

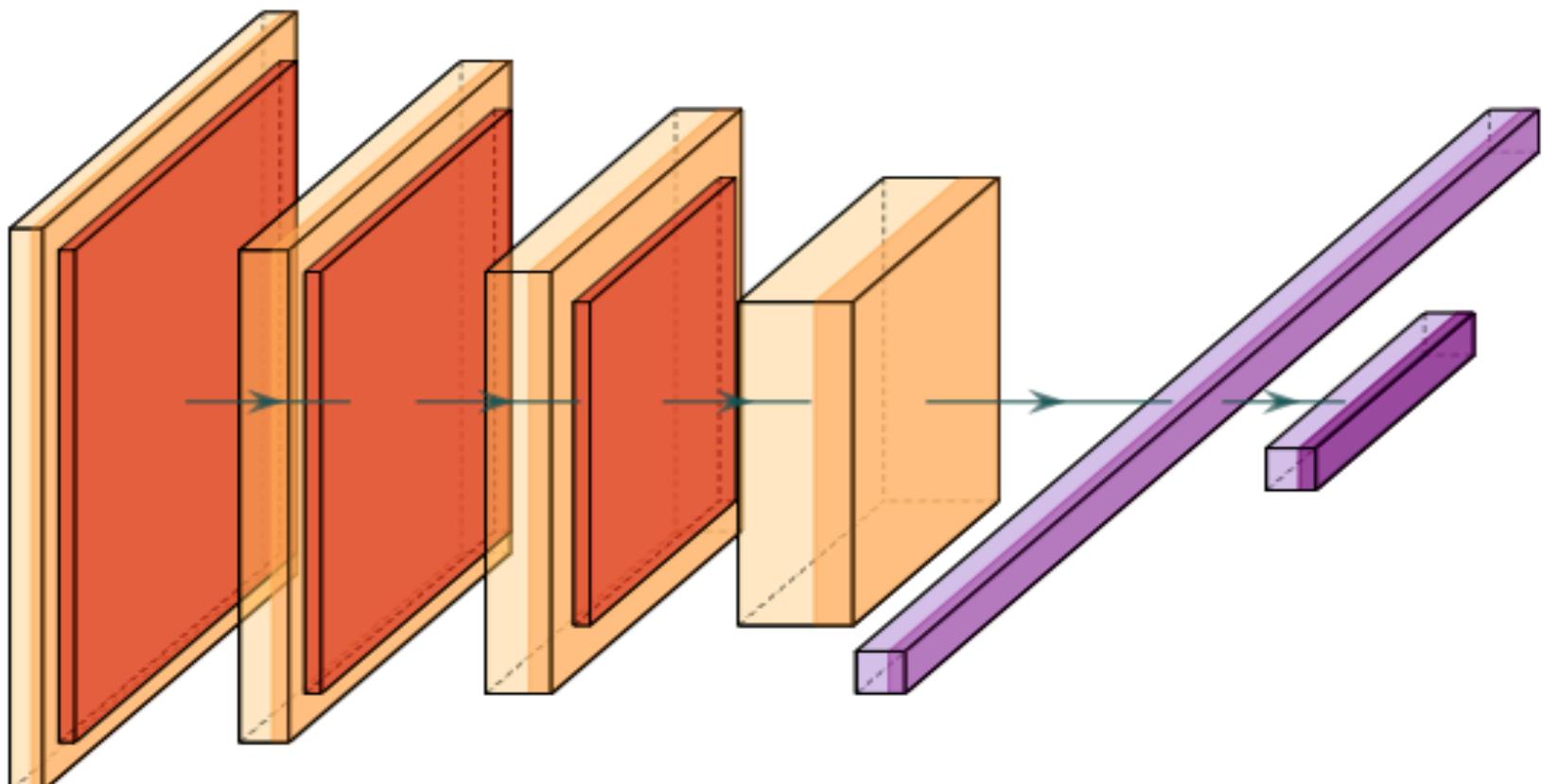
# PGrad: Learning Principal Gradients for Domain Generalization

Zhe Wang, Jake Grigsby, Yanjun Qi

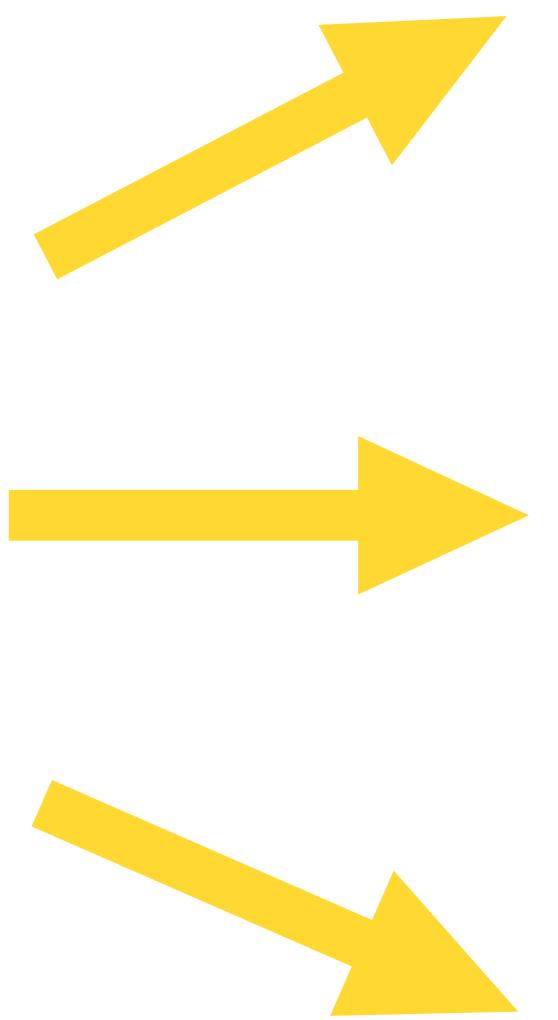
University of Virginia



# Domain Generalization Challenge



$$f_{\theta^*}$$



$T_1$

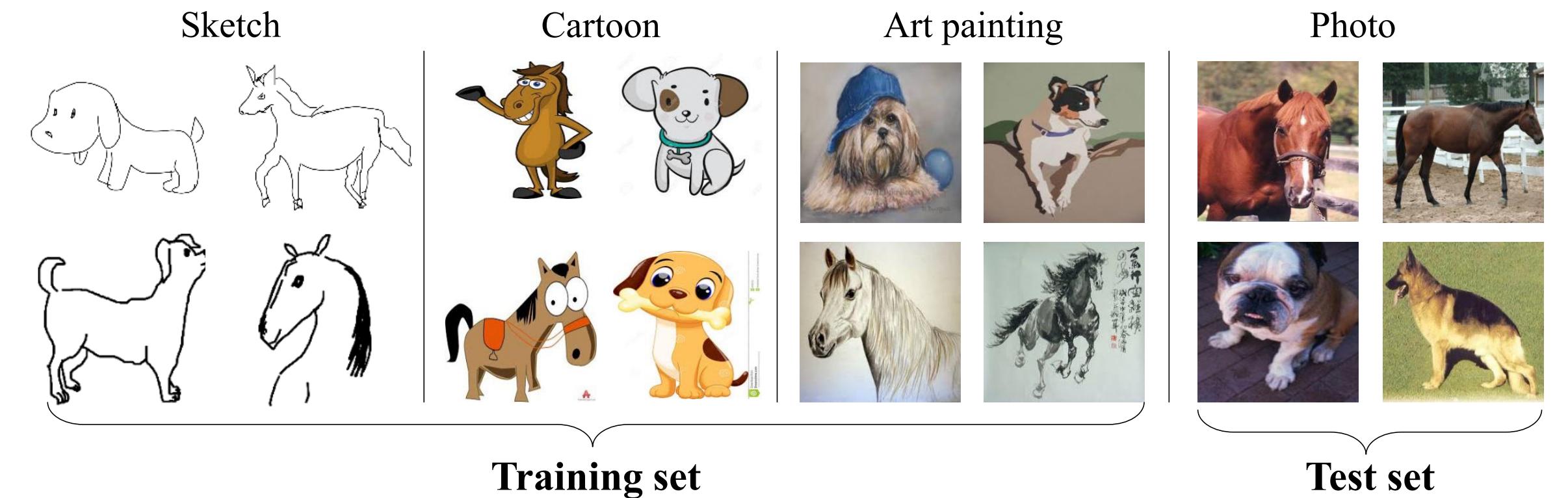
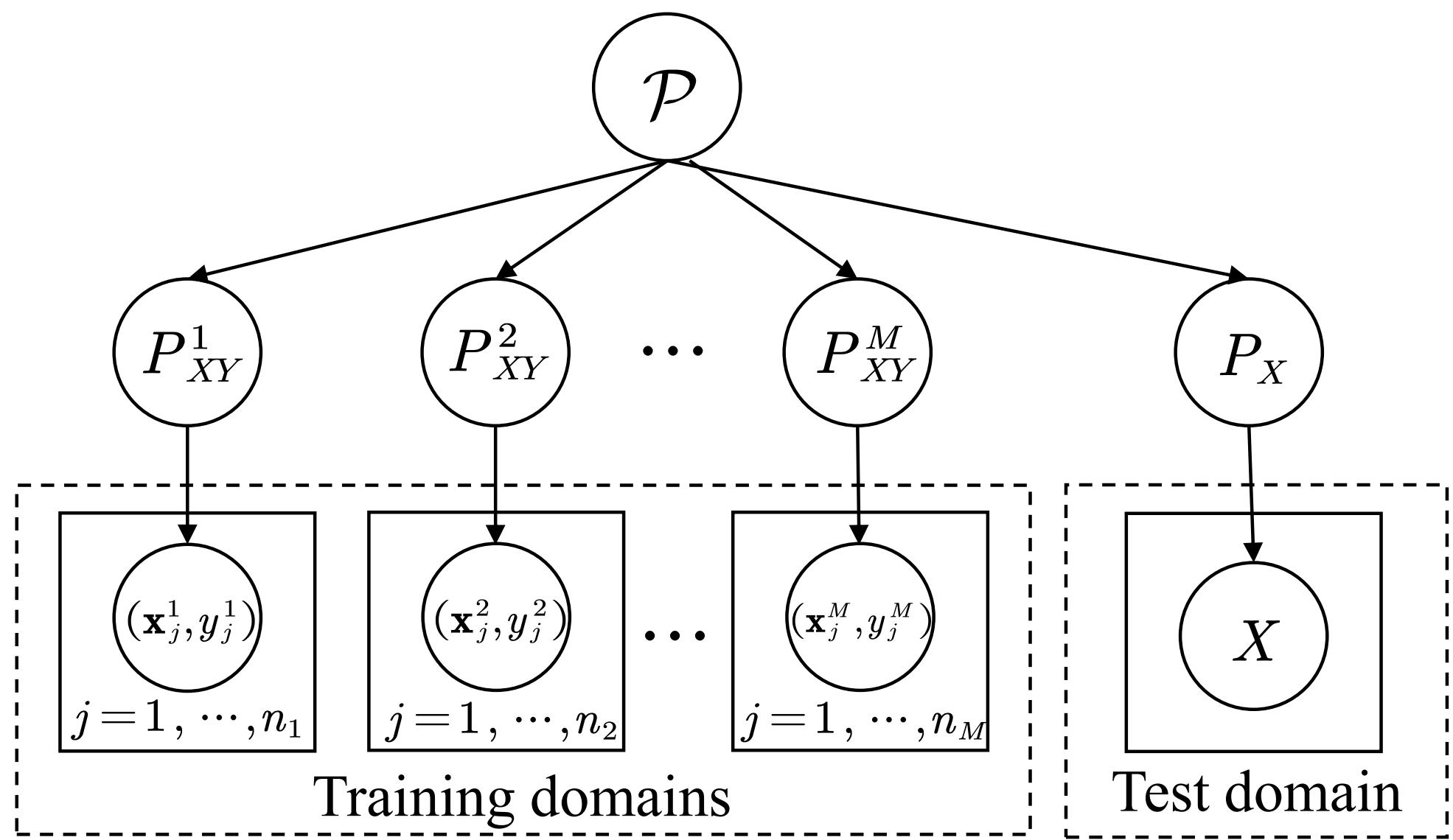


