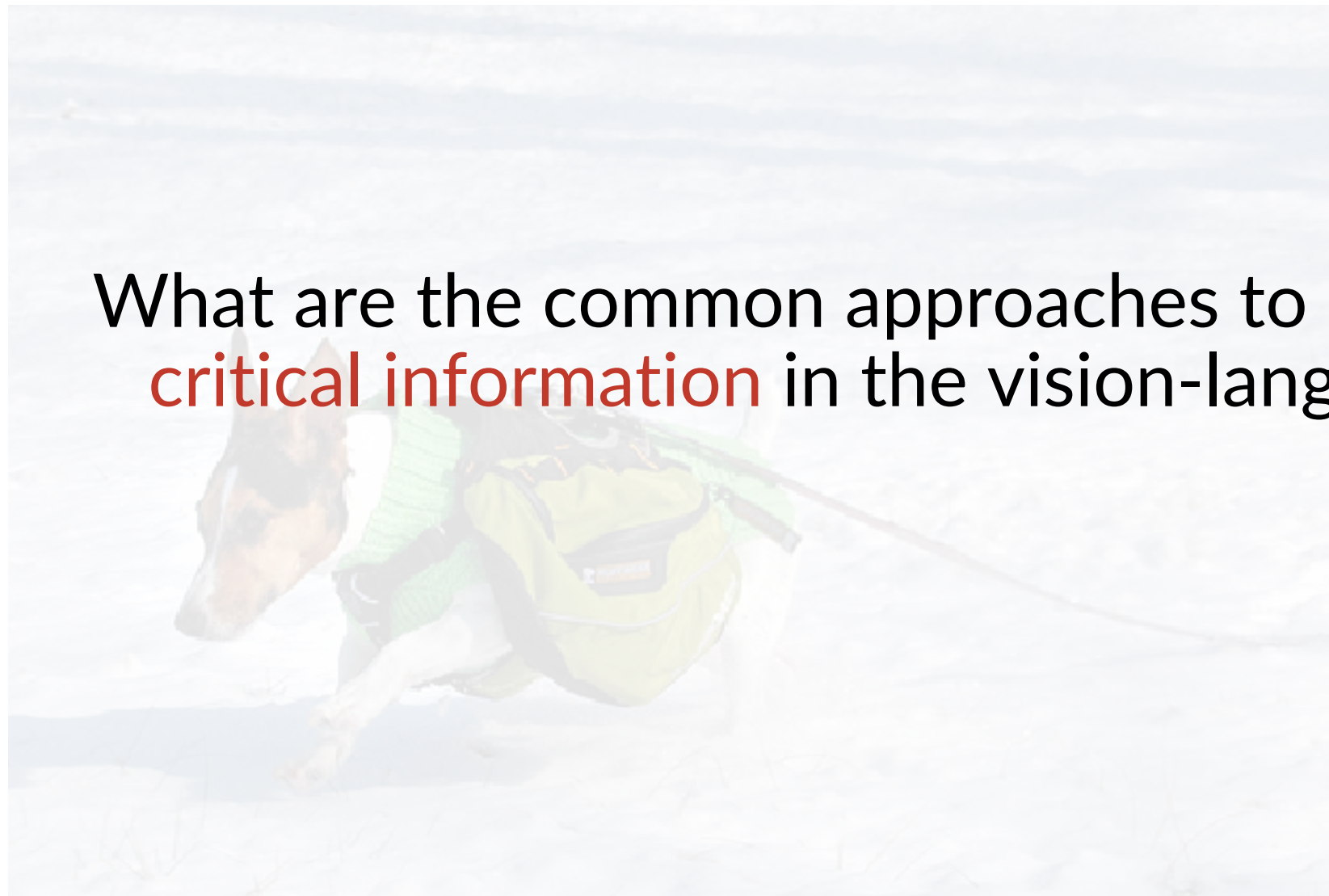
A **dog** wearing a **green sweater** and a **backpack** walking on **snow**.

A **dog** wearing a **green sweater** and **backpack** running in the **snow**.

A small **dog** wearing a **sweater** walking in the **snow**.

Comparison

image



text

A dog wearing a green sweater and fanny pack walks on a snow-covered field.

A small dog wearing a green sweater and fanny pack walks through snow.

A dog wearing a green sweater and backpack walking on snow.

A dog wearing a green sweater and backpack running in the snow.

A small dog wearing a sweater walking in the snow.

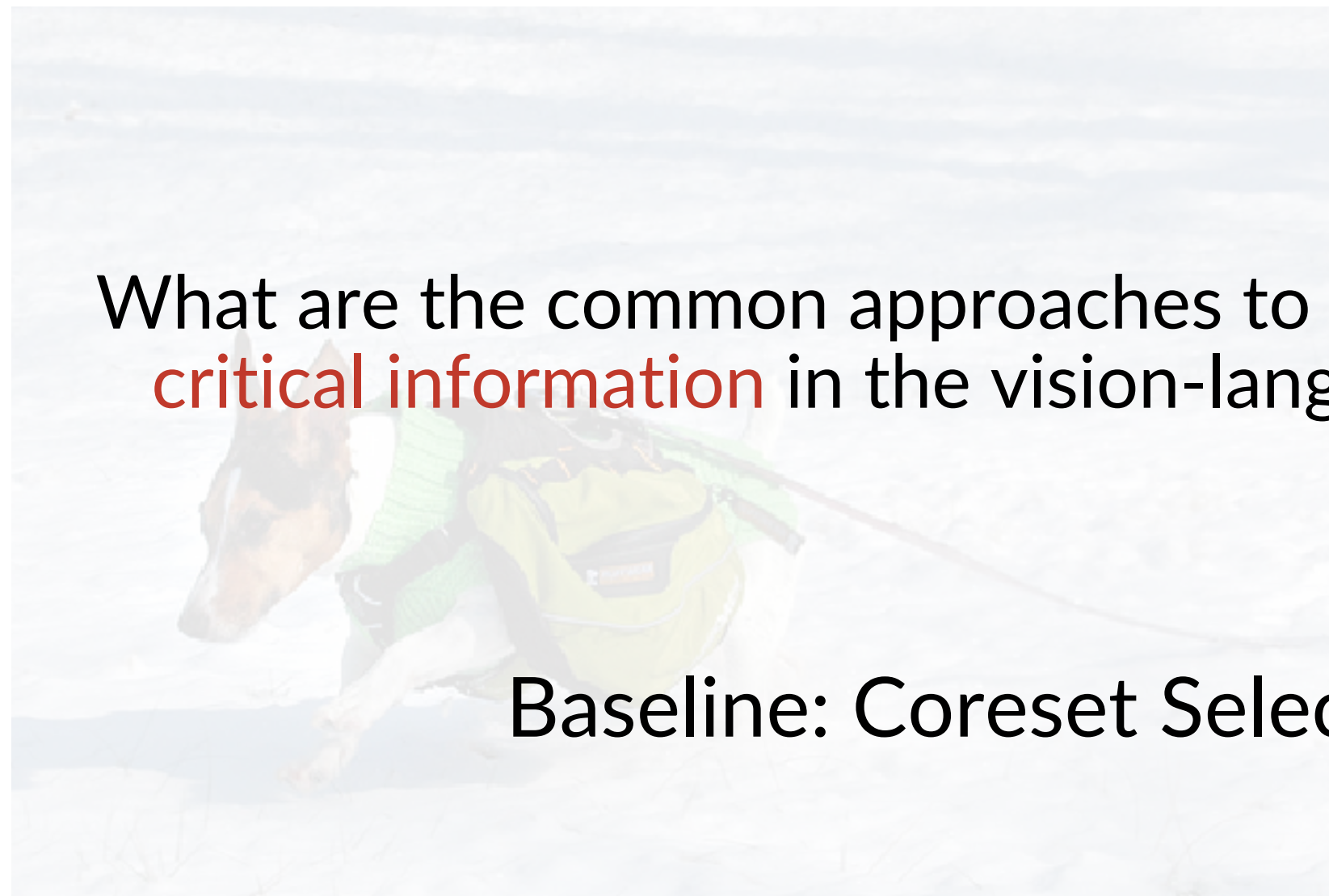


What are the common approaches to select the most **critical information** in the vision-language dataset?

Comparison

image

text



What are the common approaches to select the most **critical information** in the vision-language dataset?

Baseline: Coreset Selection

A dog wearing a green sweater and **fanny pack** walks on a snow-covered field.

A small dog wearing a green sweater and **fanny pack** walks through snow.

A dog wearing a green sweater and **backpack** walking on snow.

A dog wearing a green sweater and **backpack** running in the snow.

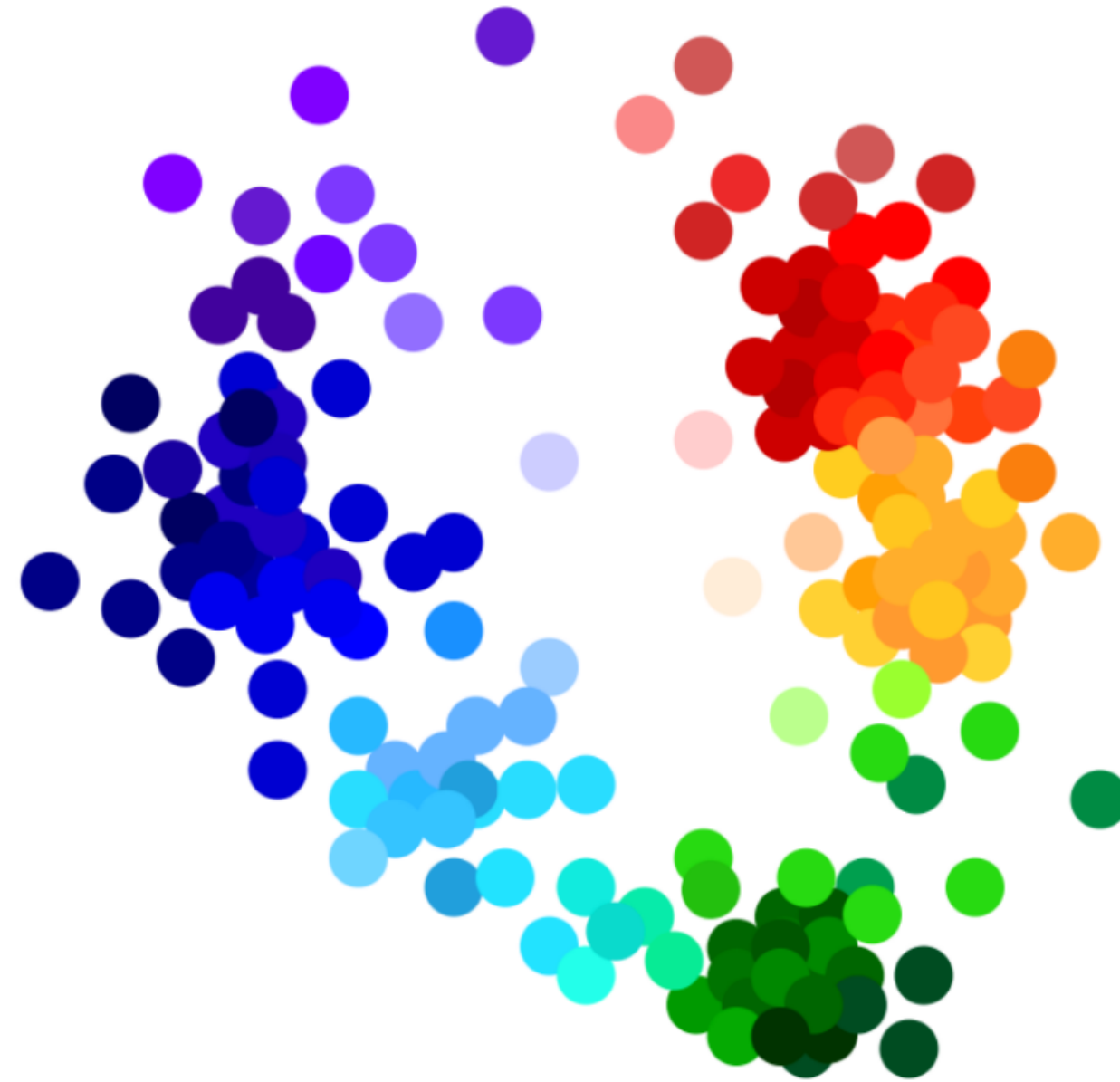
A small dog wearing a sweater walking in the snow.