$T_2$



$T_3$

# Domain Generalization Challenge



# Math Formulation

**Domain generalization** assumes no access to instances from future test domains.

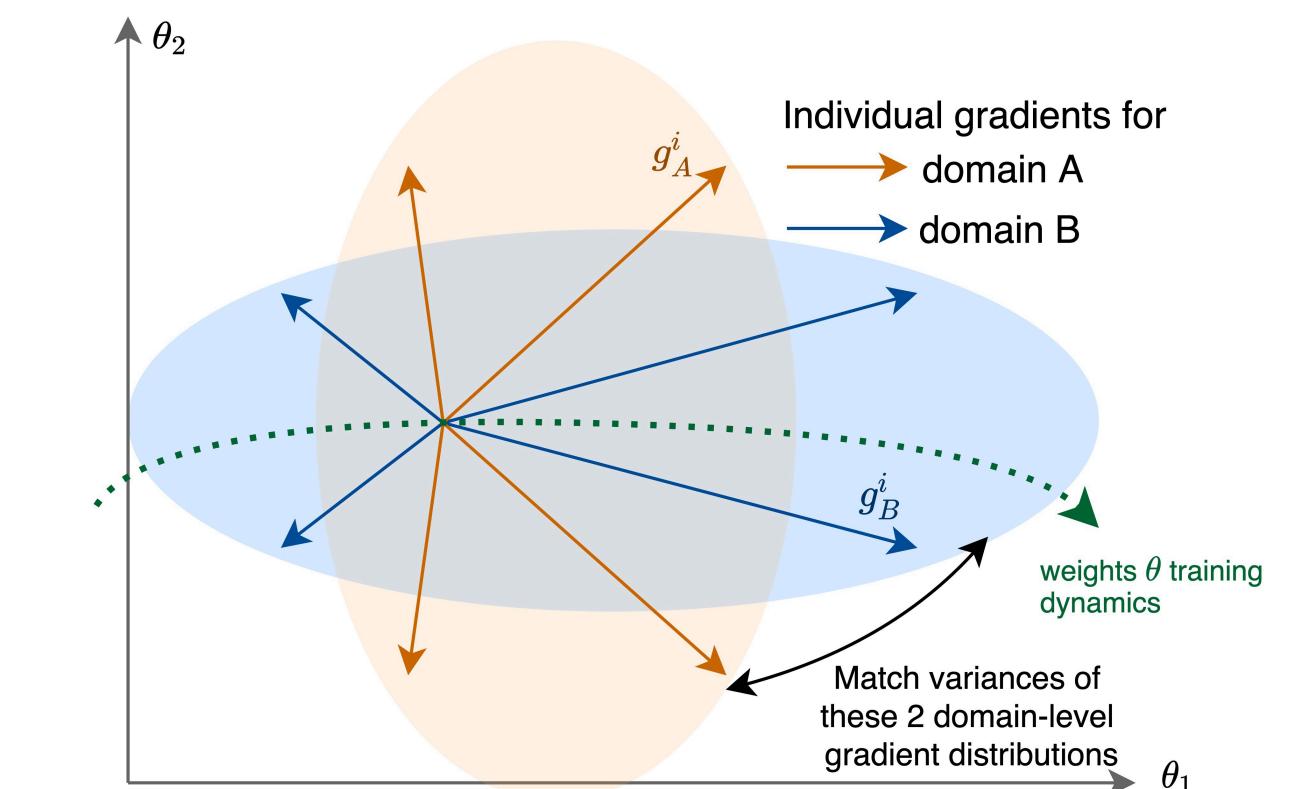
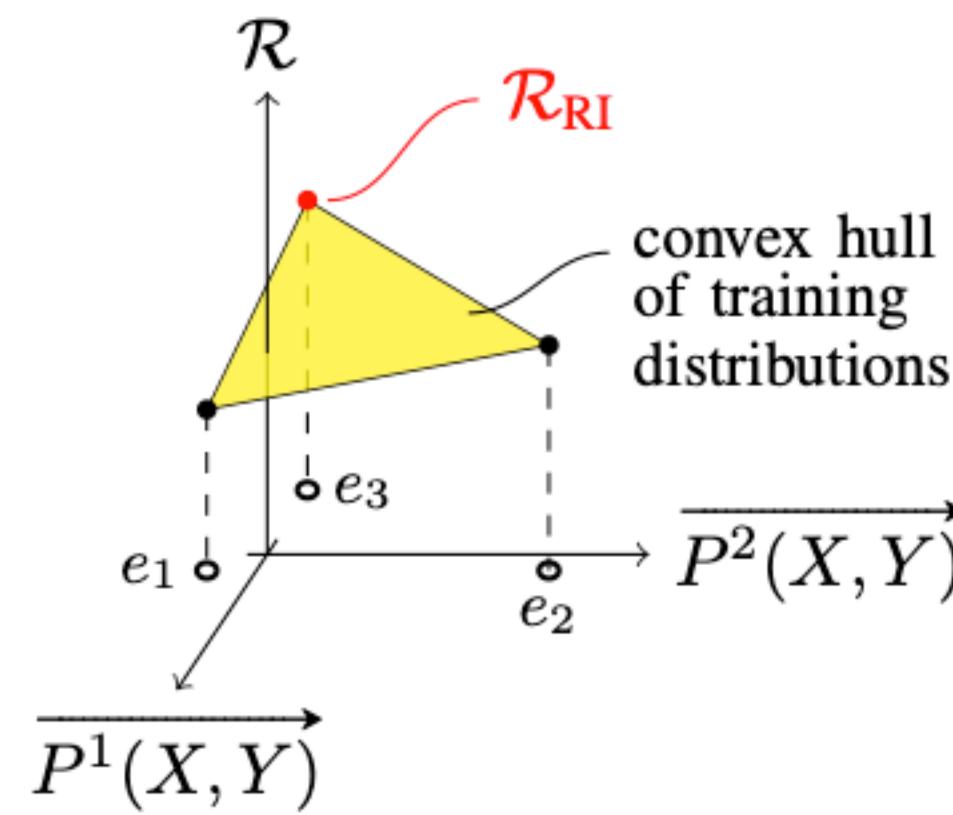
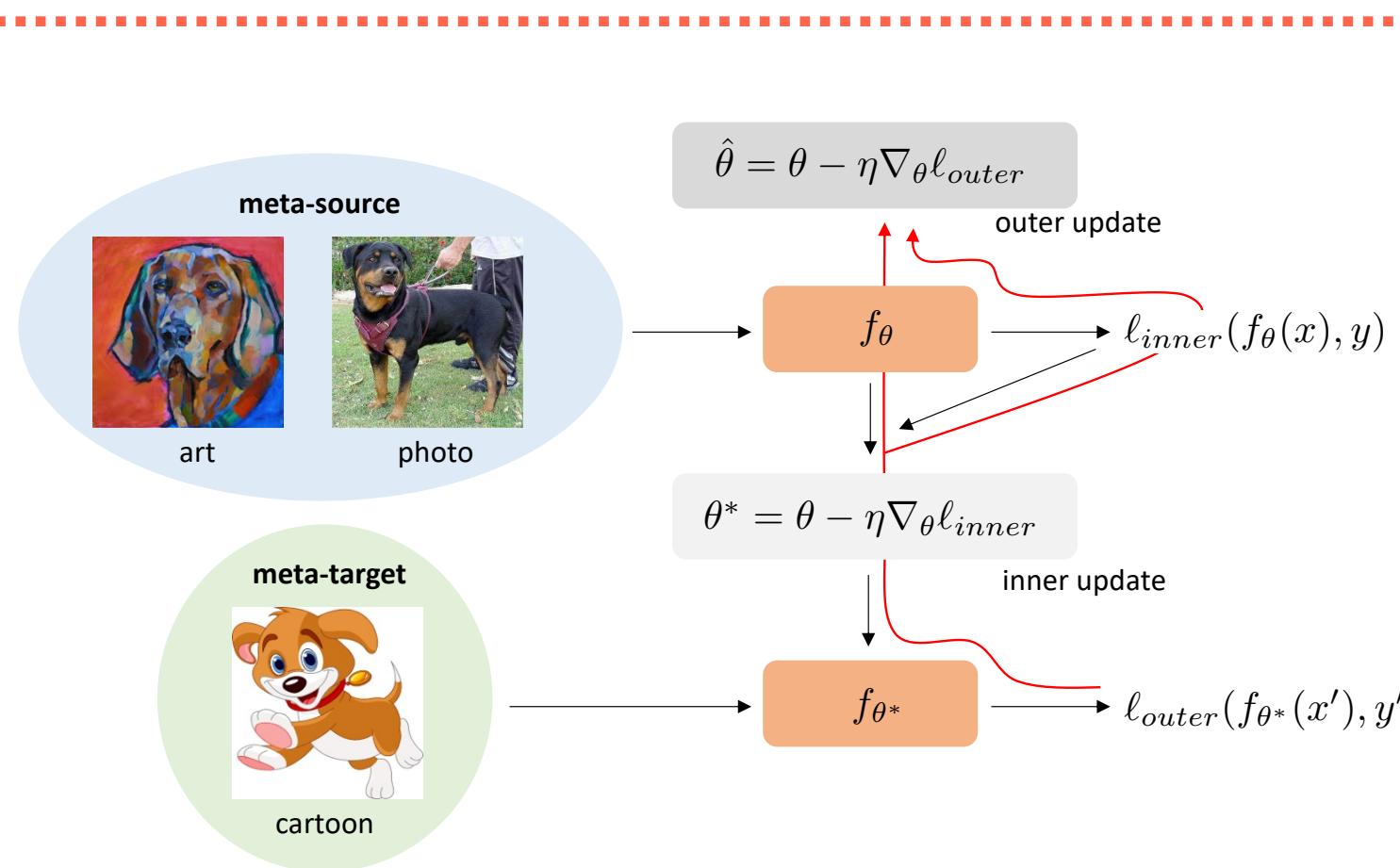
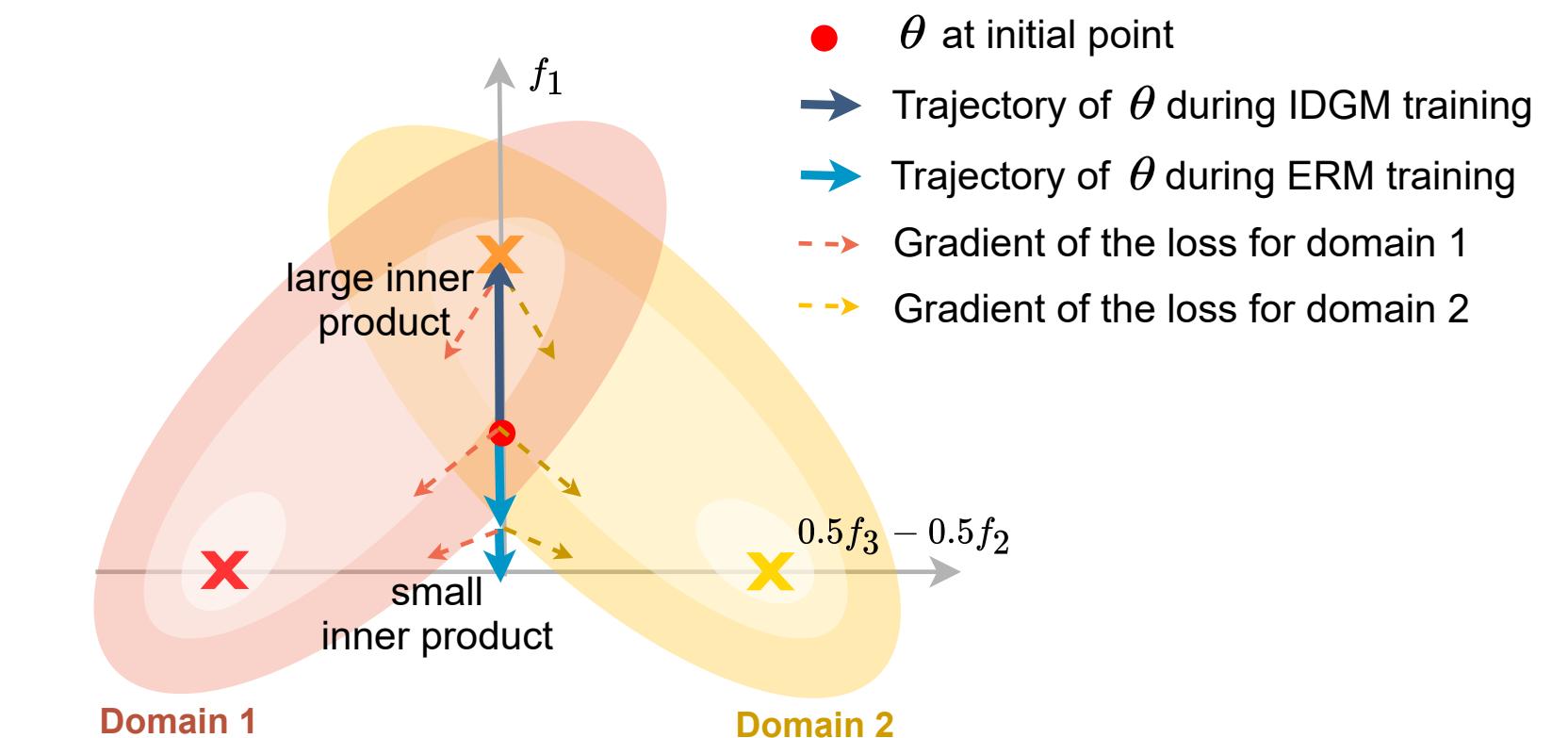
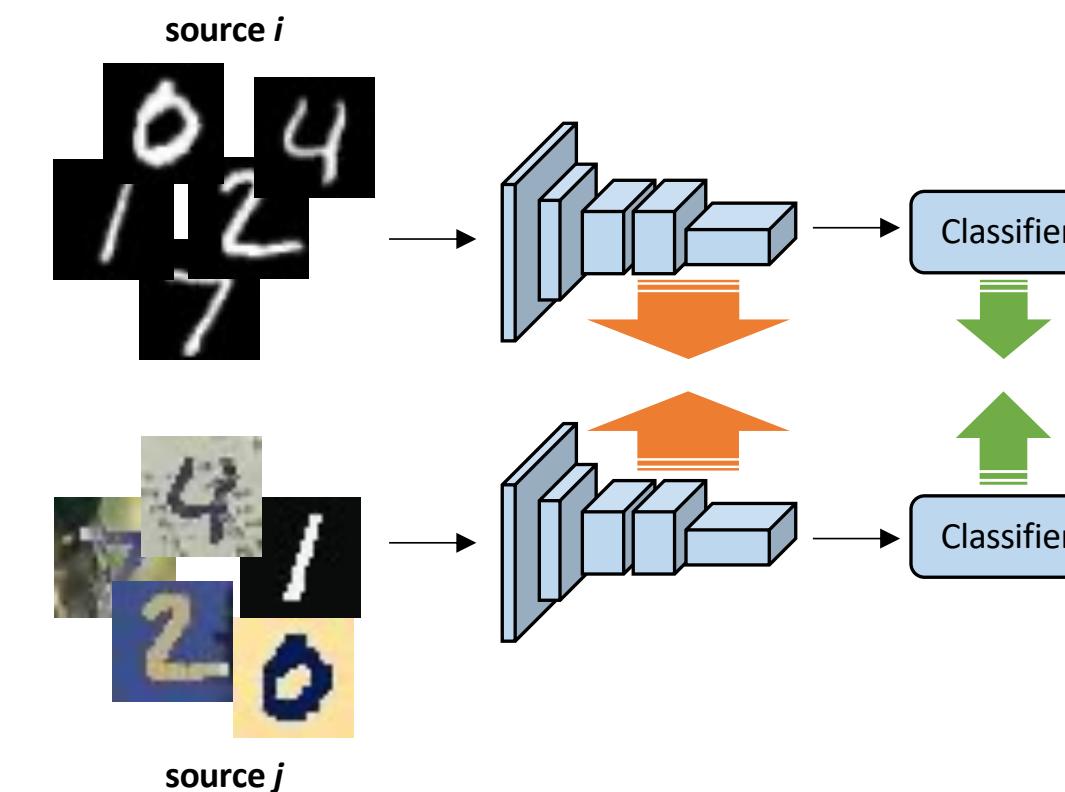
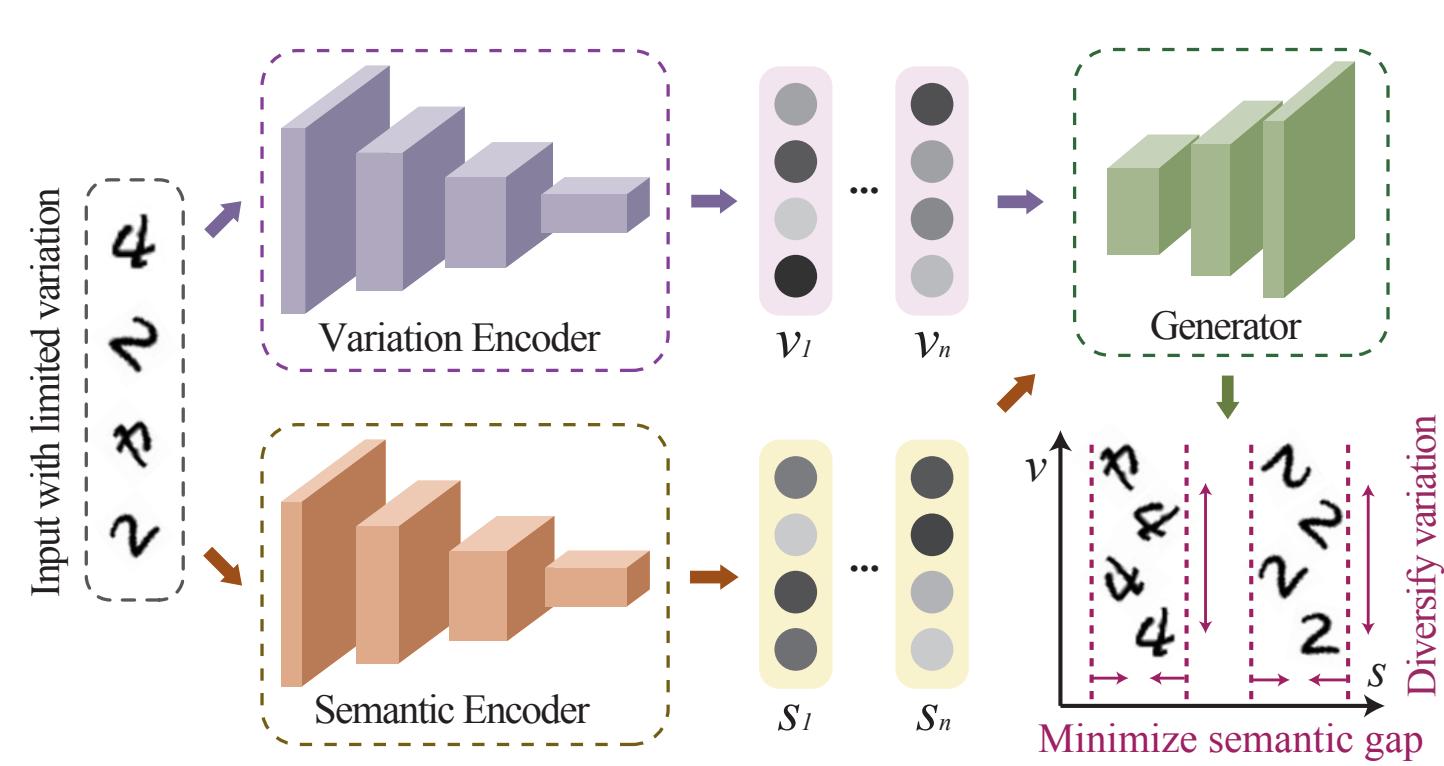
**Training:**