Baseline

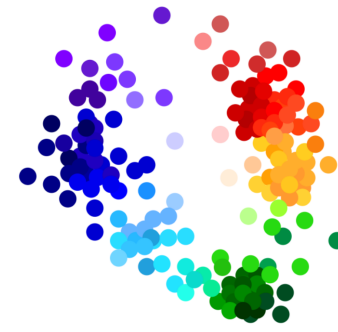
Coreset Selection





Baseline

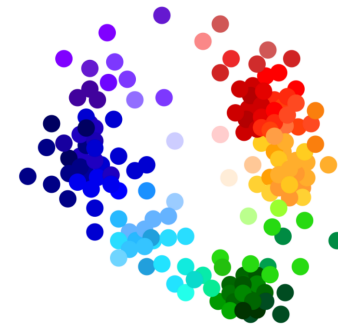
Coreset Selection



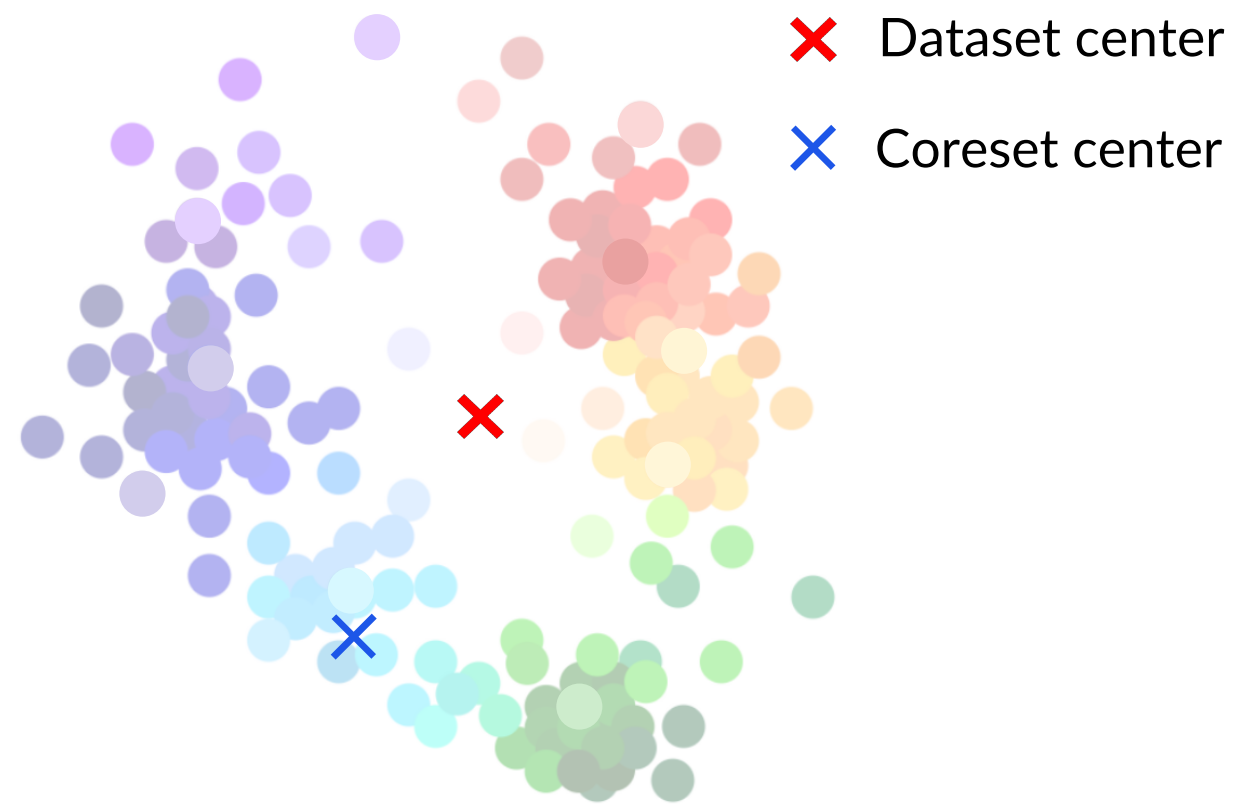
Herding

Baseline

Coreset Selection

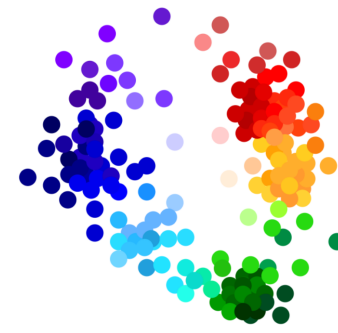


Herding

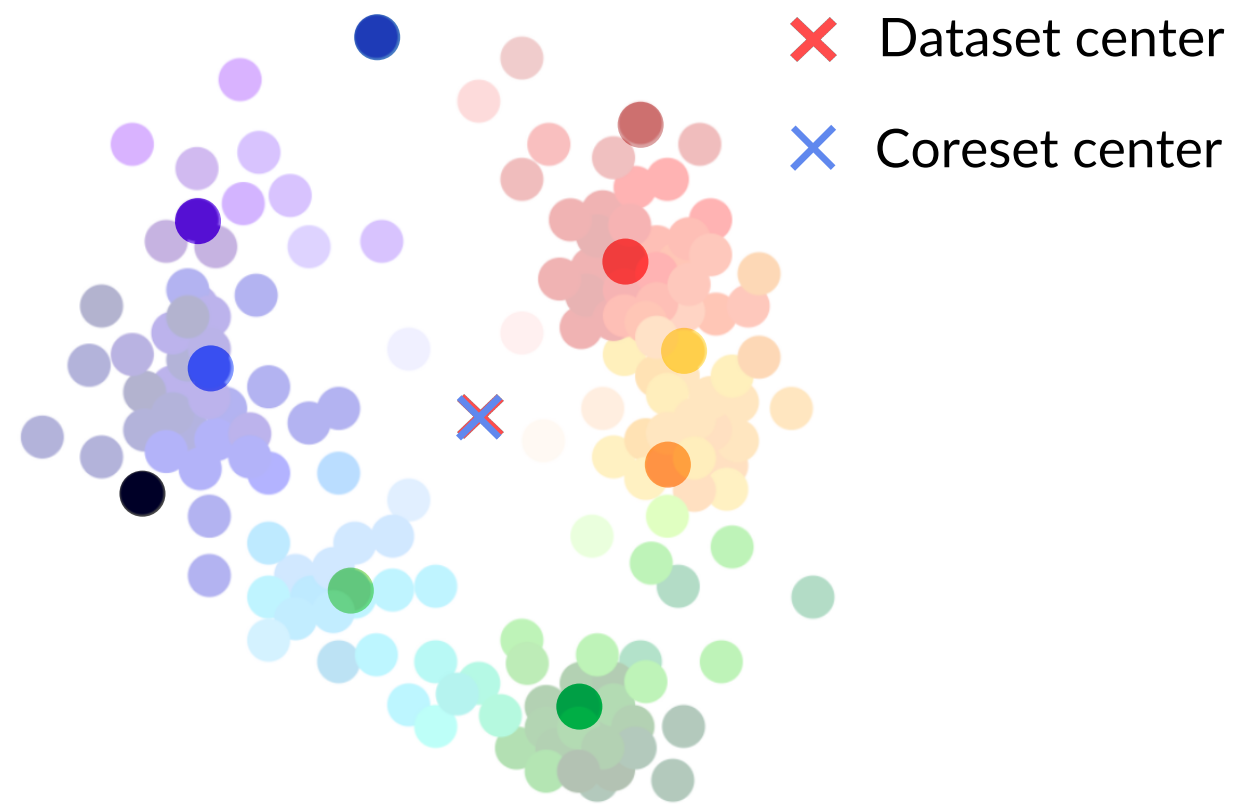


Baseline

Coreset Selection

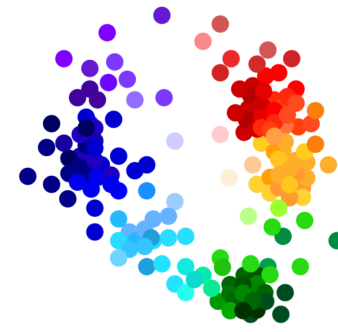


Herding

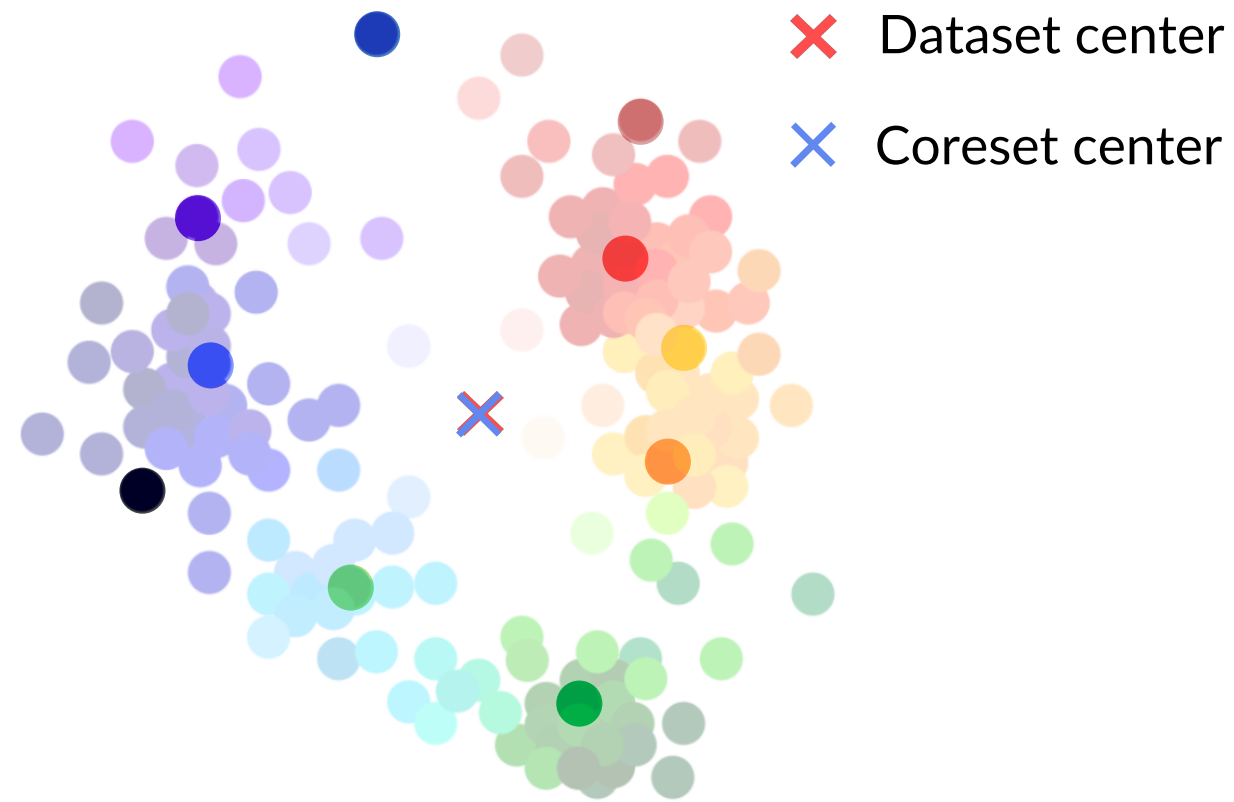


Baseline

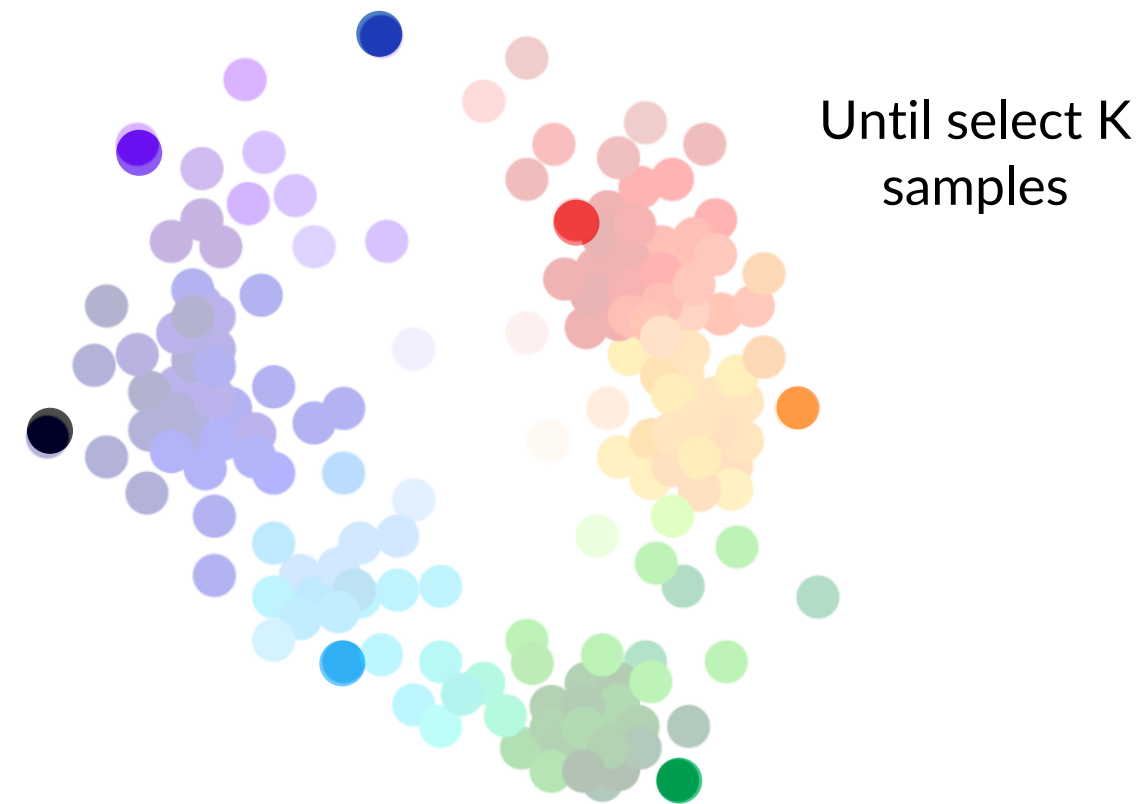
Coreset Selection



Herding

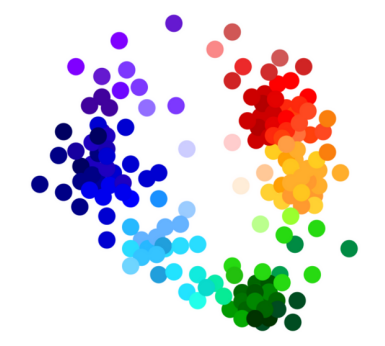


K-center

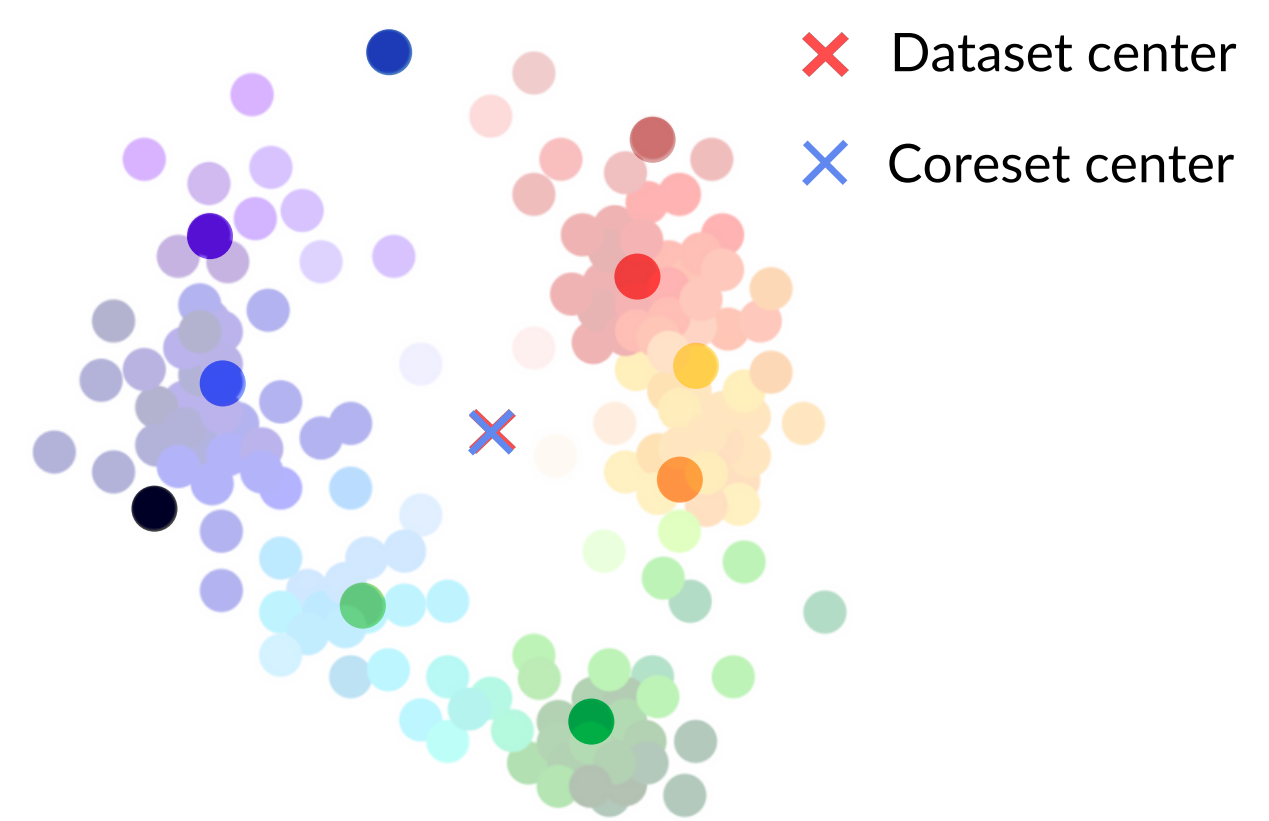


Baseline

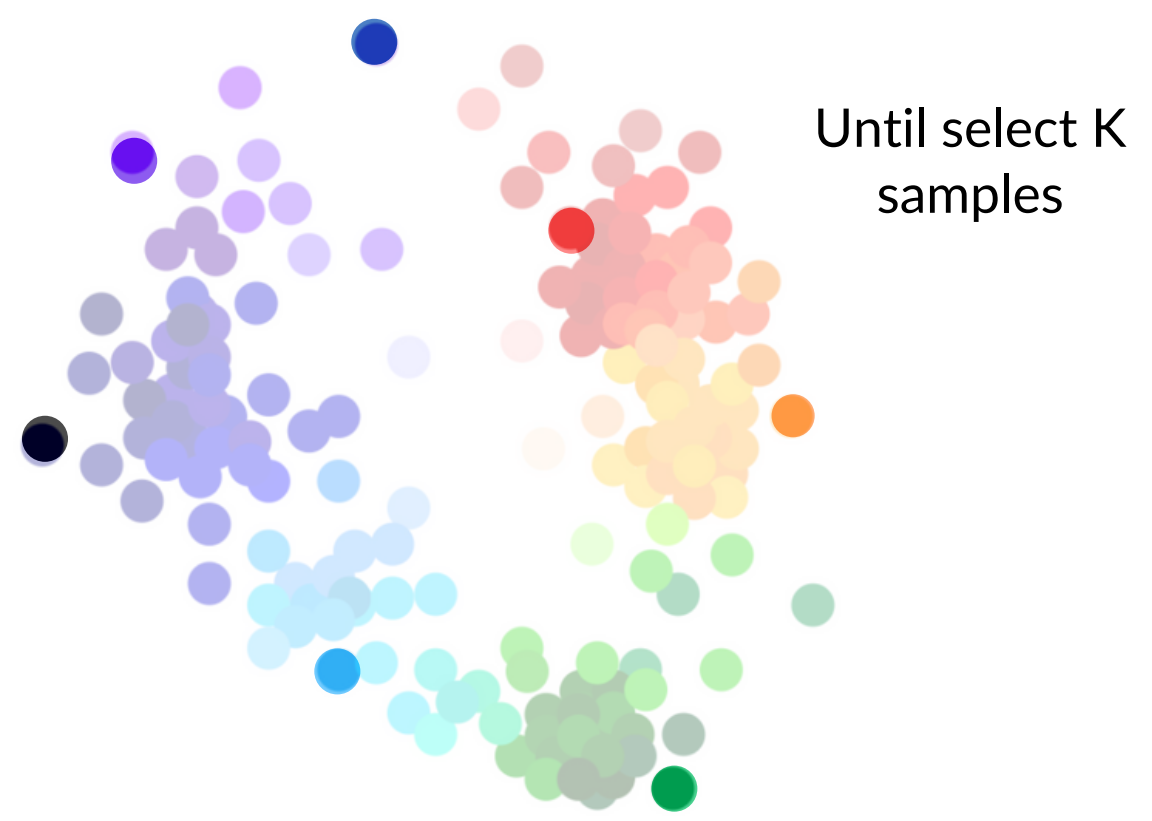
Coreset Selection



Herding



K-center



Until select K samples

Forgetting

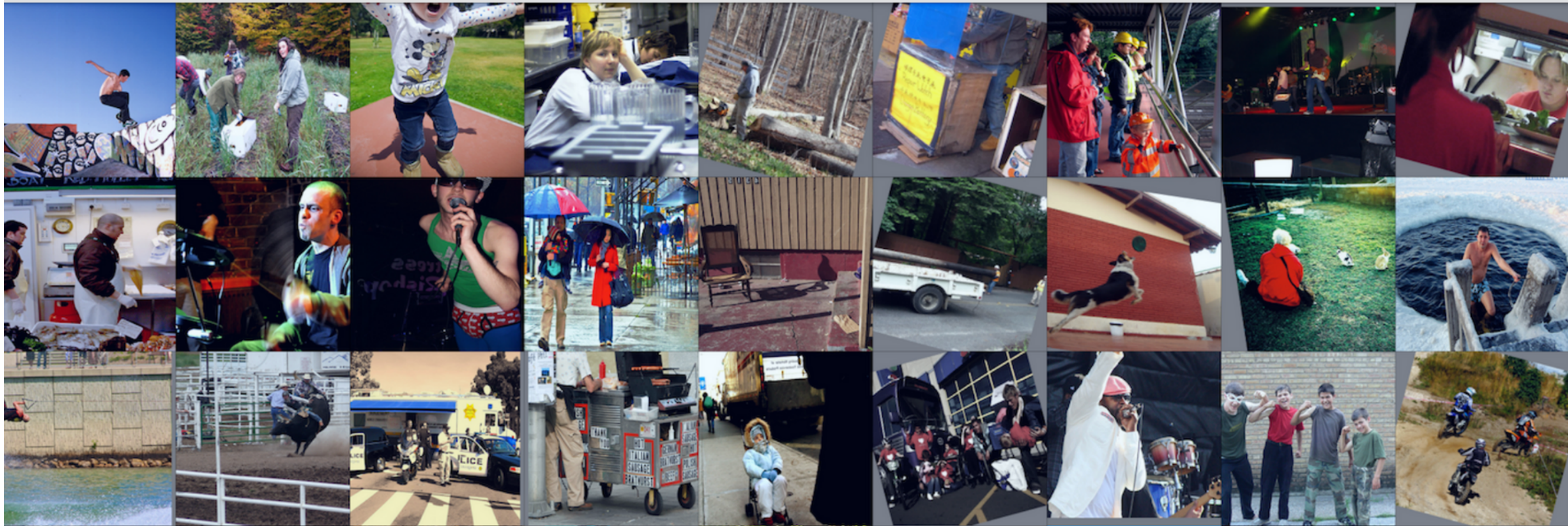
Forgetting event:
When the model correctly predicts the example in one epoch but fails in the next.

Track these forgetting events during training and identify the ones with the least forgetting events.