- We have a set of training domains  $\mathcal{D}_{tr} = \{D_i\}_{i=1}^n$ ;
- $D_i$  is characterized by i.i.d samples  $\{\mathbf{x}_k^i, \mathbf{y}_k^i\}$
- $D_i$  is associated with a joint distribution  $P_{X \times Y}^{D_i}$ , and  $P_{X \times Y}^{D_i} \neq P_{X \times Y}^{D_j}$

**Testing:**

- we evaluate the model on test domains  $\mathcal{T}_{te} = \{T_j\}_{j=1}^m$ , where  $P_{X \times Y}^{D_i} \neq P_{X \times Y}^{T_j}$

# Previous Work: Feature, Optimization, and Beyond



Piratla, et.al, "Efficient domain generalization via common-specific low-rank decomposition," in ICML, 2020.

Li et.al, "Learning to generalize: Meta-learning for domain generalization," in AAAI, 2018.

Shi et.al, "Gradient Matching for Domain Generalization," in ICLR, 2022.

Chattopadhyay et.al "Learning to balance specificity and invariance for in and out of domain generalization," in ECCV, 2020.

Krueger et.al "Out-of-Distribution Generalization via Risk Extrapolation (REx)" in Arxiv, 2021.

Rame et.al, "Fishr: Invariant Gradient Variances for Out-of-Distribution Generalization," in ICML 2022.

# Limitations

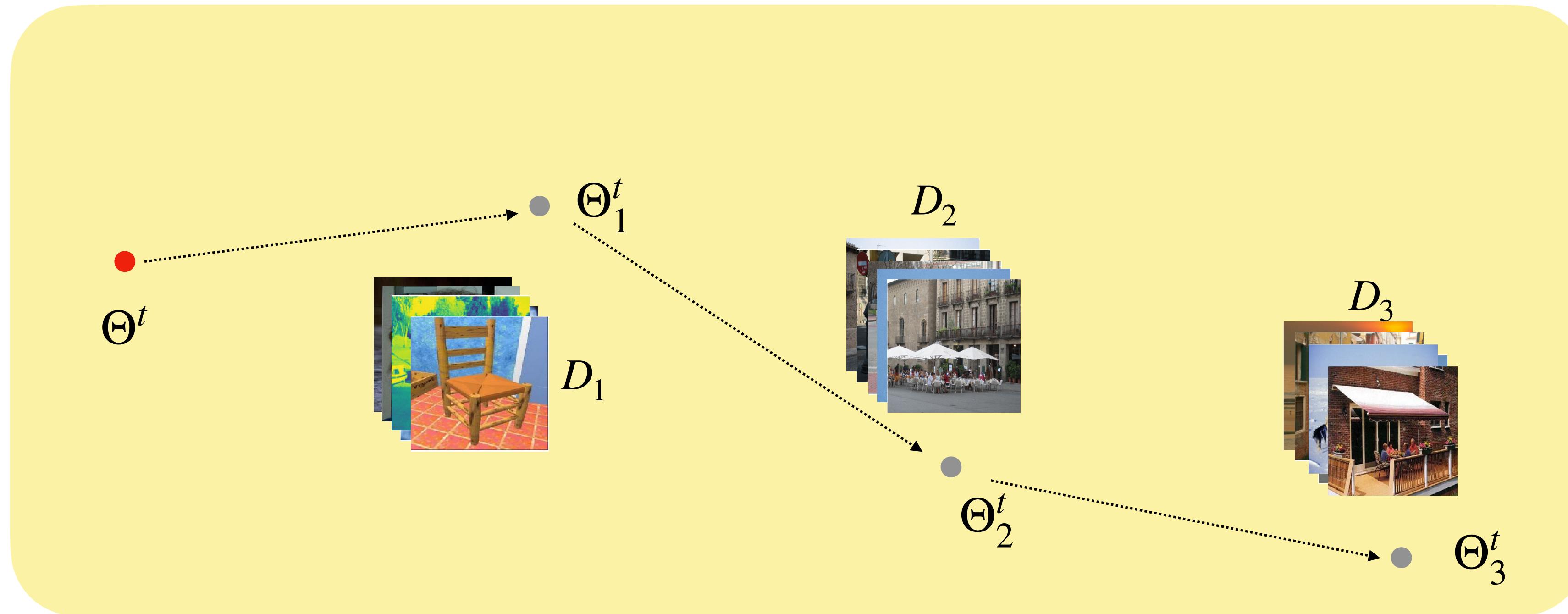
- Unsatisfied Empirical Performance;
  - *Empirical Risk Minimization* remains the strong baseline.
- Strong Distributional Assumption;
  - Invariant marginal or conditional distributions.
- Dataset or Task Dependent.

Gulrajani et.al, “In Search Of Lost Domain Generalization,” in ICLR, 2021.

Galstyan et.al “Failure Mode of Domain Generalization Algorithms” in CVPR, 2022.