Our Task

Large Vision-Language Dataset

N pairs of image + text



Our Task

Large Vision-Language Dataset

N pairs of image + text



30K pairs

Flickr30K

Our Task

Large Vision-Language Dataset

N pairs of image + text

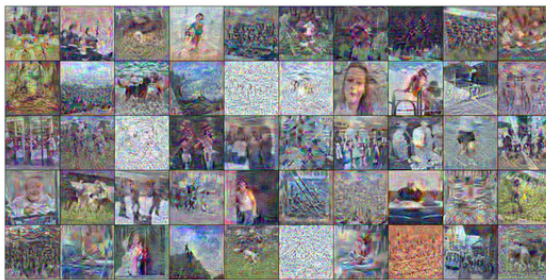


30K pairs

Flickr30K

$N \gg M$

Small Synthetic Dataset



M pairs of image + text

Our Task

Large Vision-Language Dataset

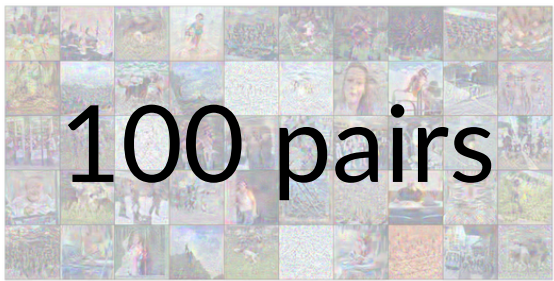
N pairs of image + text



Flickr30K

$N \gg M$

Small Synthetic Dataset

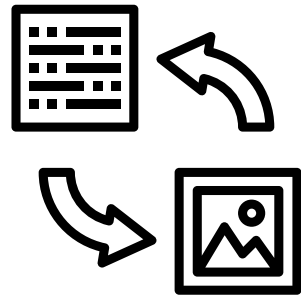


M pairs of image + text

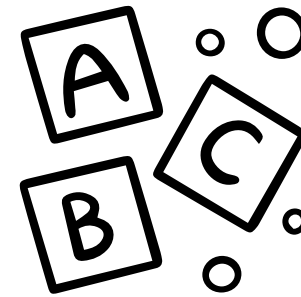
Challenges

Vision-Language Dataset Distillation

Complex Cross-Modal Relationships



Discrete text optimization issue



Heavy computational cost

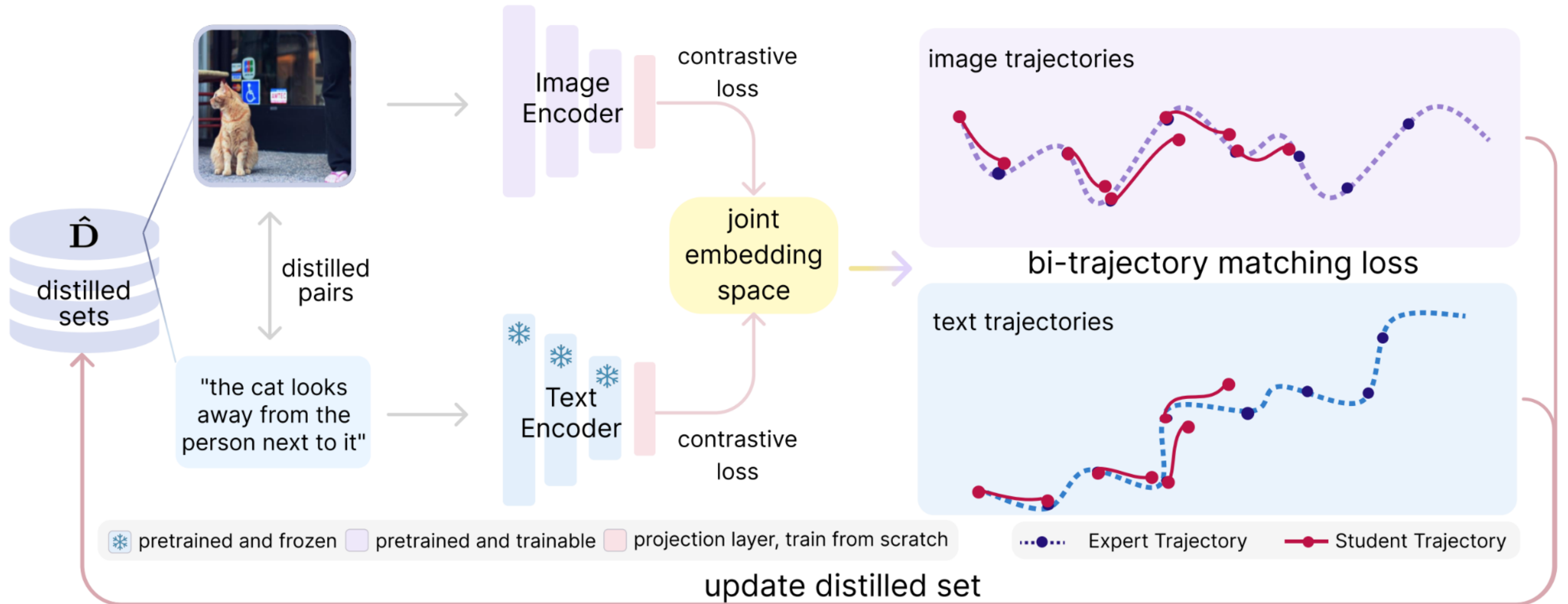




| Method

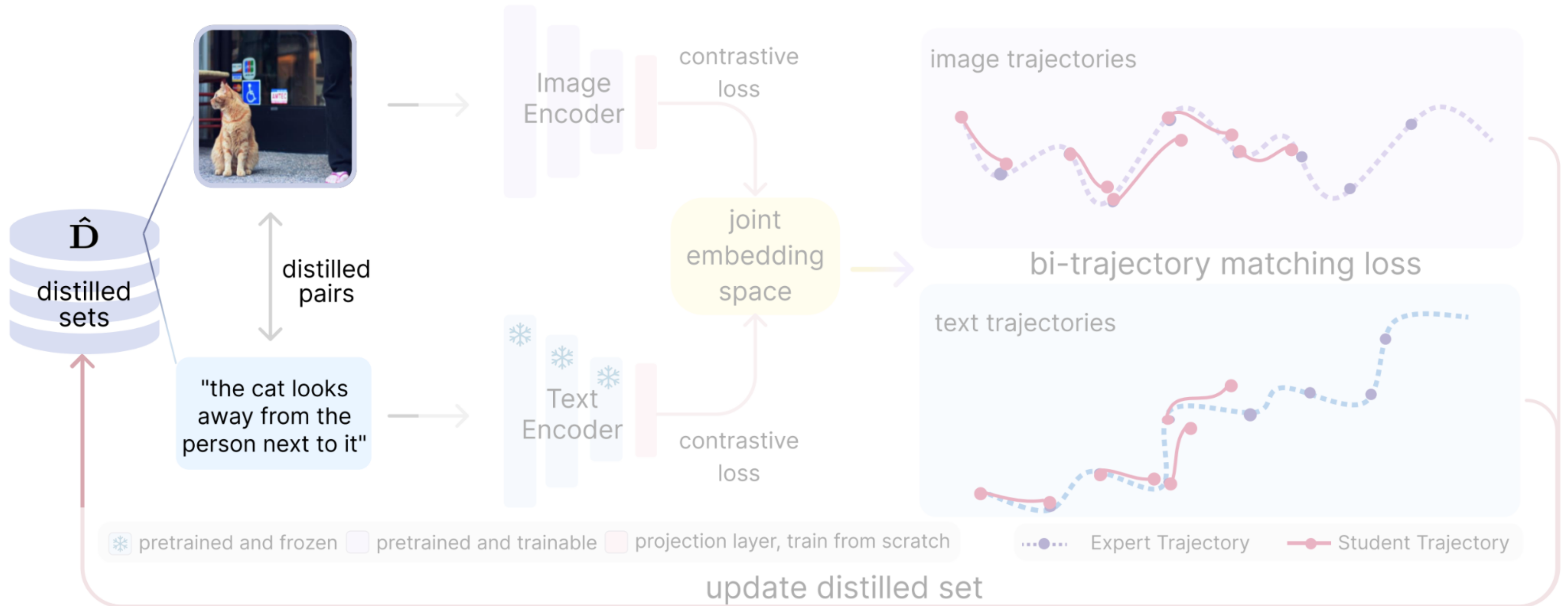
Method

Pipeline Overview



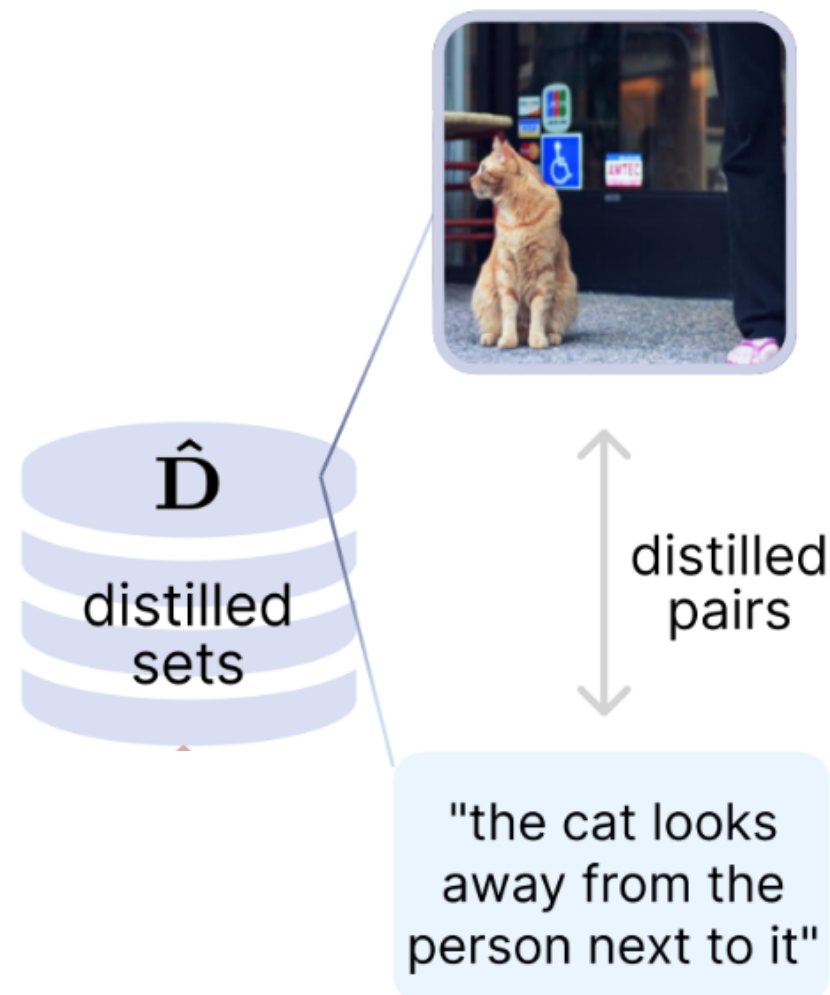
Method

Problem Formulation



Method

Problem Formulation



Large Scale Dataset

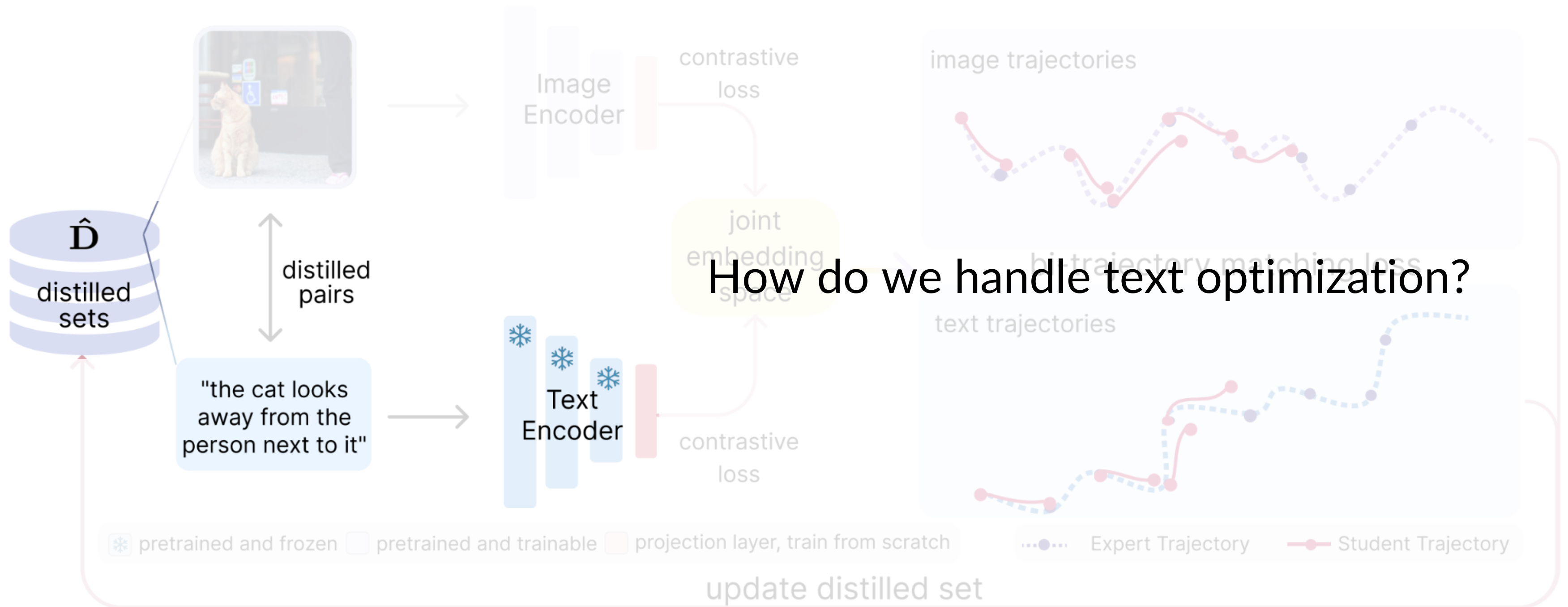
$$\mathbf{D} = \{(x_i, y_i)\}_{i=1}^N$$

Small Synthetic Dataset

$$\hat{\mathbf{D}} = \{(\hat{x}_j, \hat{y}_j)\}_{j=1}^M$$

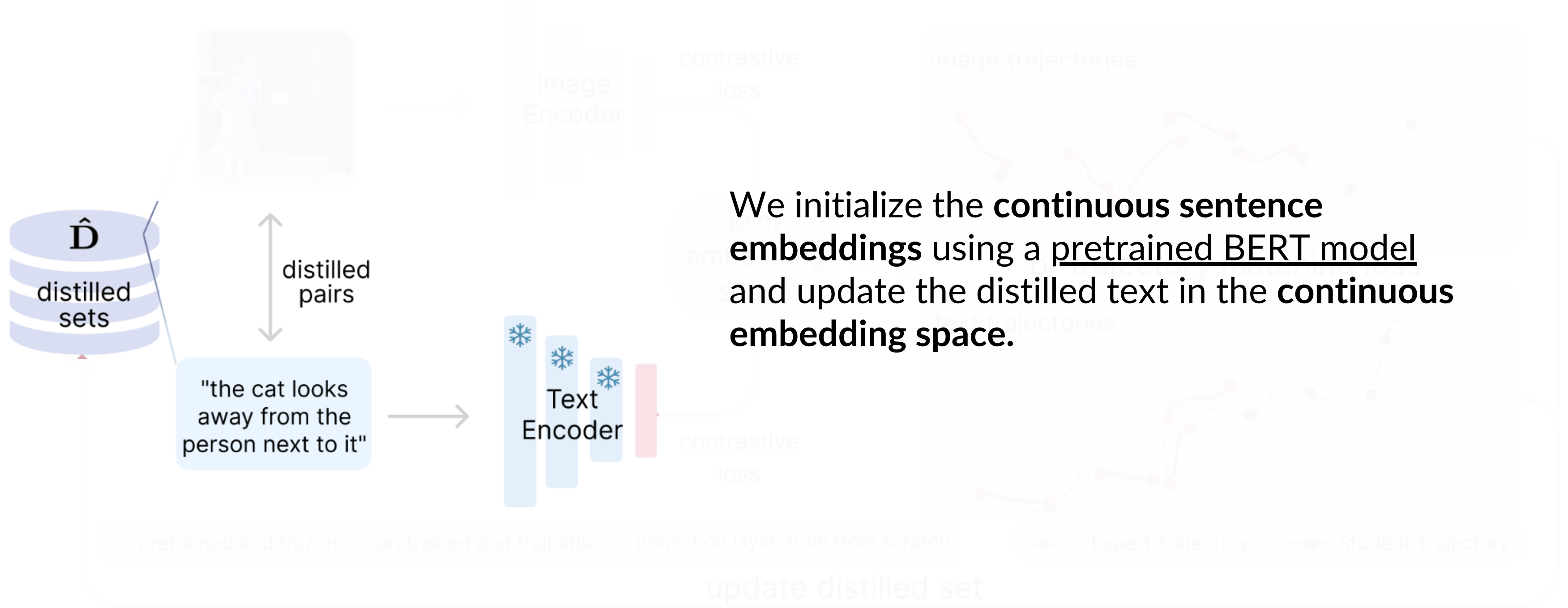
$$M \ll N$$

Method



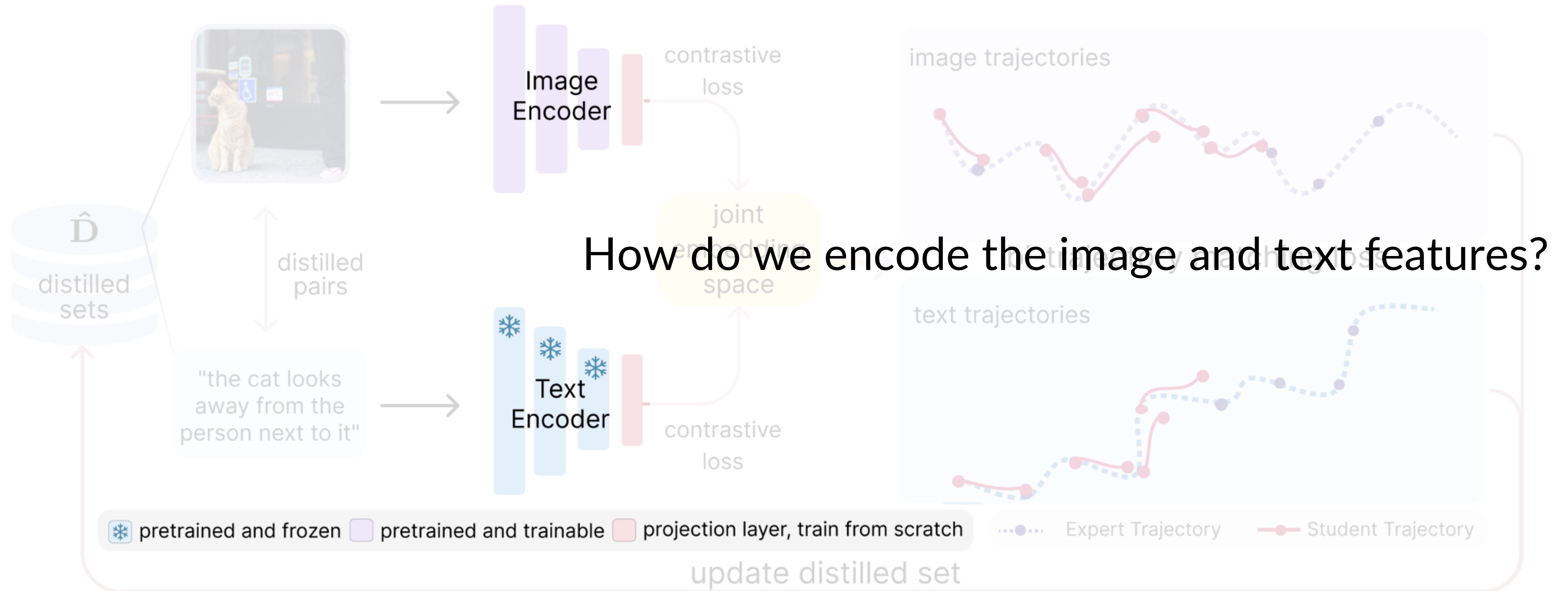
Method

How do we handling text optimization?