Blanchard et.al “Generalizing from several related classification tasks to a new unlabeled sample”, in NeurIPS, 2011

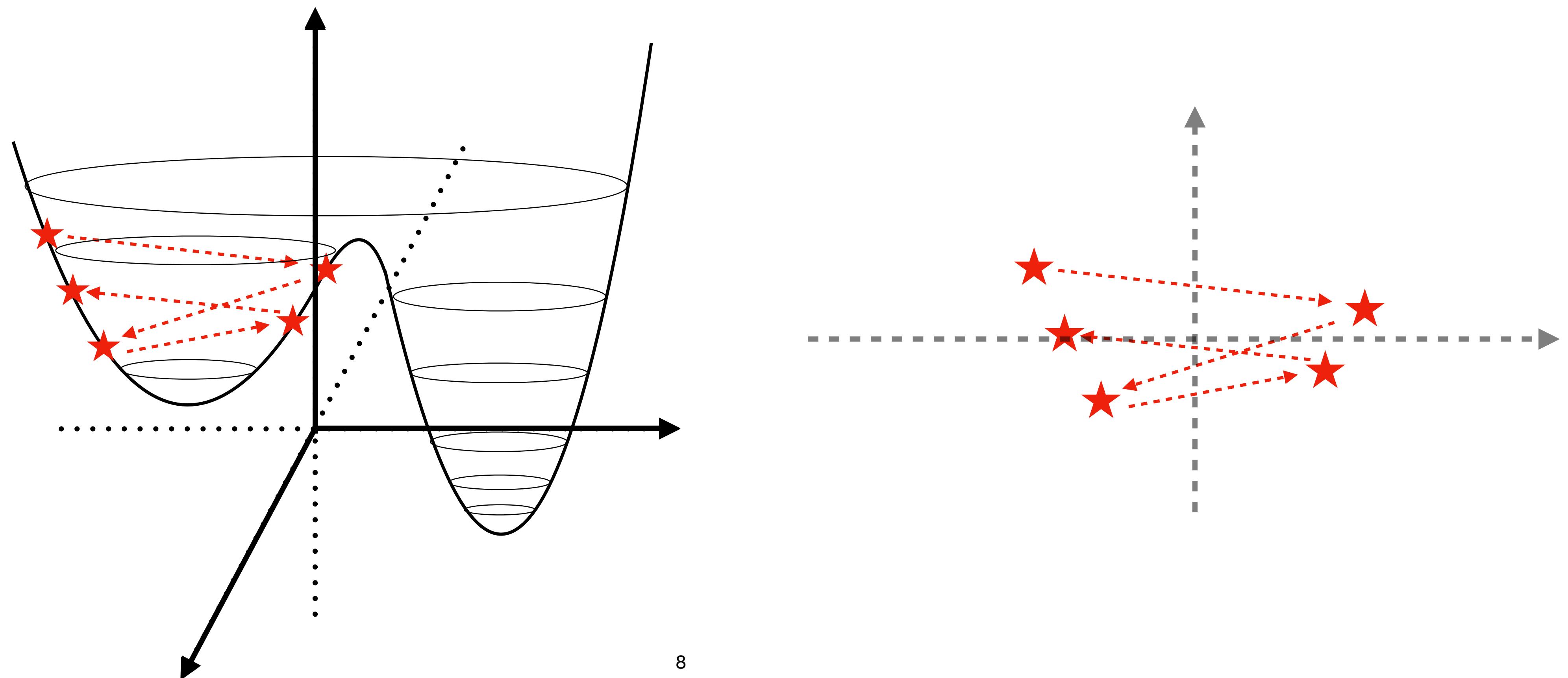
# Proposed: Trajectory Based Direction Learning



Learning From Optimization Trajectories

# Why Trajectory Based?

- Disentangle dominant directions for analysis. For example:
  - Zeroing out the insignificant directions.
  - Adding weights for significant directions.



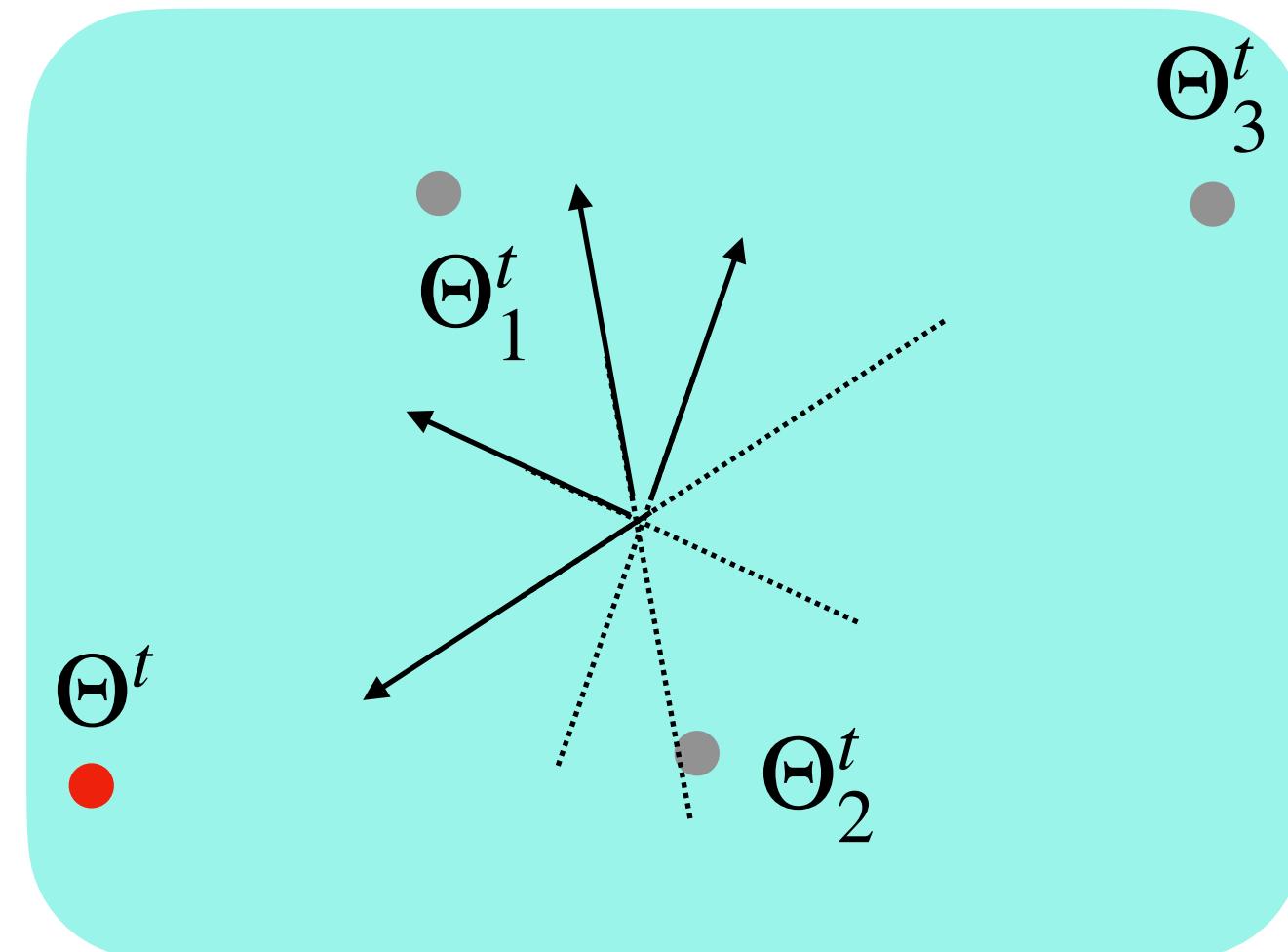
# Principal Direction Learning (Our Main Idea)

1. Trajectory Sampling;

$$S = \Theta^t \rightarrow \Theta_1^t \rightarrow \Theta_2^t \rightarrow \dots \rightarrow \Theta_n^t.$$

2. Local Principal Coordinate Construction

$$\max_{V_z} \text{Var}([Sv_z]), \quad \text{s.t. } V^T V = \mathbf{I}_d.$$



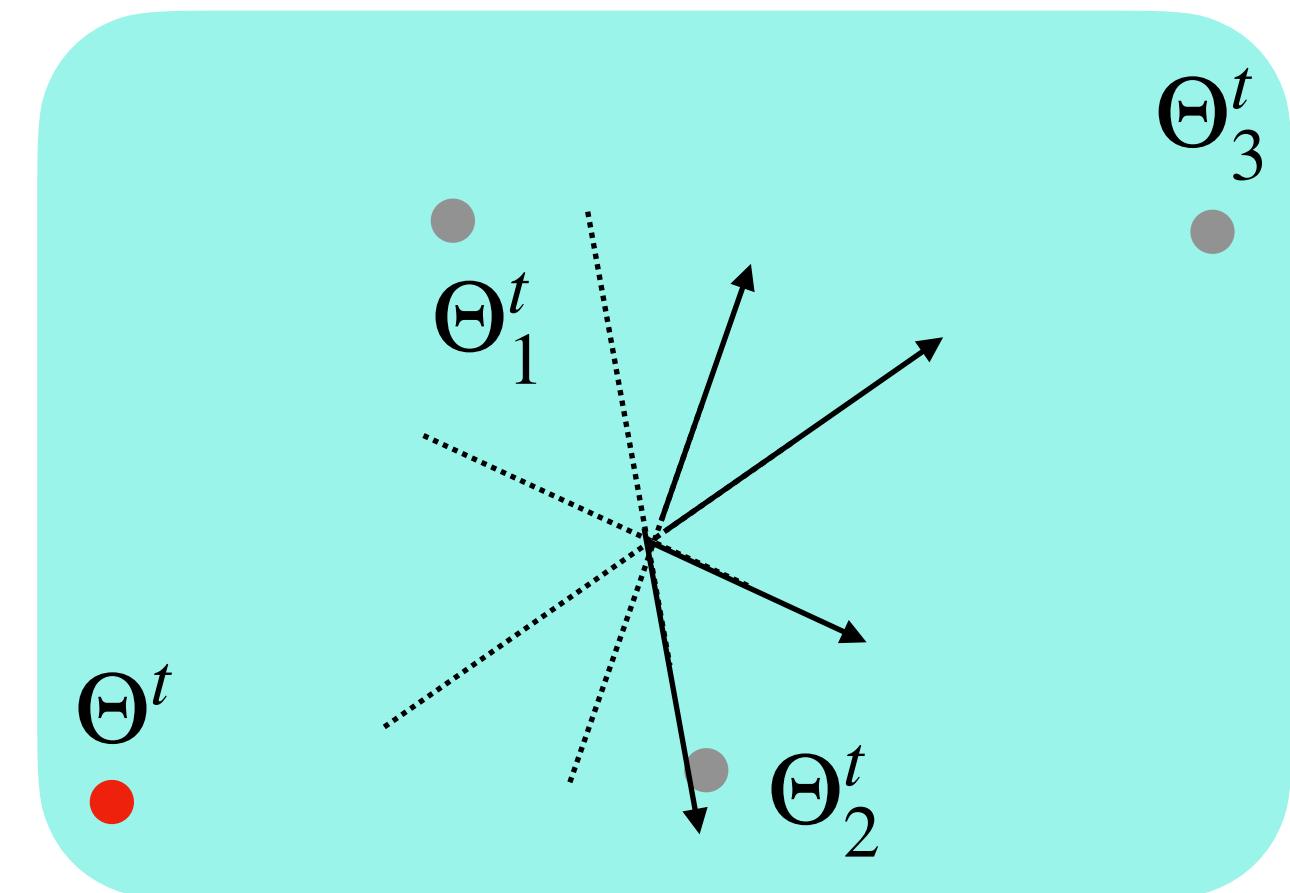
# Principal Direction Learning (Our Main Idea)