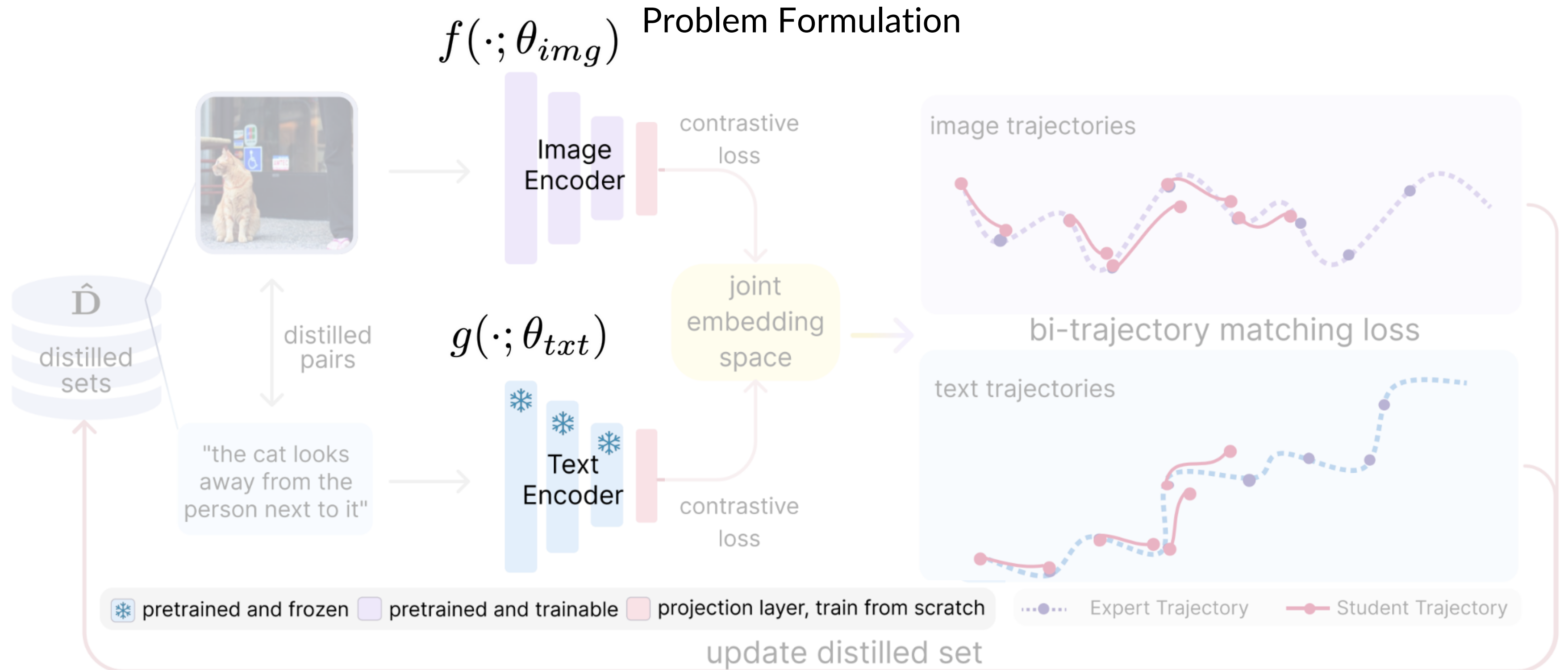


Method

Problem Formulation

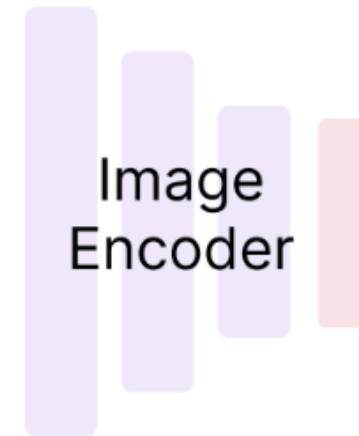


Method

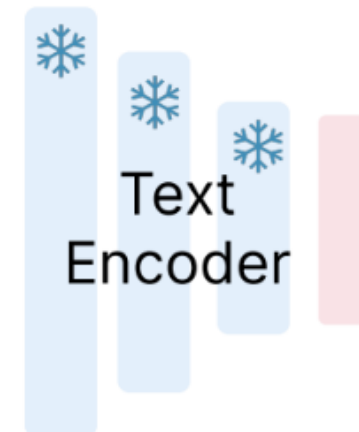


Method

$f(\cdot; \theta_{img})$ Problem Formulation



$g(\cdot; \theta_{txt})$

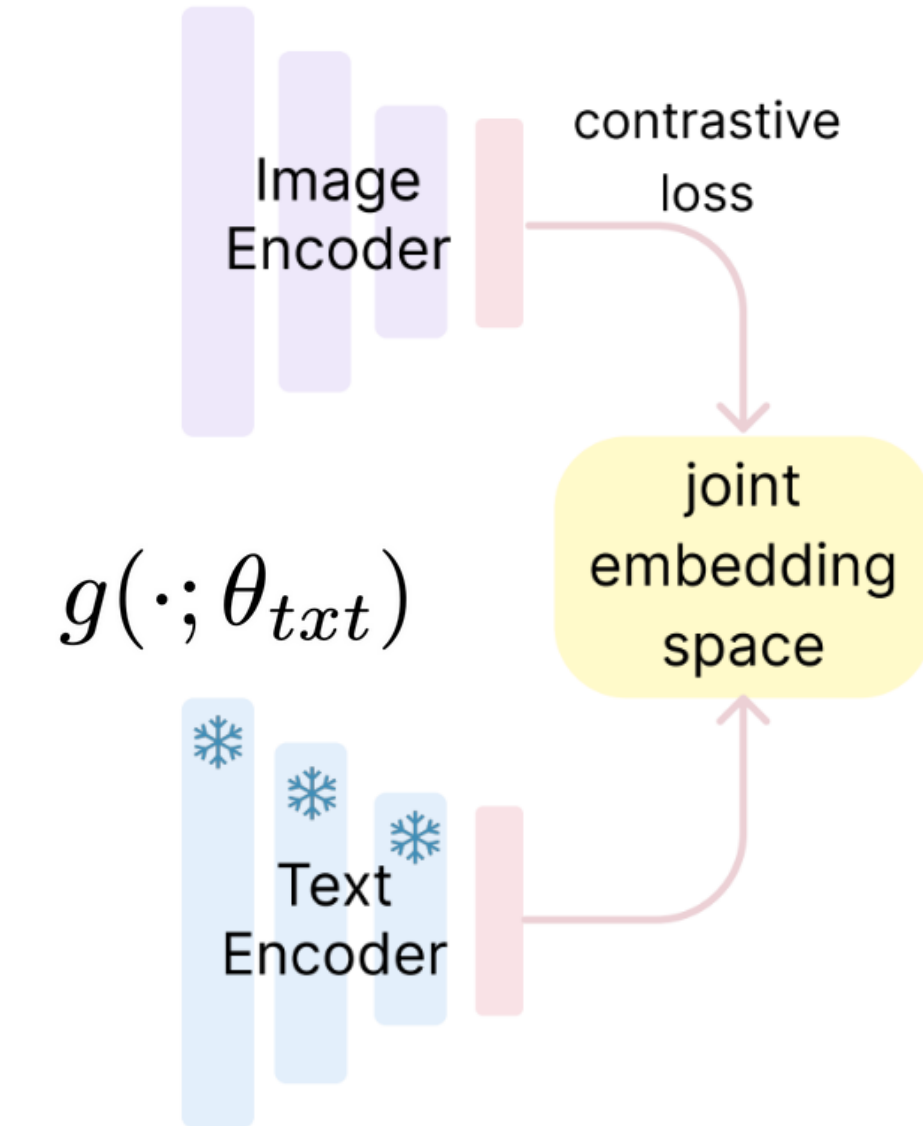


$$\theta^* \approx \arg \min_{\theta} \frac{1}{N} \sum_{i=1}^N \ell (f(x_i; \theta_{img}), g(y_i; \theta_{txt}))$$

pretrained and frozen pretrained and trainable projection layer, train from scratch

Method

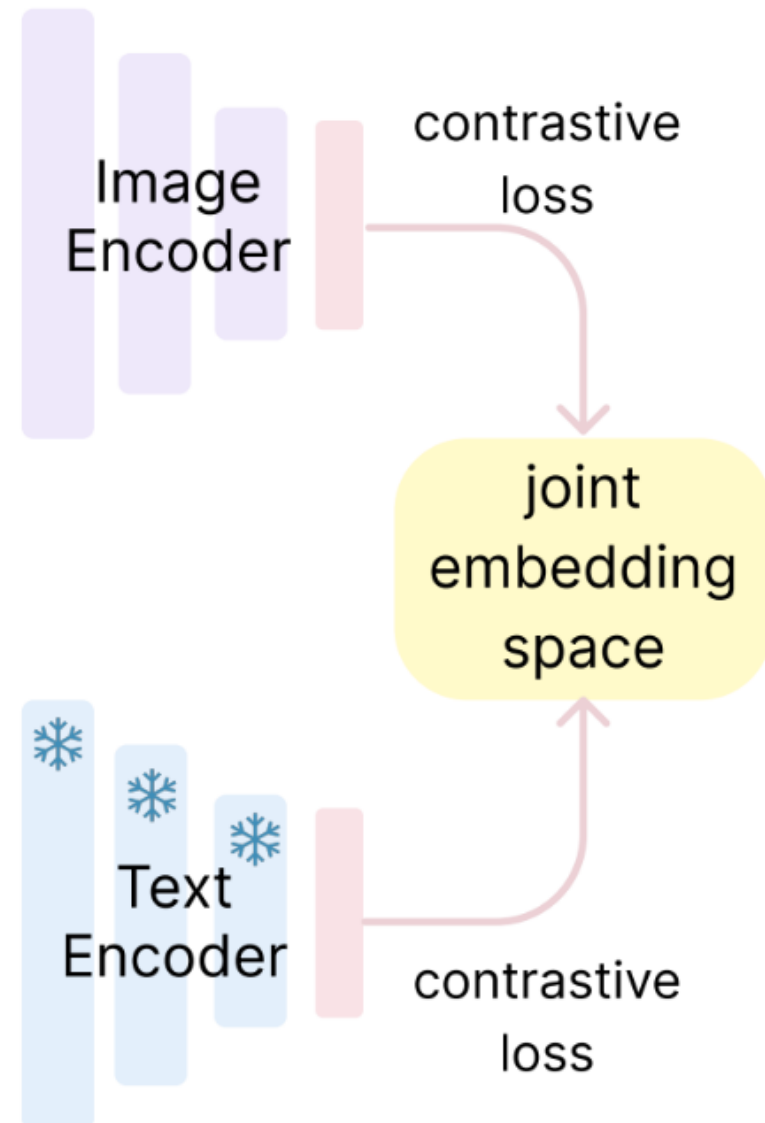
$f(\cdot; \theta_{img})$ Problem Formulation



pretrained and frozen pretrained and trainable projection layer, train from scratch

Method

Problem Formulation



$$\alpha(x, y) = \frac{\langle f(x; \theta_{img}), g(y; \theta_{txt}) \rangle}{\|f(x; \theta_{img})\| \|g(y; \theta_{txt})\|}$$

$$\ell_{contrastive} = -\frac{1}{2n} \sum_{(x,y) \text{ in batch}} \left(\log \frac{\exp(\alpha(x, y))}{\sum_{y' \neq y} \exp(\alpha(x, y'))} + \log \frac{\exp(\alpha(x, y))}{\sum_{x' \neq x} \exp(\alpha(x', y))} \right)$$

❄ pretrained and frozen
 ▭ pretrained and trainable
 ▭ projection layer, train from scratch



Method

Bi-Trajectory Guided Vision-Language Co-Distillation

Concretely, our approach consists of two stages:

Stage 1 Expert training

Training multiple models for T epochs on the full dataset \mathbf{D}

Obtaining the expert training trajectories $\{\tau^*\}$, with each trajectory $\tau^* = \{\theta_t^*\}_{t=0}^T$

For our multimodal setting, the models are trained with contrastive loss.

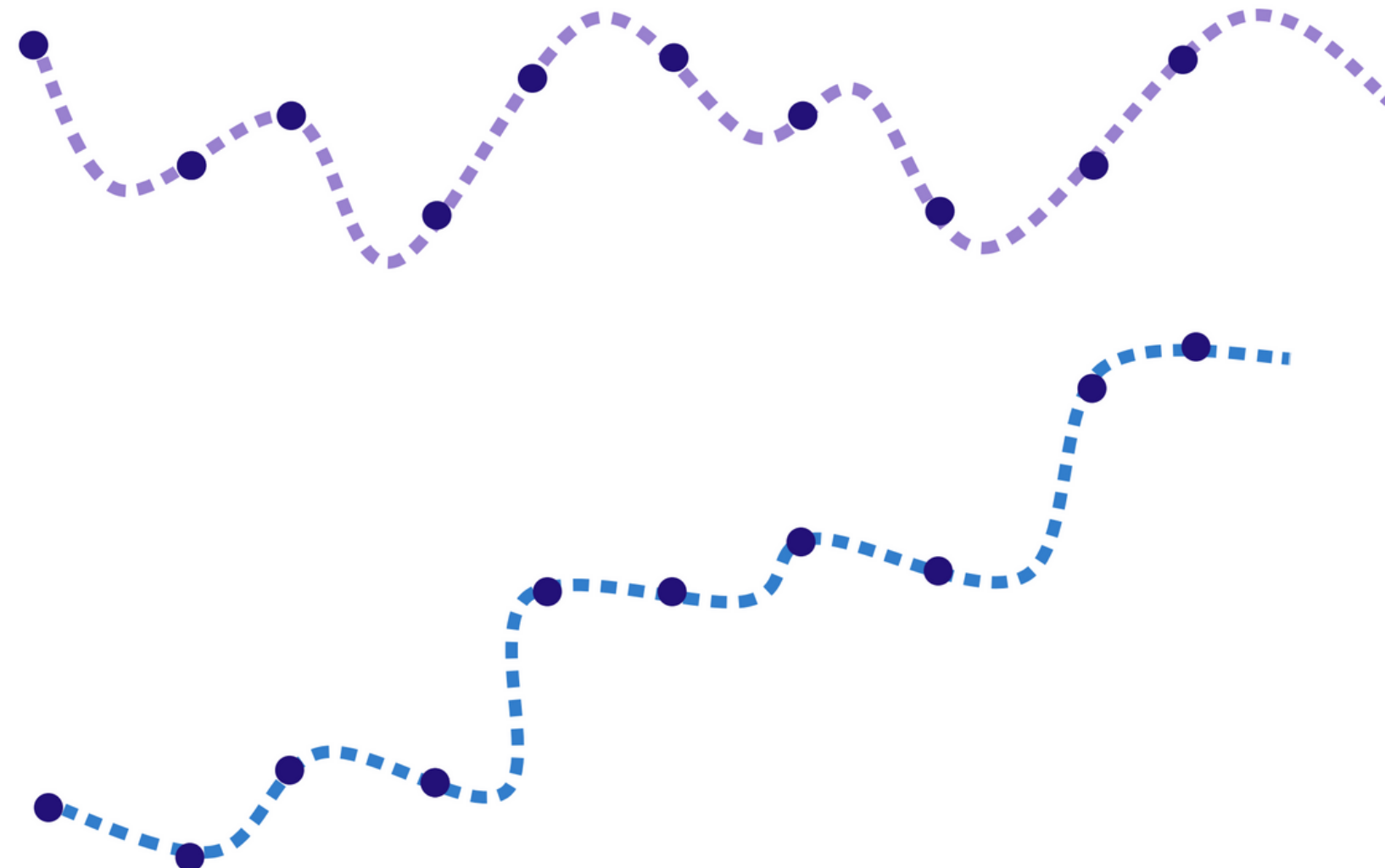


Method

Bi-Trajectory Guided Vision-Language Co-Distillation

Concretely, our approach consists of two stages:

Stage 1 Expert training



$$\tau^* = \{\theta_t^*\}_{t=0}^T$$



Method

Bi-Trajectory Guided Vision-Language Co-Distillation

Concretely, our approach consists of two stages:

Stage 2 Distillation

Training a set of student models on the current distilled dataset $\hat{\mathbf{D}} = \{(\hat{x}_j, \hat{y}_j)\}_{j=1}^M$ using the same contrastive loss.

Update $\hat{\mathbf{D}}$ based on the **bi-trajectory matching loss** of the student models' parameter trajectories and the expert trajectories \mathcal{T}^* .



Method

Bi-Trajectory Guided Vision-Language Co-Distillation

Bi-trajectory matching loss:

$$\ell_{trajectory} = \frac{\|\hat{\theta}_{img,s+\hat{R}} - \theta_{img,s+R}^*\|_2^2}{\|\theta_{img,s}^* - \theta_{img,s+R}^*\|_2^2} + \frac{\|\hat{\theta}_{txt,s+\hat{R}} - \theta_{txt,s+R}^*\|_2^2}{\|\theta_{txt,s}^* - \theta_{txt,s+R}^*\|_2^2}$$



Method

Bi-Trajectory Guided Vision-Language Co-Distillation

Bi-trajectory matching loss:

$$\tau^* = \{\theta_t^*\}_{t=0}^T$$

$$\ell_{trajectory} = \frac{\|\hat{\theta}_{img,s+\hat{R}} - \theta_{img,s+R}^*\|_2^2}{\|\theta_{img,s}^* - \theta_{img,s+R}^*\|_2^2} + \frac{\|\hat{\theta}_{txt,s+\hat{R}} - \theta_{txt,s+R}^*\|_2^2}{\|\theta_{txt,s}^* - \theta_{txt,s+R}^*\|_2^2}$$



Method

Bi-Trajectory Guided Vision-Language Co-Distillation

Bi-trajectory matching loss:

$$\tau^* = \{\theta_t^*\}_{t=0}^T$$

$$\ell_{trajectory} = \frac{\|\hat{\theta}_{img, s+\hat{R}} - \theta_{img, s+R}^*\|_2^2}{\|\theta_{img, s}^* - \theta_{img, s+R}^*\|_2^2} + \frac{\|\hat{\theta}_{txt, s+\hat{R}} - \theta_{txt, s+R}^*\|_2^2}{\|\theta_{txt, s}^* - \theta_{txt, s+R}^*\|_2^2}$$



Method

Bi-Trajectory Guided Vision-Language Co-Distillation

Bi-trajectory matching loss:

$$\tau^* = \{\theta_t^*\}_{t=0}^T$$

$$\ell_{trajectory} = \frac{\|\hat{\theta}_{img,s+\hat{R}} - \theta_{img,s+R}^*\|_2^2}{\|\theta_{img,s}^* - \theta_{img,s+R}^*\|_2^2} + \frac{\|\hat{\theta}_{txt,s+\hat{R}} - \theta_{txt,s+R}^*\|_2^2}{\|\theta_{txt,s}^* - \theta_{txt,s+R}^*\|_2^2}$$

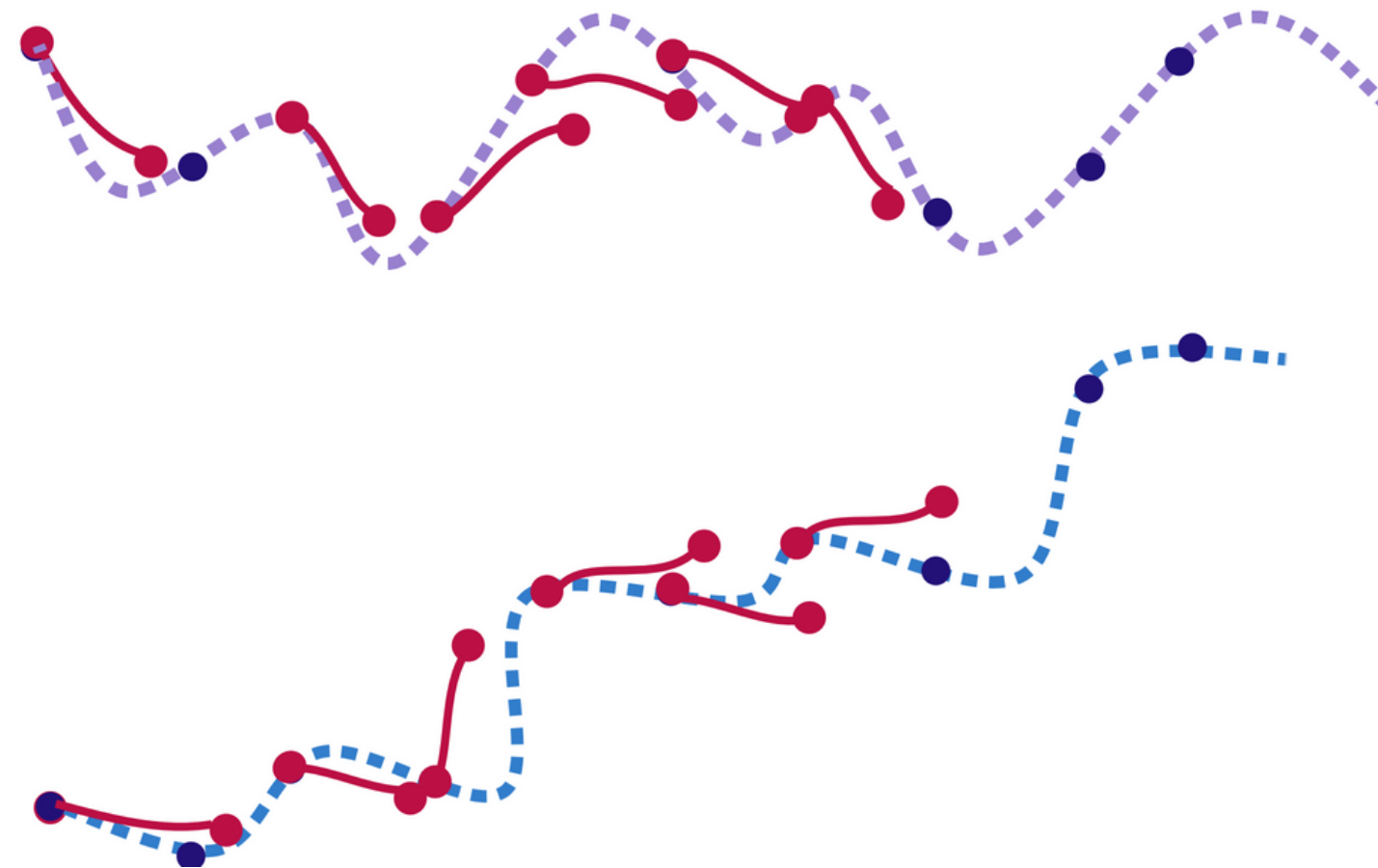


Method

Bi-Trajectory Guided Vision-Language Co-Distillation

Concretely, our approach consists of two stages:

Stage 2 Distillation



Method

Low-Rank Adaptation Matching

Vision Transformers (ViTs)

1. High dimensionality of the embeddings
2. Large number of parameters

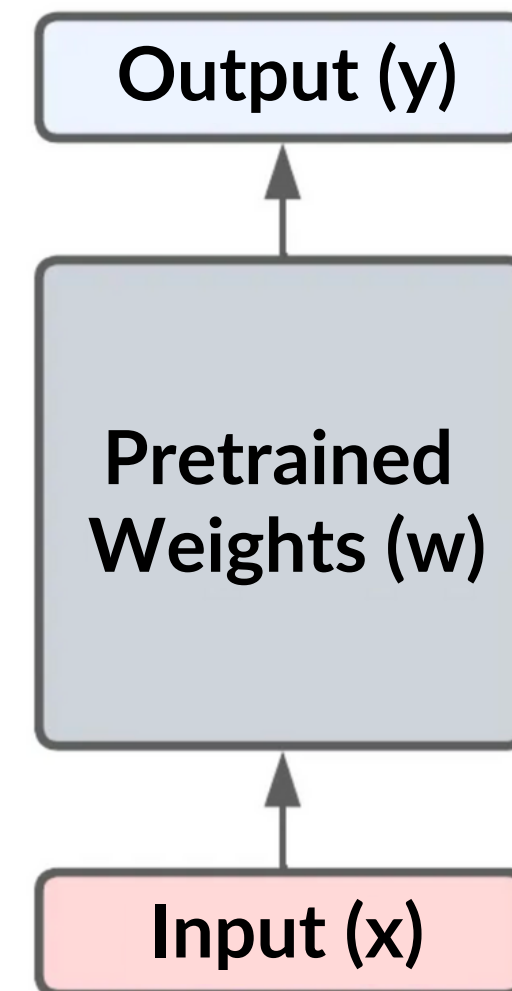
Method

Low-Rank Adaptation Matching

Consider a linear layer with
m input units, n output units

The weight matrix for this layer has dimensions m x n

$$Y = W X$$



Method

Low-Rank Adaptation Matching

Consider a linear layer with
m input units, n output units

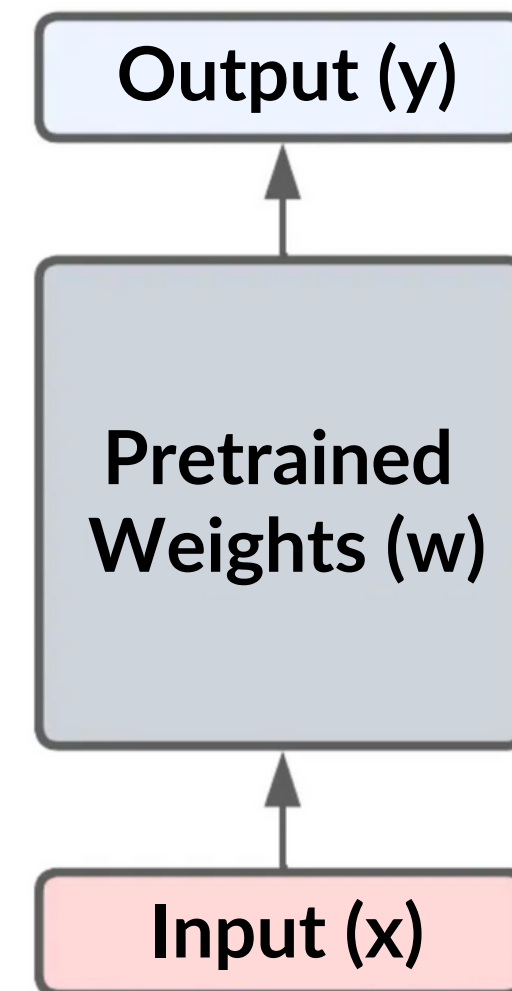
The weight matrix for this layer has dimensions m x n

$$Y = W X$$

m = 800

n = 3200

W: 800 x 3200 = 2,560,000 weights.



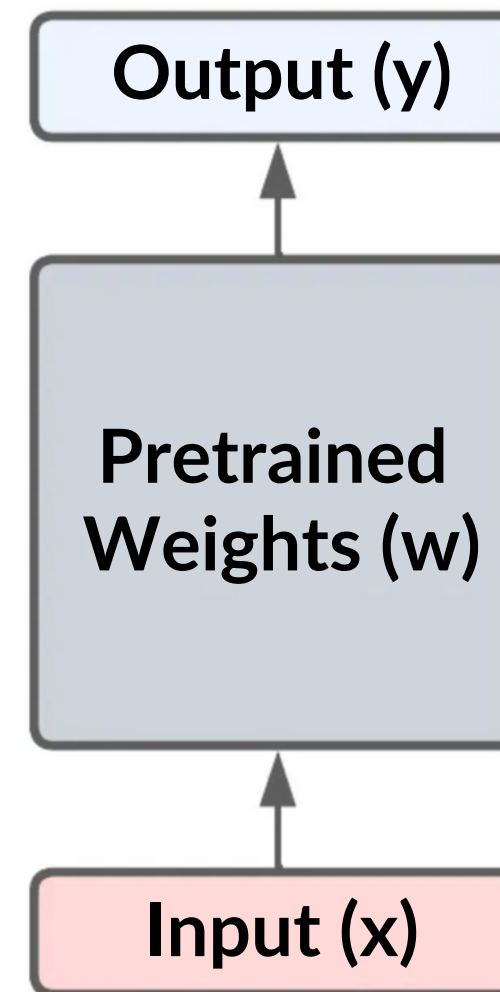
Method

Low-Rank Adaptation Matching

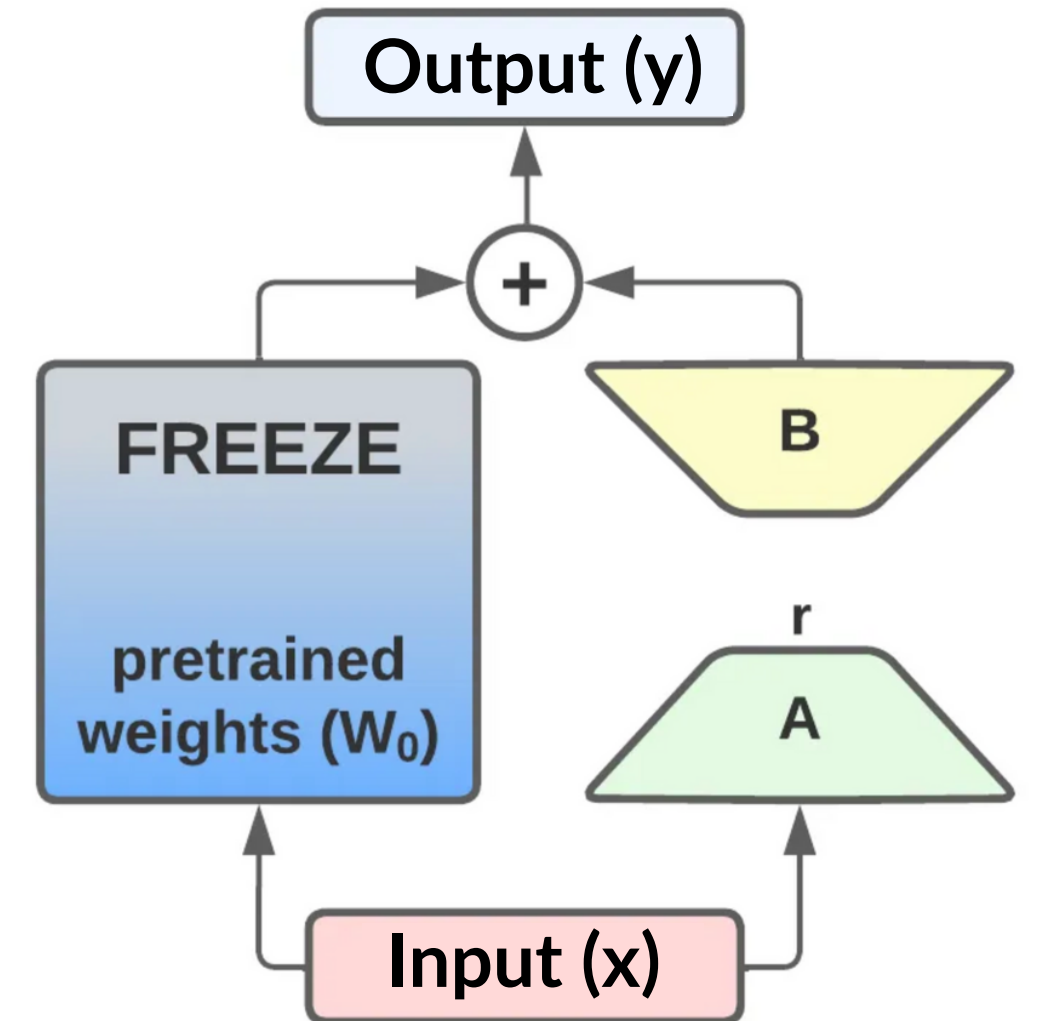
Low-Rank Adaptation (LoRA)

We keep W fixed and introduce two matrices, A and B

$$Y = W X + A * B X$$



Low-Rank Adaptation



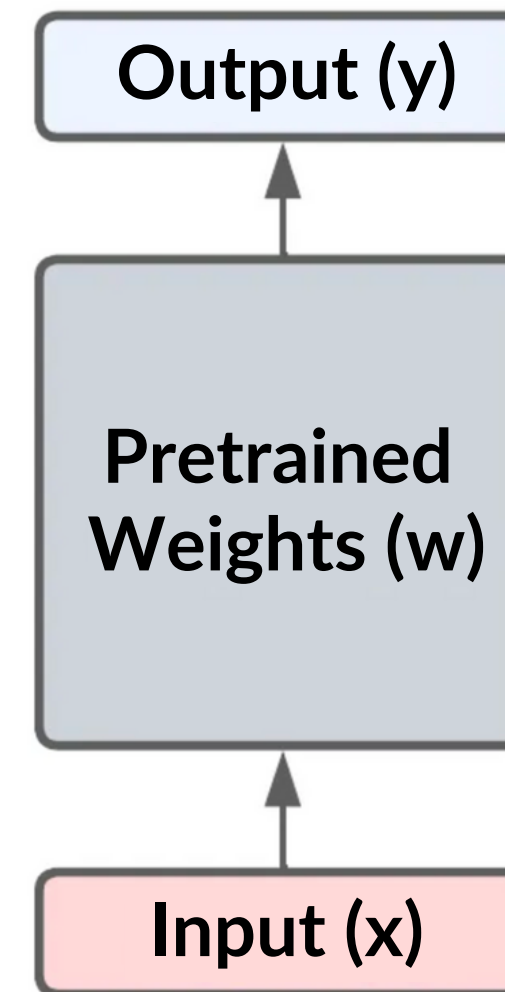
Method

Low-Rank Adaptation Matching

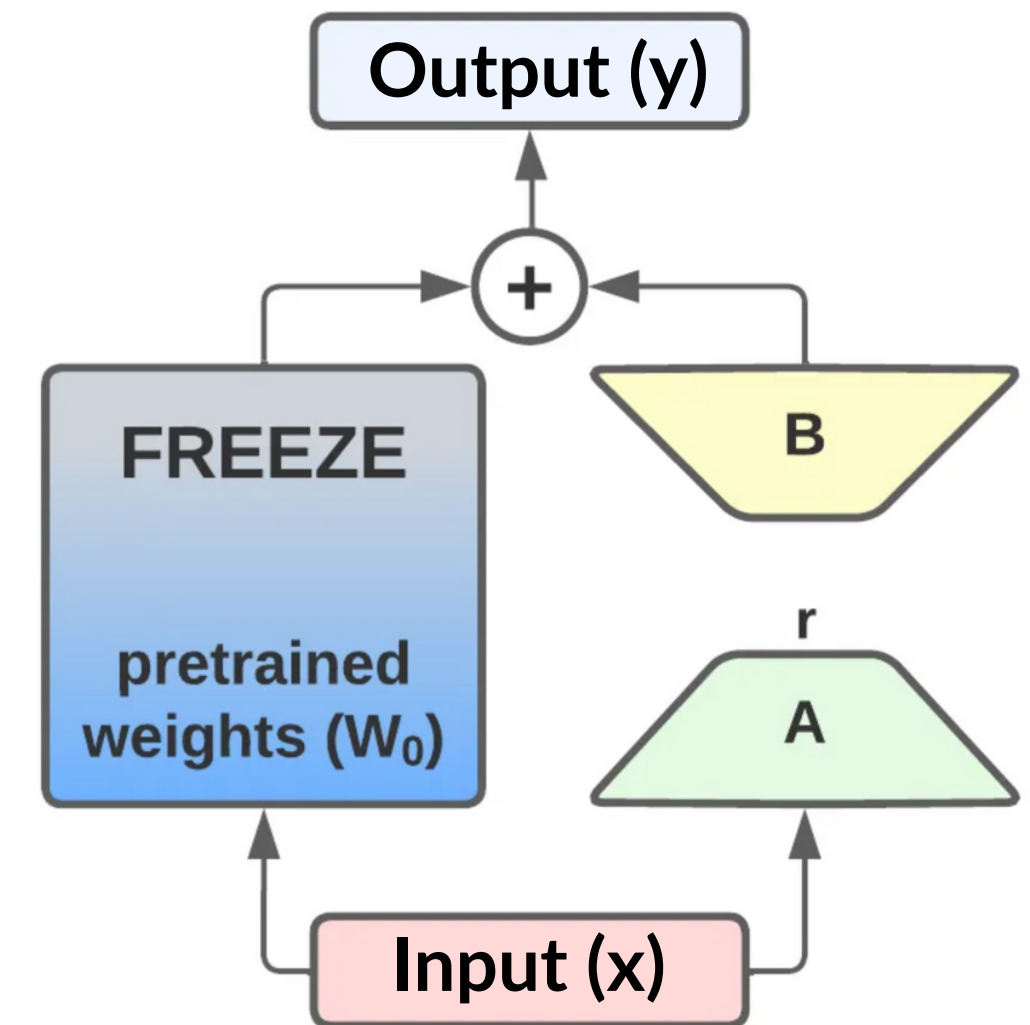
Low-Rank Adaptation (LoRA)