## 3. Direction Calibration

Teach the model to climb down the loss landscape with:

$$\nabla_r = \Theta^t - \Theta_n^t.$$

$$\mathbf{w}_z = r_z \mathbf{v}_z, \quad r_z = \begin{cases} 1, & \text{if } \langle \mathbf{v}_z, \nabla_r \rangle \geq 0 \\ -1, & \text{otherwise} \end{cases}$$



## 4. Principal Direction Learning

$$\nabla_p = \sum_{z=0}^n \frac{\lambda_z}{\|\lambda\|_2} w_z. \quad \Rightarrow \quad \|\nabla_p\| = 1.$$

# Principal Direction Learning (Our Main Idea)

3. Direction Calibration

$$\nabla_r = \Theta^t - \Theta_n^t.$$

4. Principal Direction Learning

$$\nabla_p = \sum_{z=0}^n \frac{\lambda_z}{\|\lambda\|_2} w_z. \implies \|\nabla_p\| = 1.$$



5. Optimization Process Aware LR Tuning

$$\nabla_p = \sum_{z=0}^k \frac{\lambda_z \|\nabla_r\|}{\|\lambda\|_2} w_z. \implies \text{As Training Continues, } \|\nabla_r\| \rightarrow 0$$

## Principal Direction Learning (Computational Challenge in Step 2)

$$\max_{V_z} \text{Var}([Sv_z]), \quad \text{s.t. } V^T V = \mathbf{I}_d.$$

- Closed form solution:

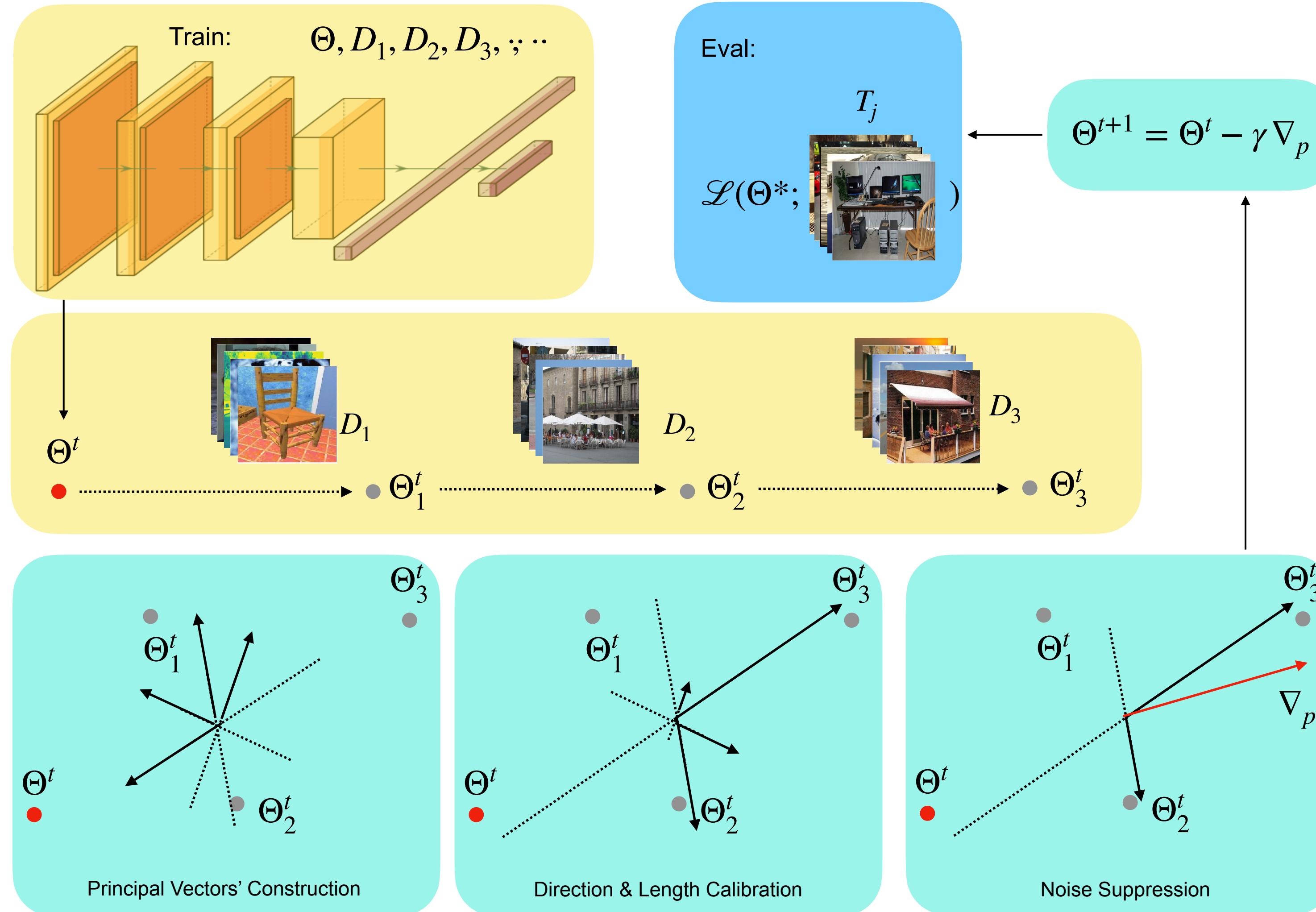
$$\lambda_z, \mathbf{v}_z = \text{SVD}_z\left(\frac{1}{n}\hat{S}^T \hat{S}\right), \quad \text{where } \hat{S}^T \hat{S} \in \mathbb{R}^{d \times d}$$

- **Challenge**: The computational complexity of the SVD  $\mathcal{O}(d^3)$ , where  $d \sim \mathcal{O}(10^{10})$
- **Solution**: Transpose trick.

$$\hat{S} \hat{S}^T \mathbf{e}_z = \lambda_z \mathbf{e}_z \implies \hat{S}^T \hat{S} \hat{S}^T \mathbf{e}_z = \lambda_z \hat{S}^T \mathbf{e}_z \implies \mathbf{v}_z = \hat{S}^T \mathbf{e}_z$$