We keep W fixed and introduce two matrices, A and B

$$Y = W X + A * B X$$



Low-Rank Adaptation



Matrix A has a shape of $800 \times r$, and matrix B has a shape of $r \times 3200$.

$$m = 800, n = 3200, r = 1$$

$$(800 \times 1) + (1 \times 3200) = 4000$$

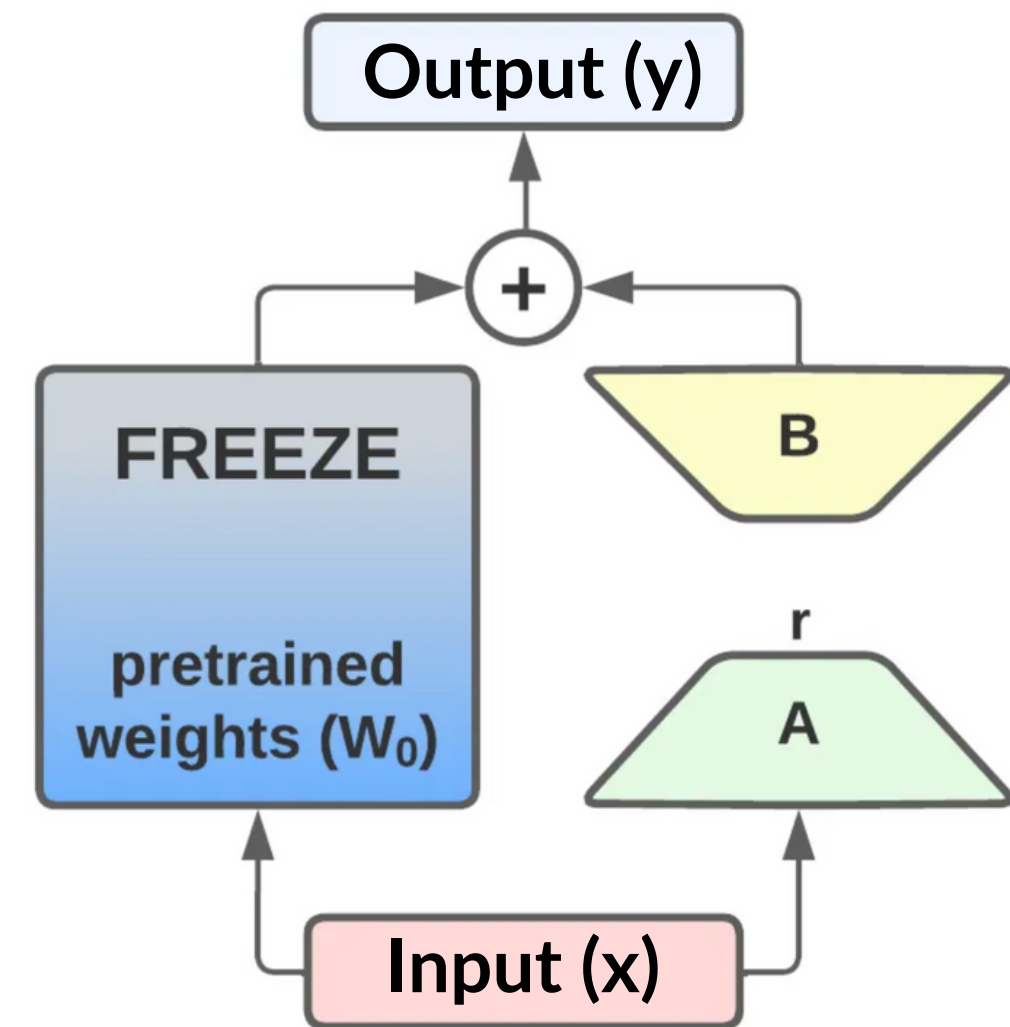
Method

Low-Rank Adaptation Matching

Low-Rank Adaptation (LoRA) Matching

LoRA matching optimizes the trajectories of low rank adapters instead of the full parameters.

Low-Rank Adaptation



Method

Low-Rank Adaptation Matching

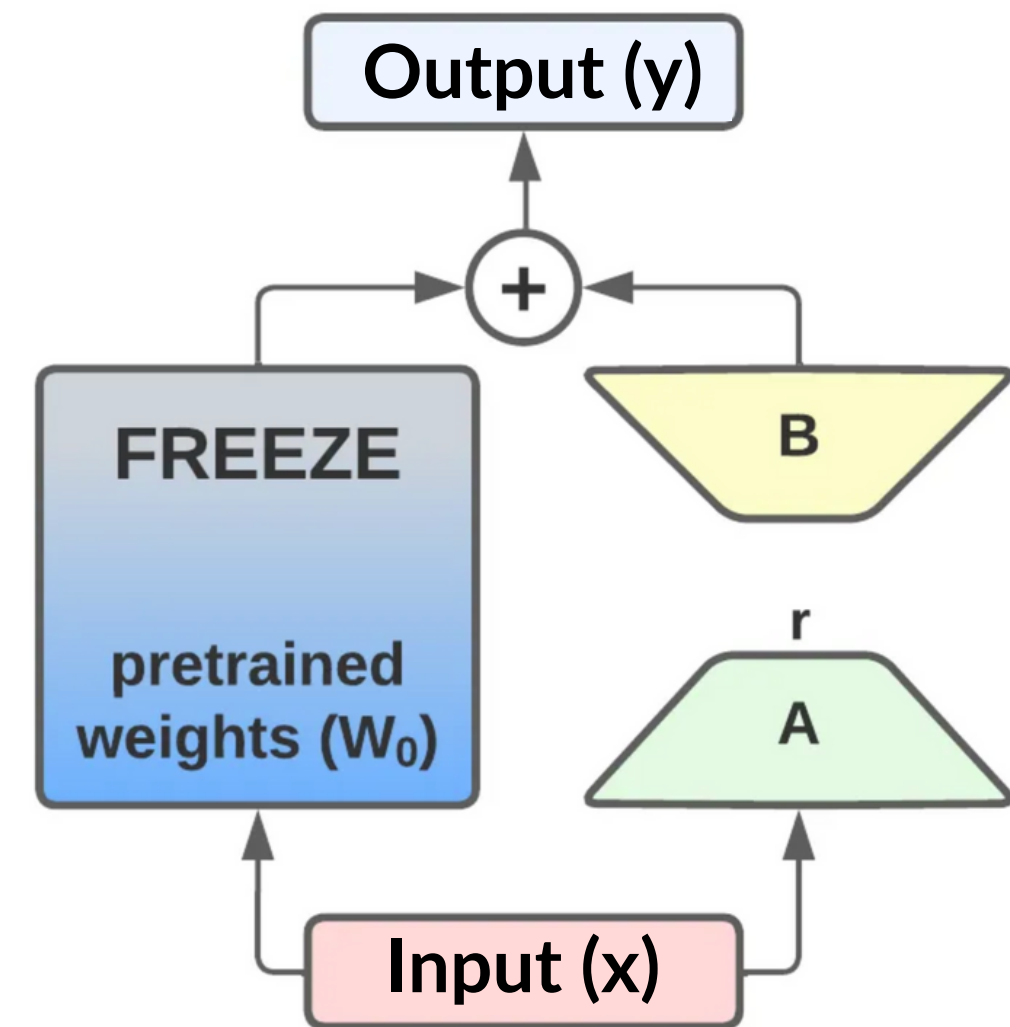
Low-Rank Adaptation (LoRA) Matching

LoRA matching optimizes the trajectories of low rank adapters instead of the full parameters.

VIT_base

86 million

Low-Rank Adaptation



Method

Low-Rank Adaptation Matching

Low-Rank Adaptation (LoRA) Matching

LoRA matching optimizes the trajectories of low rank adapters instead of the full parameters.

VIT_base

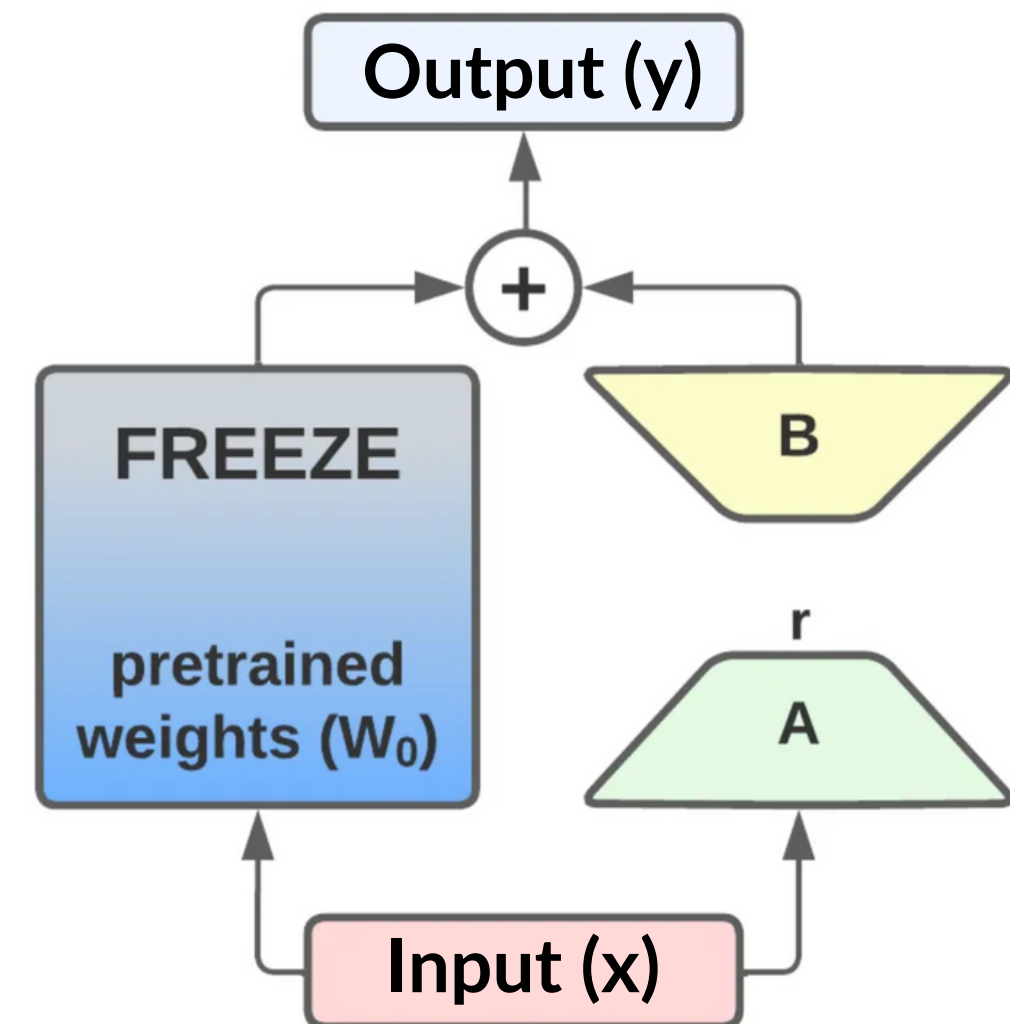
86 million

LoRA_VIT_base

18 million

$r=4$

Low-Rank Adaptation



Method

Low-Rank Adaptation Matching

Low-Rank Adaptation (LoRA) Matching

LoRA matching optimizes the trajectories of low rank adapters instead of the full parameters.
Efficient model adaptation with minimal extra parameters.
Focus on optimizing a smaller parameter set.
Efficiently optimize trajectory loss during distillation.

VIT_base

86 million

LoRA_VIT_base

18 million

$r=4$

Low-Rank Adaptation

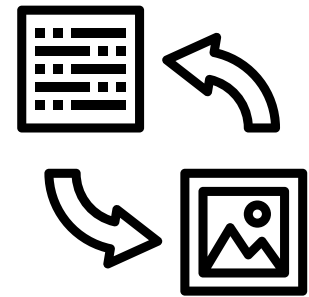




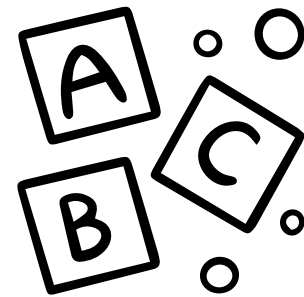
Challenges <> Solutions

Vision-Language Dataset Distillation

Complex Cross-Modal Relationships



Discrete text optimization issue



Heavy computational cost

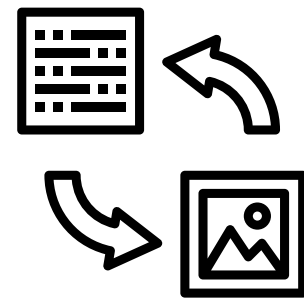




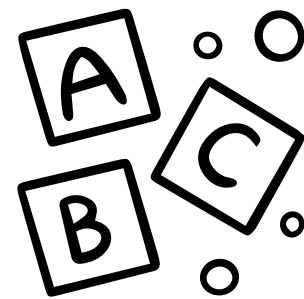
Challenges <> Solutions

Vision-Language Dataset Distillation

Complex Cross-Modal Relationships



Discrete text optimization issue



Heavy computational cost

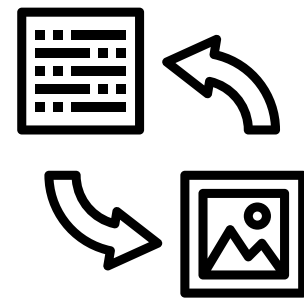


Bi-trajectory Matching with Contrastive Loss

Challenges <> Solutions

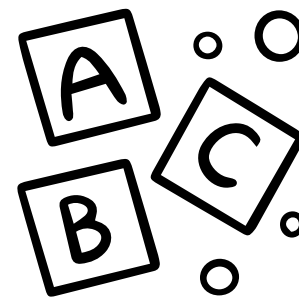
Vision-Language Dataset Distillation

Complex Cross-Modal Relationships



Bi-trajectory Matching with Contrastive Loss

Discrete text optimization issue



Continuous Sentence Embeddings

Heavy computational cost

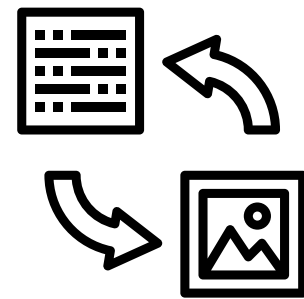




Challenges <> Solutions

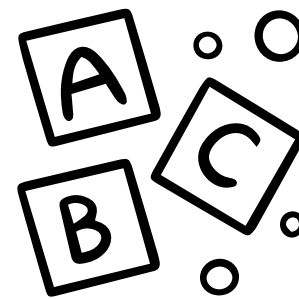
Vision-Language Dataset Distillation

Complex Cross-Modal Relationships



Bi-trajectory Matching with Contrastive Loss

Discrete text optimization issue



Continuous Sentence Embeddings

Heavy computational cost



Finetuning with Low-Rank Adaptation

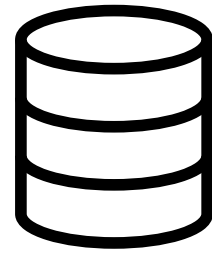


| Experiments

Experiments

Evaluation Setting

Dataset



Flickr30K

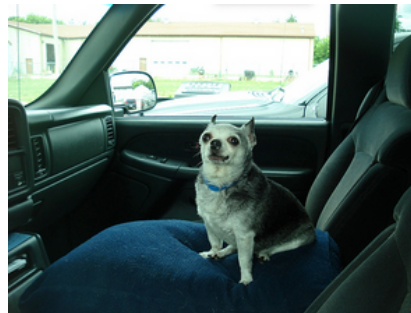


A small dog wearing a sweater walking in the snow.

....

30K

COCO



a grey dog seated on a chair of a vehicle

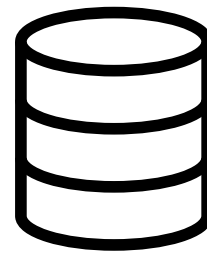
....

328K

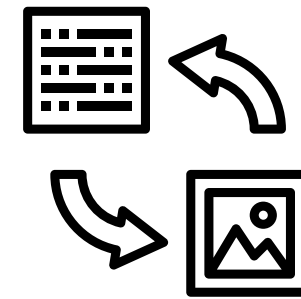
Experiments

Evaluation Setting

Dataset



Task



Flickr30K

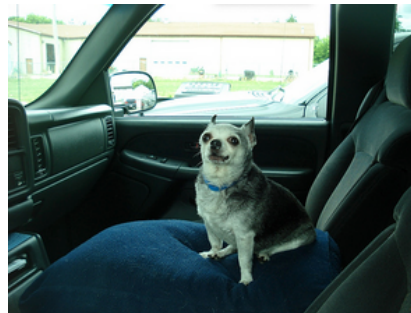


A small dog wearing a sweater walking in the snow.