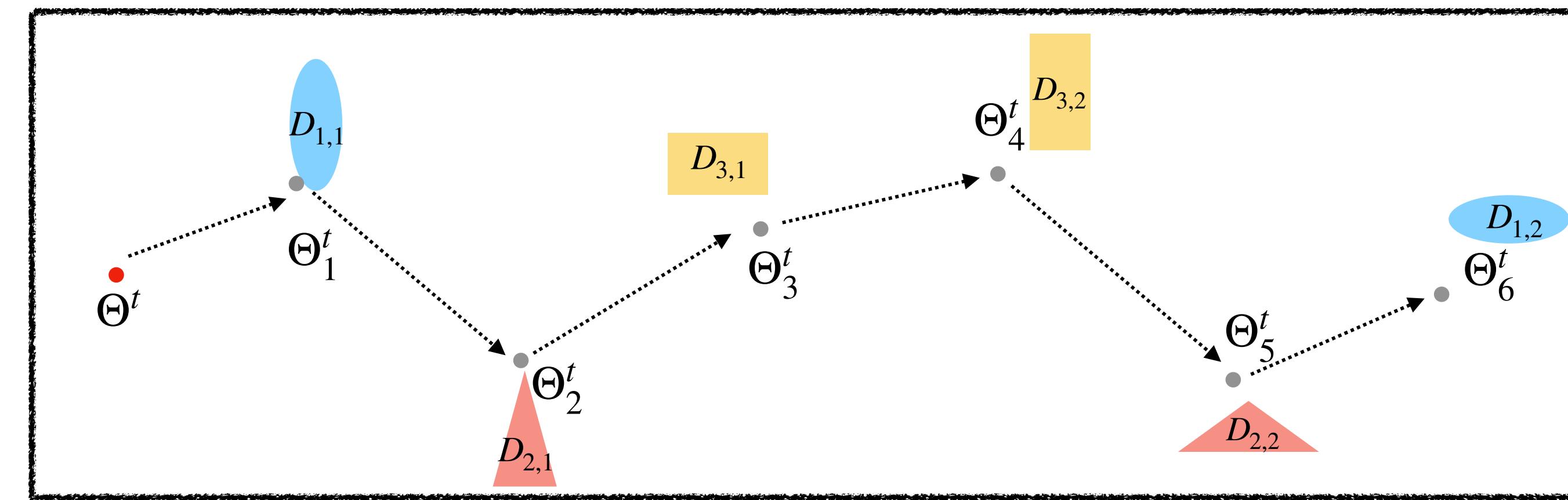
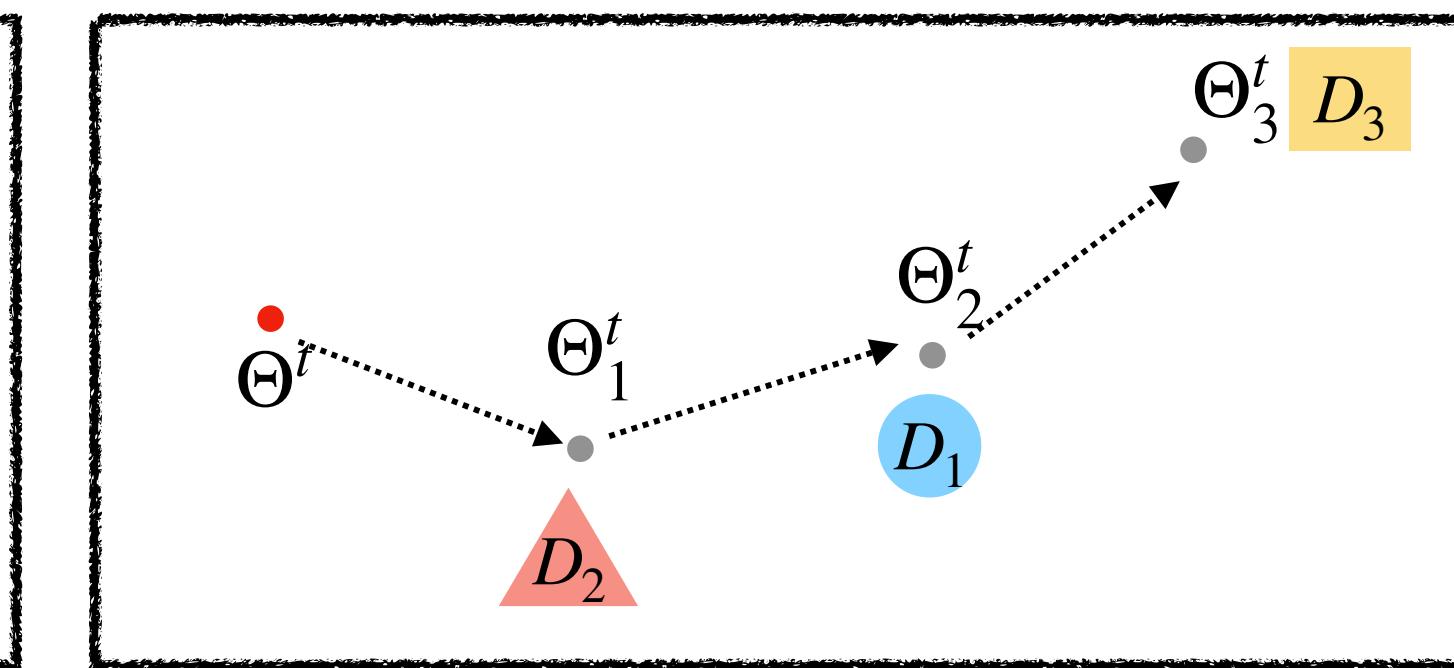
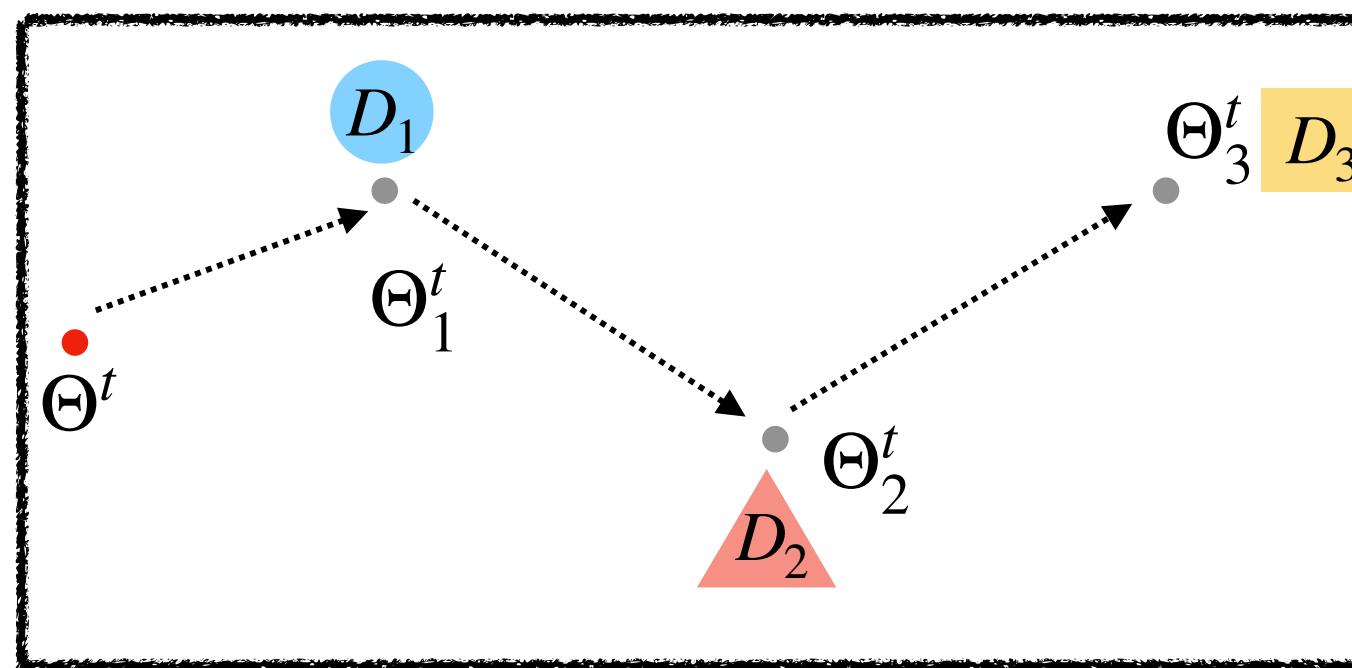
- **Analysis**: The computational complexity is  $\mathcal{O}(n^3)$ , where  $n = 3, 4, \dots$

# PGrad: Learning Principal Gradients for Domain Generalization



# Principal Direction Learning (Optional)

- Can we sample more flexible and longer trajectories?



## PGrad & Variants

- *PGrad*;
  - Random shuffling the domain order during training.
- *PGrad-B*(default);
  - Split each training batch into  $B$  smaller batches.
- *PGrad-BMix*
  - Combine PGrad-B and MixUp

Zhang et.al, “Mixup: Beyond Empirical Risk Minimization,” in ICLR, 2018.

## One of the Benchmark Suits – DomainBed

- Image Dataset - Classification - ResNet-50

Table 1: A summary on DOMAINBED dataset, metrics, and architectures we used.

Dataset	# of Images	Domains	# of Classes
PACS (Li et al., 2017)	9,991	Artpaint, Cartoon, Sketches, Photo	7
VLCS (Fang et al., 2013)	10,729	PASCAL VOC 2007, LabelMe, Caltech, Sun	5
OFFICEHOME (Venkateswara et al., 2017)	15,588	Art, Clipart, Product, Real-World	65
TERRAINCOGNITA (Beery et al., 2018)	24,788	Location #100, #38, #43, #46	10
DOMAINNET (Peng et al., 2019)	586,575	Clipart, Infograph, Painting, Quickdraw, Real, Sketch	345

- Evaluation protocol
  - We select  $k - 1$  as training domains and the remaining for testing.
  - For each combination, we repeat 3(random seeds)  $\times$  2(hyperparameters) running.

# DomainBed Results – Prediction Acc Improvement

Table 2: Test accuracy (%) on five datasets from the DomainBed benchmark. We group 20% data from each training domain to construct validation set for model selection.