....

30K

COCO



a grey dog seated on a chair of a vehicle

....

328K

Image-text retrieval

image-to-text retrieval (TR) & text-to-image retrieval (IR)

R@K (with K = 1, 5, 10)

compute the fraction of times the correct result is found among the top K items.

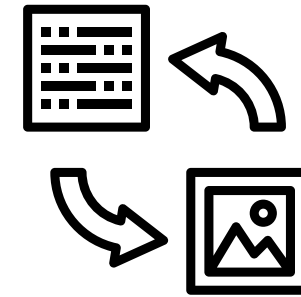
Experiments

Evaluation Setting

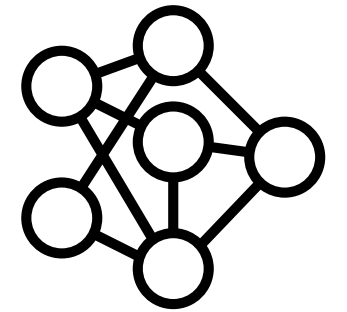
Dataset



Task



Model



Flickr30K

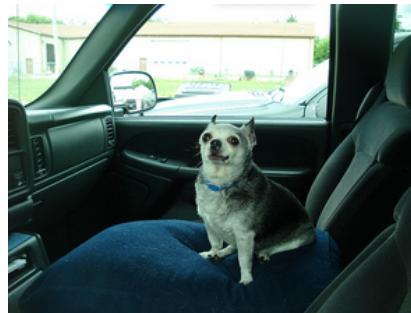


A small dog wearing a sweater walking in the snow.

....

30K

COCO



a grey dog seated on a chair of a vehicle

....

328K

Image-text retrieval

image-to-text retrieval (TR) & text-to-image retrieval (IR)

$R@K$ (with $K = 1, 5, 10$)

compute the fraction of times the correct result is found among the top K items.

Image encoder:

- NFFNet
- ViT

Language encoder:

- BERT



Experiments

Quantitative Results | [Baseline Comparisons](#)

Dataset	#pairs	ratio %	TR					IR				
			Coreset Selection				Dist (ours)	Coreset Selection				Dist (ours)
			R	H	K	F		R	H	K	F	
Flickr30K	100	0.34	1.3	1.1	0.6	1.2	9.9 ± 0.3	1.0	0.7	0.7	0.7	4.7 ± 0.2
	200	0.68	2.1	2.3	2.2	1.5	10.2 ± 0.8	1.1	1.5	1.5	1.2	4.6 ± 0.9
	500	1.67	5.2	5.1	4.9	3.6	13.3 ± 0.6	2.4	3.0	3.5	1.8	6.6 ± 0.3
	1000	3.45	5.2	5	5.6	3.1	13.3 ± 1.0	3.8	4.1	4.4	3.2	7.9 ± 0.8
COCO	100	0.08	0.8	0.8	1.4	0.7	2.5 ± 0.3	0.3	0.5	0.4	0.3	1.3 ± 0.1
	200	0.17	1.0	1.0	1.2	1.1	3.3 ± 0.2	0.6	0.9	0.7	0.6	1.7 ± 0.1
	500	0.44	1.9	1.9	2.5	2.1	5.0 ± 0.4	1.1	1.7	1.1	0.8	2.5 ± 0.5
	1000	0.88	1.9	2.4	2.4	1.9	6.8 ± 0.4	1.5	1.3	1.5	0.7	3.3 ± 0.1

Random selection of training examples (**R**) | Herding (**H**) K-center (**K**) | Forgetting (**F**)



Experiments

Quantitative Results | [Baseline Comparisons](#)

Dataset	#pairs	ratio %	TR					IR				
			Coreset Selection				Dist (ours)	Coreset Selection				Dist (ours)
			R	H	K	F		R	H	K	F	
Flickr30K	100	0.34	1.3	1.1	0.6	1.2	9.9 ± 0.3	1.0	0.7	0.7	0.7	4.7 ± 0.2
	200	0.68	2.1	2.3	2.2	1.5	10.2 ± 0.8	1.1	1.5	1.5	1.2	4.6 ± 0.9
	500	1.67	5.2	5.1	4.9	3.6	13.3 ± 0.6	2.4	3.0	3.5	1.8	6.6 ± 0.3
	1000	3.45	5.2	5	5.6	3.1	13.3 ± 1.0	3.8	4.1	4.4	3.2	7.9 ± 0.8
COCO	100	0.08	0.8	0.8	1.4	0.7	2.5 ± 0.3	0.3	0.5	0.4	0.3	1.3 ± 0.1
	200	0.17	1.0	1.0	1.2	1.1	3.3 ± 0.2	0.6	0.9	0.7	0.6	1.7 ± 0.1
	500	0.44	1.9	1.9	2.5	2.1	5.0 ± 0.4	1.1	1.7	1.1	0.8	2.5 ± 0.5
	1000	0.88	1.9	2.4	2.4	1.9	6.8 ± 0.4	1.5	1.3	1.5	0.7	3.3 ± 0.1

Random selection of training examples (**R**) | Herding (**H**) K-center (**K**) | Forgetting (**F**)



Experiments

Quantitative Results | [Baseline Comparisons](#)

Dataset	#pairs	ratio %	TR					IR				
			Coreset Selection				Dist (ours)	Coreset Selection				Dist (ours)
			R	H	K	F		R	H	K	F	
Flickr30K	100	0.34	1.3	1.1	0.6	1.2	9.9 ± 0.3	1.0	0.7	0.7	0.7	4.7 ± 0.2
	200	0.68	2.1	2.3	2.2	1.5	10.2 ± 0.8	1.1	1.5	1.5	1.2	4.6 ± 0.9
	500	1.67	5.2	5.1	4.9	3.6	13.3 ± 0.6	2.4	3.0	3.5	1.8	6.6 ± 0.3
	1000	3.45	5.2	5	5.6	3.1	13.3 ± 1.0	3.8	4.1	4.4	3.2	7.9 ± 0.8
COCO	100	0.08	0.8	0.8	1.4	0.7	2.5 ± 0.3	0.3	0.5	0.4	0.3	1.3 ± 0.1
	200	0.17	1.0	1.0	1.2	1.1	3.3 ± 0.2	0.6	0.9	0.7	0.6	1.7 ± 0.1
	500	0.44	1.9	1.9	2.5	2.1	5.0 ± 0.4	1.1	1.7	1.1	0.8	2.5 ± 0.5
	1000	0.88	1.9	2.4	2.4	1.9	6.8 ± 0.4	1.5	1.3	1.5	0.7	3.3 ± 0.1

Random selection of training examples (**R**) | Herding (**H**) K-center (**K**) | Forgetting (**F**)

Experiments

Qualitative Results | [Distilled Examples](#)



a man in a black wet suit is surfing a huge wave in the beautiful blue water

before the distillation



a man surfs over a huge wave in the ocean

after 2000 distillation steps

Note that the texts visualized here are nearest sentence decodings in the training set corresponding to the distilled text embeddings.

Experiments

Qualitative Results | Before and After Distillation



a newly married couple sharing a kiss in front of a convertible



a couple kissed in front of a beautiful three-tiered cake with blue ribbon and pink accents



a man in a black wetsuit is surfing a huge wave in the beautiful blue water



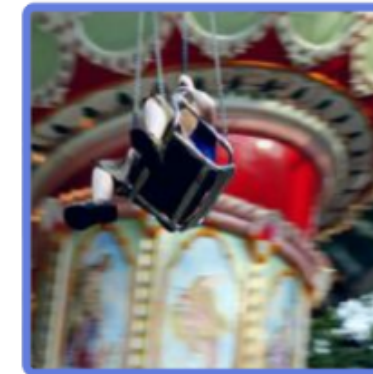
a man surfs over a huge wave in the ocean



a woman in a toboggan sledding with a child that is in a toboggan also



skateboarder in jeans and t-shirt performing jump



man is sitting in a swing on a carnival ride



a little boy in a toy room wearing batman pajamas is building something with toys



a little boy in a white shirt is rock climbing



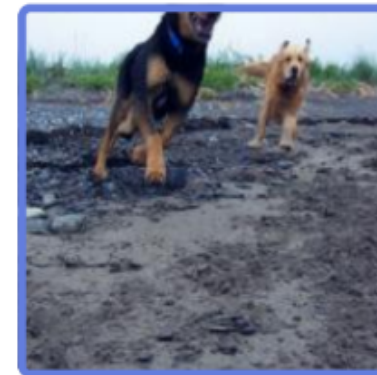
a boy with a blue helmet and gray pants is rock-climbing



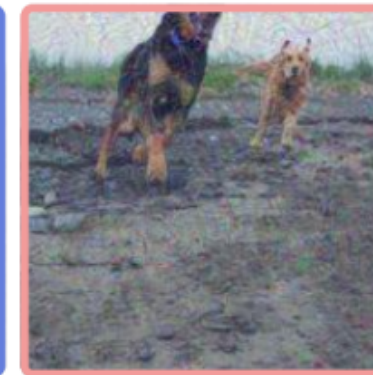
two boys are watching another boy perform a jump on his bmx bike



kid in hoodie jumps a ramp



two dogs run through mud



the brown and white dog nips at the yellow dog



two men are competing for the ball in a game of soccer



four football players look on from the sideline while two teams are in formation at the line of scrimmage

Experiments

Qualitative Results | Before and After Distillation





Experiments

Ablation Studies

With and without LoRA on ViT

Dataset	#Pairs	Without LoRA		With LoRA	
		TR	IR	TR	IR
Flickr30K	100	1.5	0.6	10.4 ± 0.8	5.4 ± 0.2
	1000	3.3	1.5	15.8 ± 1.4	8.1 ± 0.7



Experiments

Ablation Studies

With and without LoRA on ViT

Dataset	#Pairs	Without LoRA		With LoRA	
		TR	IR	TR	IR
Flickr30K	100	1.5	0.6	10.4 ± 0.8	5.4 ± 0.2
	1000	3.3	1.5	15.8 ± 1.4	8.1 ± 0.7

(Here we only reports R@1, more details are in the paper.)



Experiments

Ablation Studies

With and without LoRA on ViT

Dataset	#Pairs	Without LoRA		With LoRA	
		TR	IR	TR	IR
Flickr30K	100	1.5	0.6	10.4 ± 0.8	5.4 ± 0.2
	1000	3.3	1.5	15.8 ± 1.4	8.1 ± 0.7

Different vis./lan. encoders

Language Model	TR	IR
BERT	9.9	4.7
CLIP	31.4	17.1

Vision Model	TR	IR
NFNet	9.9	4.7
VIT_LoRA	10.4	5.4
NF_ResNet50	6.5	3.46
NF_RegNet	7.8	3.28



Experiments

Ablation Studies

With and without LoRA on ViT

Dataset	#Pairs	Without LoRA		With LoRA	
		TR	IR	TR	IR
Flickr30K	100	1.5	0.6	10.4 ± 0.8	5.4 ± 0.2
	1000	3.3	1.5	15.8 ± 1.4	8.1 ± 0.7

Different vis./lan. encoders

Language Model	TR	IR
BERT	9.9	4.7
CLIP	31.4	17.1

Vision Model	TR	IR
NFNet	9.9	4.7
VIT_LoRA	10.4	5.4
NF_ResNet50	6.5	3.46
NF_RegNet	7.8	3.28

(Here we only reports R@1, more details are in the paper.)



Experiments

Ablation Studies

With and without LoRA on ViT

Dataset	#Pairs	Without LoRA		With LoRA	
		TR	IR	TR	IR
Flickr30K	100	1.5	0.6	10.4 ± 0.8	5.4 ± 0.2
	1000	3.3	1.5	15.8 ± 1.4	8.1 ± 0.7

Single-modality vs. multi-modality

Takeaway:

- Distillation would be impossible if we solely optimize one modality

Different vis./lan. encoders

Language Model	TR	IR
BERT	9.9	4.7
CLIP	31.4	17.1

Vision Model	TR	IR
NFNet	9.9	4.7
VIT_LoRA	10.4	5.4
NF_ResNet50	6.5	3.46
NF_RegNet	7.8	3.28

	TR	IR
T	1.3	0.5
I	3.5	1.6
Ours	9.9	4.7

T: text-only, I: image-only



Experiments

Ablation Studies

With and without LoRA on ViT

Dataset	#Pairs	Without LoRA		With LoRA	
		TR	IR	TR	IR
Flickr30K	100	1.5	0.6	10.4 ± 0.8	5.4 ± 0.2
	1000	3.3	1.5	15.8 ± 1.4	8.1 ± 0.7

Single-modality vs. multi-modality

Takeaway:

- Distillation would be impossible if we solely optimize one modality

Different vis./lan. encoders

Language Model	TR	IR
BERT	9.9	4.7
CLIP	31.4	17.1

Vision Model	TR	IR
NFNet	9.9	4.7
VIT_LoRA	10.4	5.4
NF_ResNet50	6.5	3.46
NF_RegNet	7.8	3.28

	TR	IR
T	1.3	0.5
I	3.5	1.6
Ours	9.9	4.7

T: text-only, I: image-only



Experiments

Ablation Studies

With and without LoRA on ViT

Dataset	#Pairs	Without LoRA		With LoRA	
		TR	IR	TR	IR
Flickr30K	100	1.5	0.6	10.4 ± 0.8	5.4 ± 0.2
	1000	3.3	1.5	15.8 ± 1.4	8.1 ± 0.7

Different vis./lan. encoders

Language Model	TR	IR
BERT	9.9	4.7
CLIP	31.4	17.1

Vision Model	TR	IR
NFNet	9.9	4.7
VIT_LoRA	10.4	5.4
NF_ResNet50	6.5	3.46
NF_RegNet	7.8	3.28

Single-modality vs. multi-modality

Takeaway:

- Distillation would be impossible if we solely optimize one modality

	TR	IR
T	1.3	0.5
I	3.5	1.6
Ours	9.9	4.7

T: text-only, I: image-only

Image-Text Pair Initialization

Takeaway:

- Initializing texts from scratch
- Initializing images from scratch

Real Image	Real Text	TR	IR
		0.4	0.1
	✓	1.1	0.1
✓		9	3.9
✓	✓	9.9	4.7



Experiments

Ablation Studies

With and without LoRA on ViT

Dataset	#Pairs	Without LoRA		With LoRA	
		TR	IR	TR	IR
Flickr30K	100	1.5	0.6	10.4 ± 0.8	5.4 ± 0.2
	1000	3.3	1.5	15.8 ± 1.4	8.1 ± 0.7

Different vis./lan. encoders

Language Model	TR	IR
BERT	9.9	4.7
CLIP	31.4	17.1

Vision Model	TR	IR
NFNet	9.9	4.7
VIT_LoRA	10.4	5.4
NF_ResNet50	6.5	3.46
NF_RegNet	7.8	3.28

Single-modality vs. multi-modality

Takeaway:



- Distillation would be impossible if we solely optimize one modality

	TR	IR
T	1.3	0.5
I	3.5	1.6
Ours	9.9	4.7

T: text-only, I: image-only

Image-Text Pair Initialization

Takeaway:

-  Initializing texts from scratch
-  Initializing images from scratch

Real Image	Real Text	TR	IR
		0.4	0.1
	✓	1.1	0.1
✓		9	3.9
✓	✓	9.9	4.7



Experiments

Ablation Studies

With and without LoRA on ViT

Dataset	#Pairs	Without LoRA		With LoRA	
		TR	IR	TR	IR
Flickr30K	100	1.5	0.6	10.4 ± 0.8	5.4 ± 0.2
	1000	3.3	1.5	15.8 ± 1.4	8.1 ± 0.7

Different vis./lan. encoders

Language Model	TR	IR
BERT	9.9	4.7
CLIP	31.4	17.1

Vision Model	TR	IR
NFNet	9.9	4.7
VIT_LoRA	10.4	5.4
NF_ResNet50	6.5	3.46
NF_RegNet	7.8	3.28

Single-modality vs. multi-modality

Takeaway:

- Distillation would be impossible if we solely optimize one modality

	TR	IR
T	1.3	0.5
I	3.5	1.6
Ours	9.9	4.7

T: text-only, I: image-only

Image-Text Pair Initialization

Takeaway:

- Initializing texts from scratch
- Initializing images from scratch

image component plays a more critical role in the distilled dataset.

Real Image	Real Text	TR	IR
		0.4	0.1
	✓	1.1	0.1
✓		9	3.9
✓	✓	9.9	4.7



Summary

- In this work, we propose the first vision-language dataset distillation method.
- Our experiments show that co-distilling different modalities via bi-trajectory matching holds promise.
- We hope that the insights we gathered can serve as roadmap for future studies exploring more complex settings.



Byron Zhang



Zhiwei Deng



Olga Russakovsky



[Website](#)



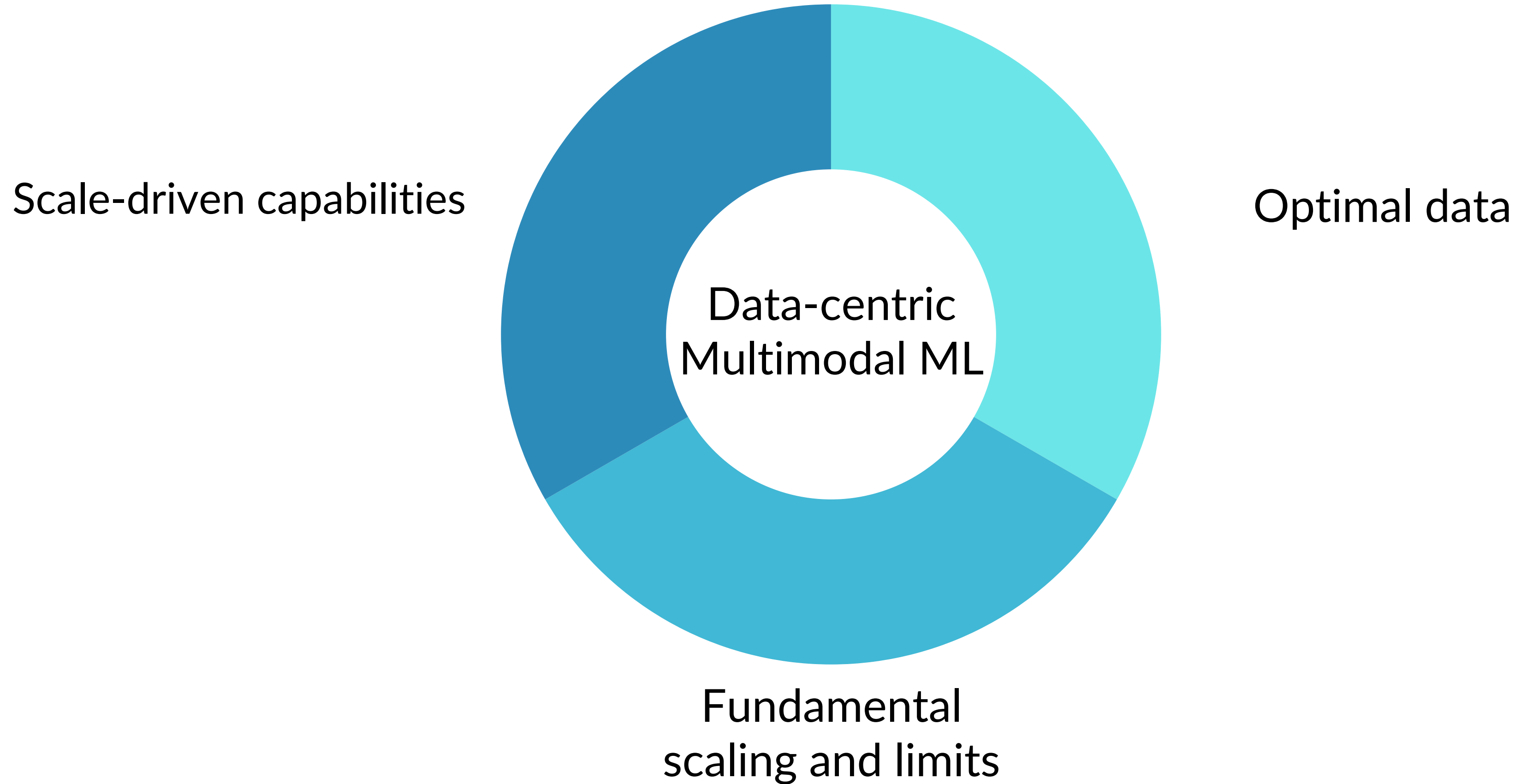
[arXiv](#)



[Code](#)



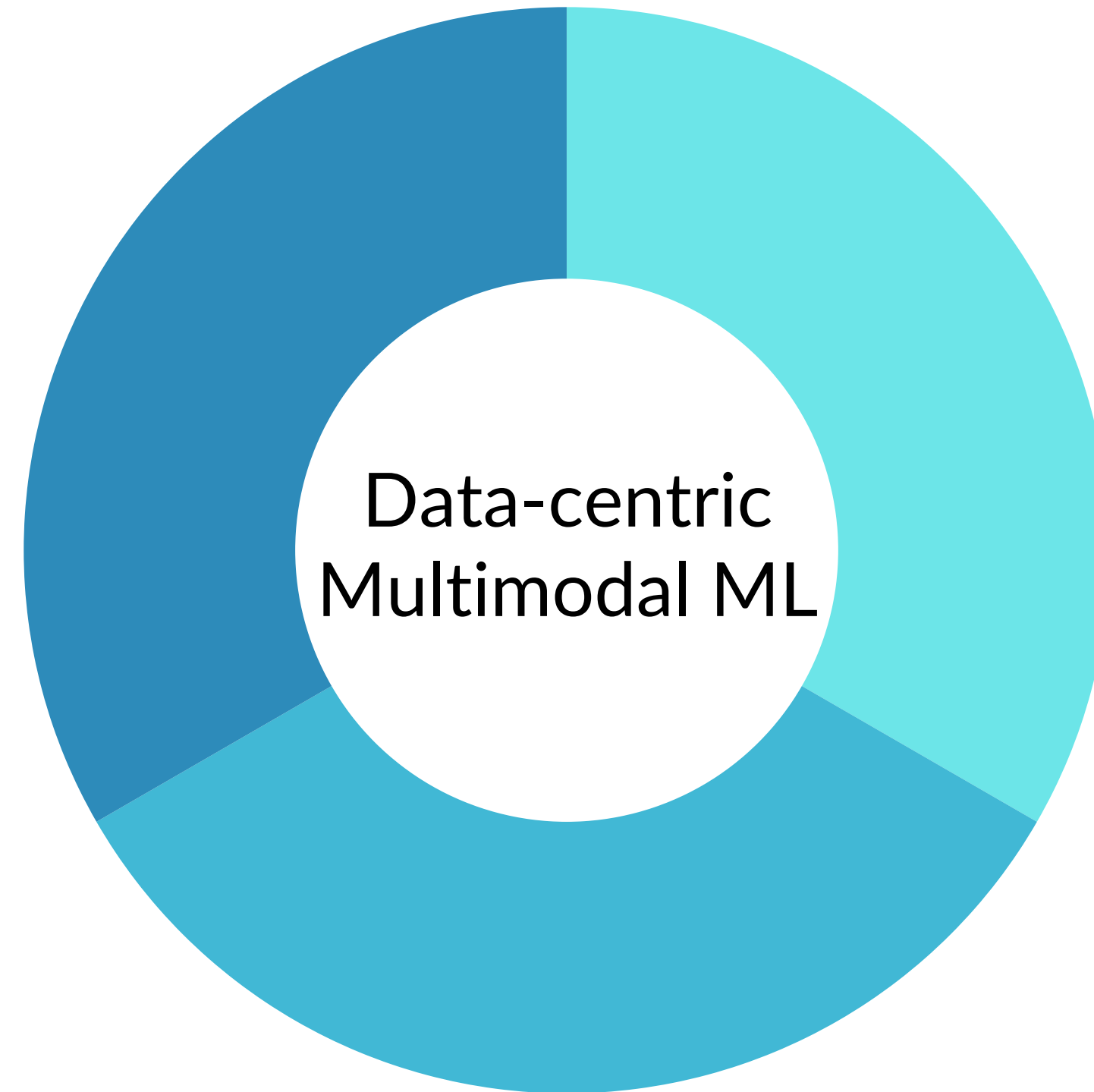
Looking Forward





Looking Forward

Scale-driven capabilities



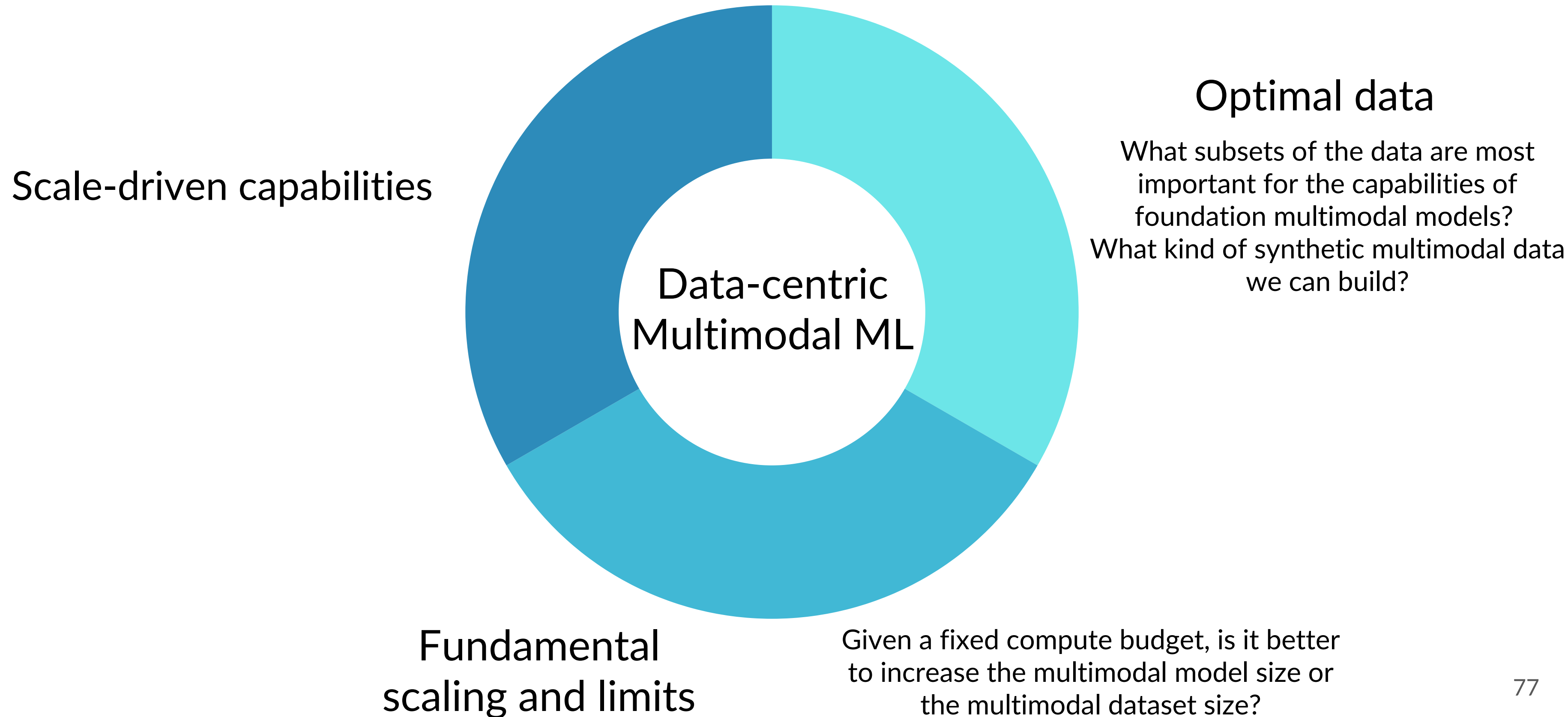
Optimal data

What subsets of the data are most important for the capabilities of foundation multimodal models?
What kind of synthetic multimodal data we can build?

Fundamental scaling and limits



Looking Forward

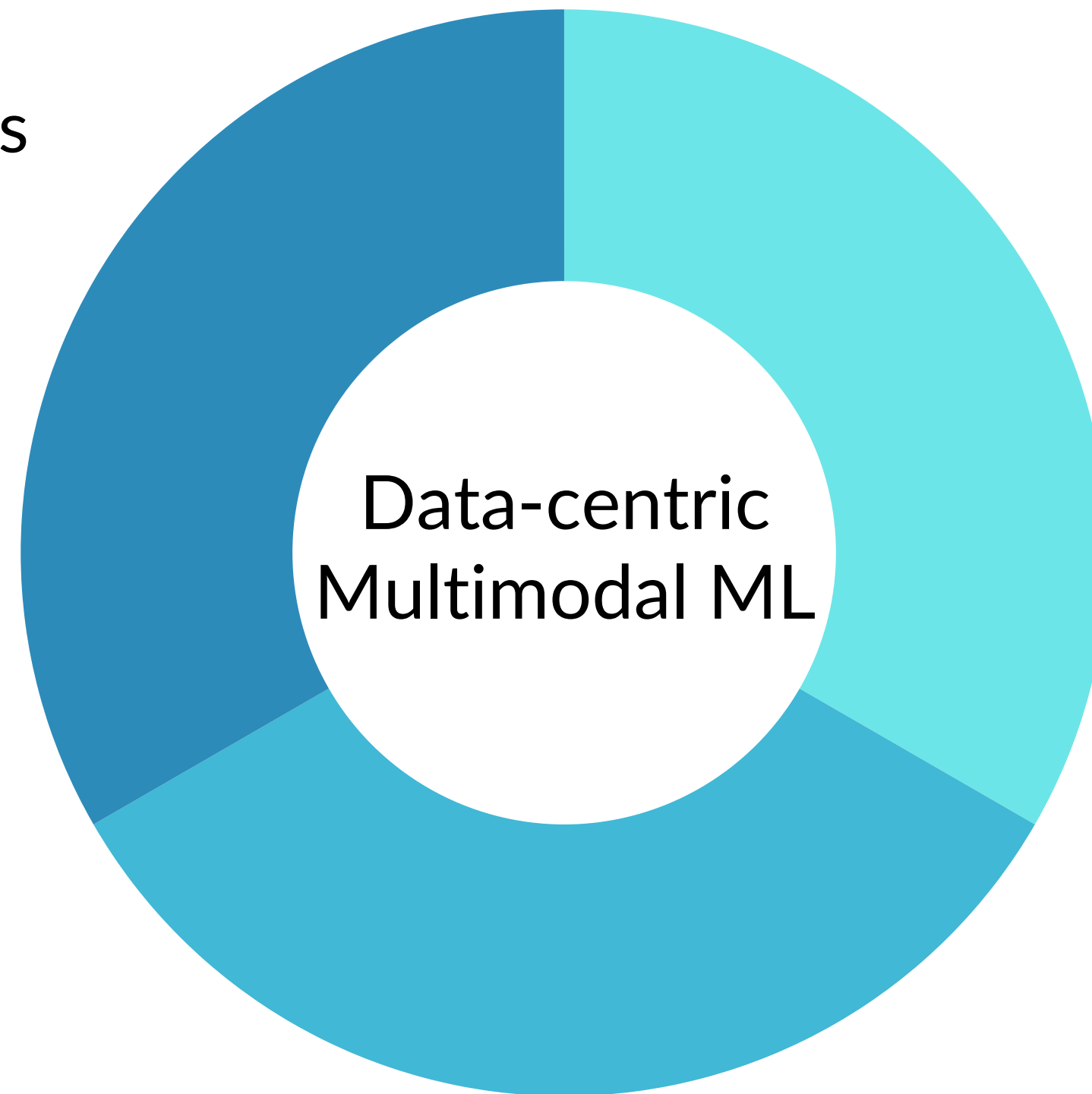




Looking Forward

Scale-driven capabilities

a more rigorous understanding of what increasing the scale does to the multimodal training procedure and how these desirable emergent capabilities come about



Optimal data

What subsets of the data are most important for the capabilities of foundation multimodal models?
What kind of synthetic multimodal data we can build?

Fundamental scaling and limits

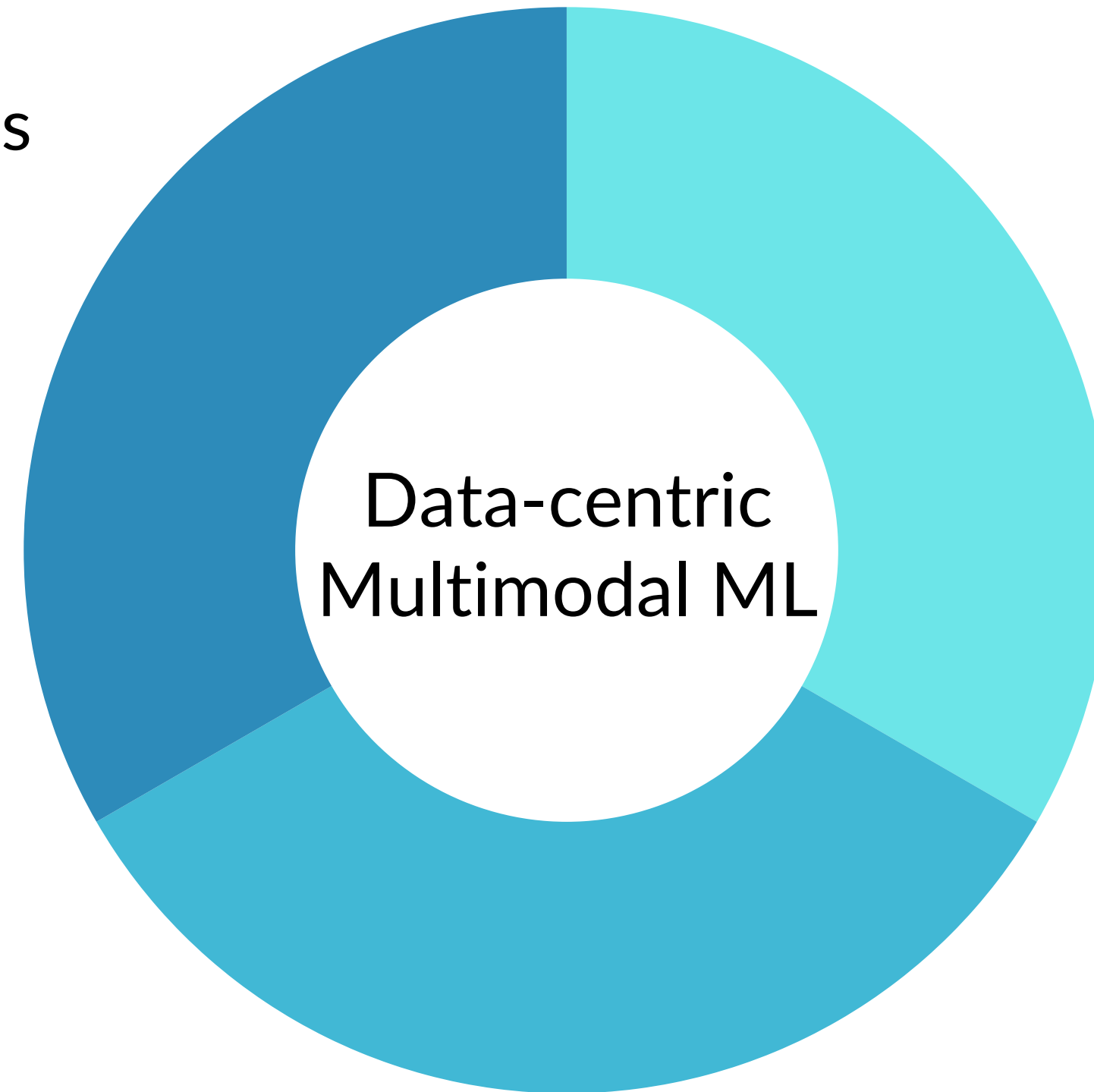
Given a fixed compute budget, is it better to increase the multimodal model size or the multimodal dataset size?



Looking Forward

Scale-driven capabilities

a more rigorous understanding of what increasing the scale does to the multimodal training procedure and how these desirable emergent capabilities come about



Optimal data

What subsets of the data are most important for the capabilities of foundation multimodal models?
What kind of synthetic multimodal data we can build?

Fundamental scaling and limits

Given a fixed compute budget, is it better to increase the multimodal model size or the multimodal dataset size?

Summary

- In this work, we propose the first vision-language dataset distillation method.
- Our experiments show that co-distilling different modalities via bi-trajectory matching holds promise.
- We hope that the insights we gathered can serve as roadmap for future studies exploring more complex settings.



Byron Zhang



Zhiwei Deng



Olga Russakovsky



[Website](#)



[arXiv](#)



[Code](#)