Categories	Algorithms	VLCS	PACS	OfficeHome	TerraInc	DomainNet	Avg
Baseline	ERM	$77.5 \pm 0.4$	$85.5 \pm 0.2$	$66.5 \pm 0.3$	$46.1 \pm 1.8$	$40.9 \pm 0.1$	63.3
Invariant	IRM	$78.5 \pm 0.5$	$83.5 \pm 0.8$	$64.3 \pm 2.2$	$47.6 \pm 0.8$	$33.9 \pm 2.8$	$61.6 \textcolor{red}{-1.7}$
	MMD	$77.5 \pm 0.9$	$84.6 \pm 0.5$	$66.3 \pm 0.1$	$42.2 \pm 1.6$	$23.4 \pm 9.5$	$58.8 \textcolor{red}{-4.5}$
	DANN	$78.6 \pm 0.4$	$83.6 \pm 0.4$	$65.9 \pm 0.6$	$46.7 \pm 0.5$	$38.3 \pm 0.1$	$62.6 \textcolor{gray}{-0.7}$
	CORAL	$78.8 \pm 0.6$	$\underline{86.2} \pm 0.3$	$68.7 \pm 0.3$	$47.6 \pm 1.0$	$41.5 \pm 0.1$	$64.5 \textcolor{blue}{+1.2}$
Optimization	GroupDRO	$76.7 \pm 0.6$	$84.4 \pm 0.8$	$66.0 \pm 0.7$	$43.2 \pm 1.1$	$33.3 \pm 0.2$	$60.7 \textcolor{red}{-2.6}$
	MLDG	$77.2 \pm 0.4$	$84.9 \pm 1.0$	$66.8 \pm 0.6$	$47.7 \pm 0.9$	$41.2 \pm 0.1$	$63.6 \textcolor{gray}{+0.3}$
Augmentation	MixUp	$77.4 \pm 0.6$	$84.6 \pm 0.6$	$68.1 \pm 0.3$	$47.9 \pm 0.8$	$39.2 \pm 0.1$	$63.4 \textcolor{gray}{+0.1}$
	ARM	$77.6 \pm 0.3$	$85.1 \pm 0.4$	$64.8 \pm 0.3$	$45.5 \pm 0.3$	$35.5 \pm 0.2$	$61.7 \textcolor{red}{-1.6}$
Gradient Manipulation	Fish	$77.8 \pm 0.3$	$85.5 \pm 0.3$	$68.6 \pm 0.4$	$45.1 \pm 1.3$	$\textbf{42.7} \pm 0.2$	$63.9 \textcolor{gray}{+0.6}$
	Fishr	$77.8 \pm 0.1$	$85.5 \pm 0.4$	$67.8 \pm 0.1$	$47.4 \pm 1.6$	$41.7 \pm 0.0$	$64.0 \textcolor{gray}{+0.7}$
	PGrad	$\textbf{79.3} \pm 0.3 \textcolor{blue}{+1.8}$	$85.1 \pm 0.3 \textcolor{gray}{-0.4}$	$69.3 \pm 0.1 \textcolor{blue}{+2.8}$	$49.0 \pm 0.3 \textcolor{blue}{+2.9}$	$41.0 \pm 0.1 \textcolor{gray}{+0.1}$	$64.7 \textcolor{blue}{+1.4}$
	PGrad-B	$\underline{78.9} \pm 0.3 \textcolor{blue}{+1.4}$	$\textbf{87.0} \pm 0.1 \textcolor{blue}{+1.5}$	$\underline{69.6} \pm 0.1 \textcolor{blue}{+3.1}$	$49.4 \pm 0.8 \textcolor{blue}{+3.3}$	$41.4 \pm 0.1 \textcolor{gray}{+0.5}$	$65.3 \textcolor{blue}{+2.0}$
	PGrad-BMix	$\underline{78.9} \pm 0.2 \textcolor{blue}{+1.4}$	$\underline{86.2} \pm 0.4 \textcolor{blue}{+0.7}$	$\textbf{69.8} \pm 0.1 \textcolor{blue}{+3.3}$	$50.7 \pm 0.6 \textcolor{blue}{+4.6}$	$42.6 \pm 0.2 \textcolor{blue}{+1.7}$	$\textbf{65.7} \textcolor{blue}{+2.4}$

# Analyzing Training Process – Loss Smoothing

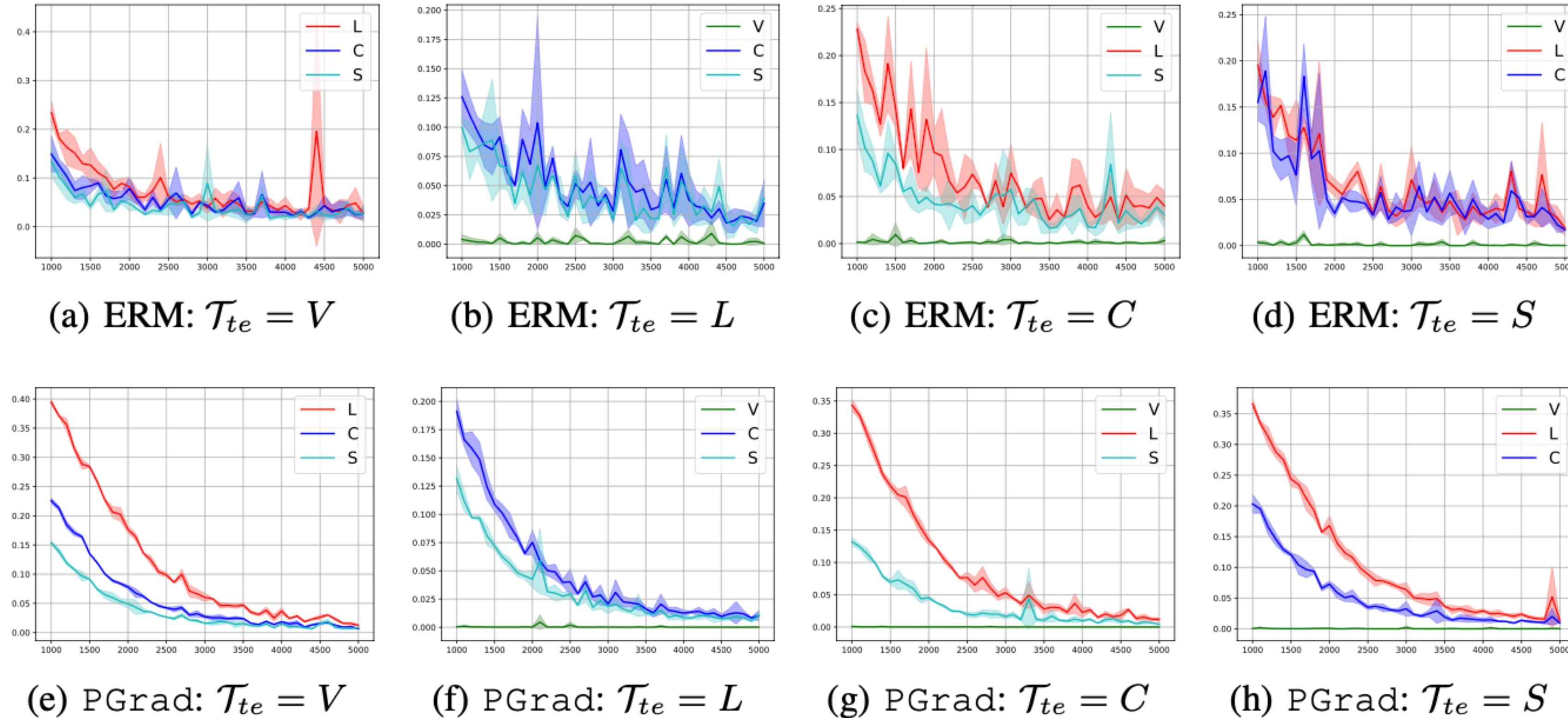
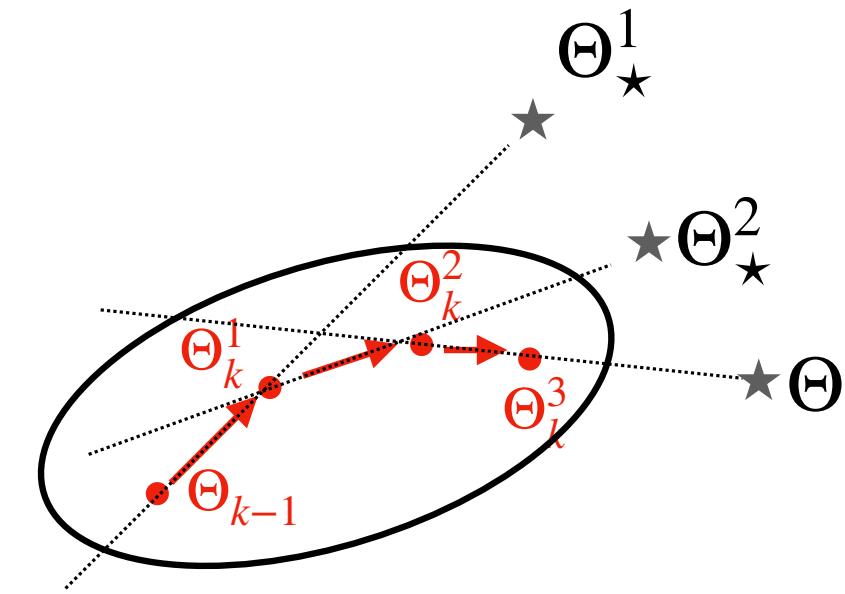
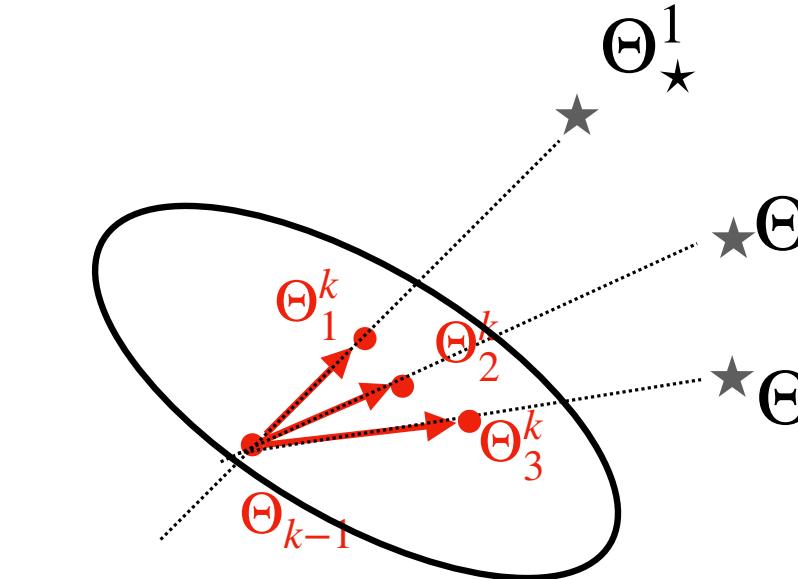


Figure 2: Visualizing domain-wise training losses on VLCS. Curves are the average over 9 runs, surrounded by  $\pm\sigma$  the standard deviation. For comparison, the loss curves start from 1,000 epochs.

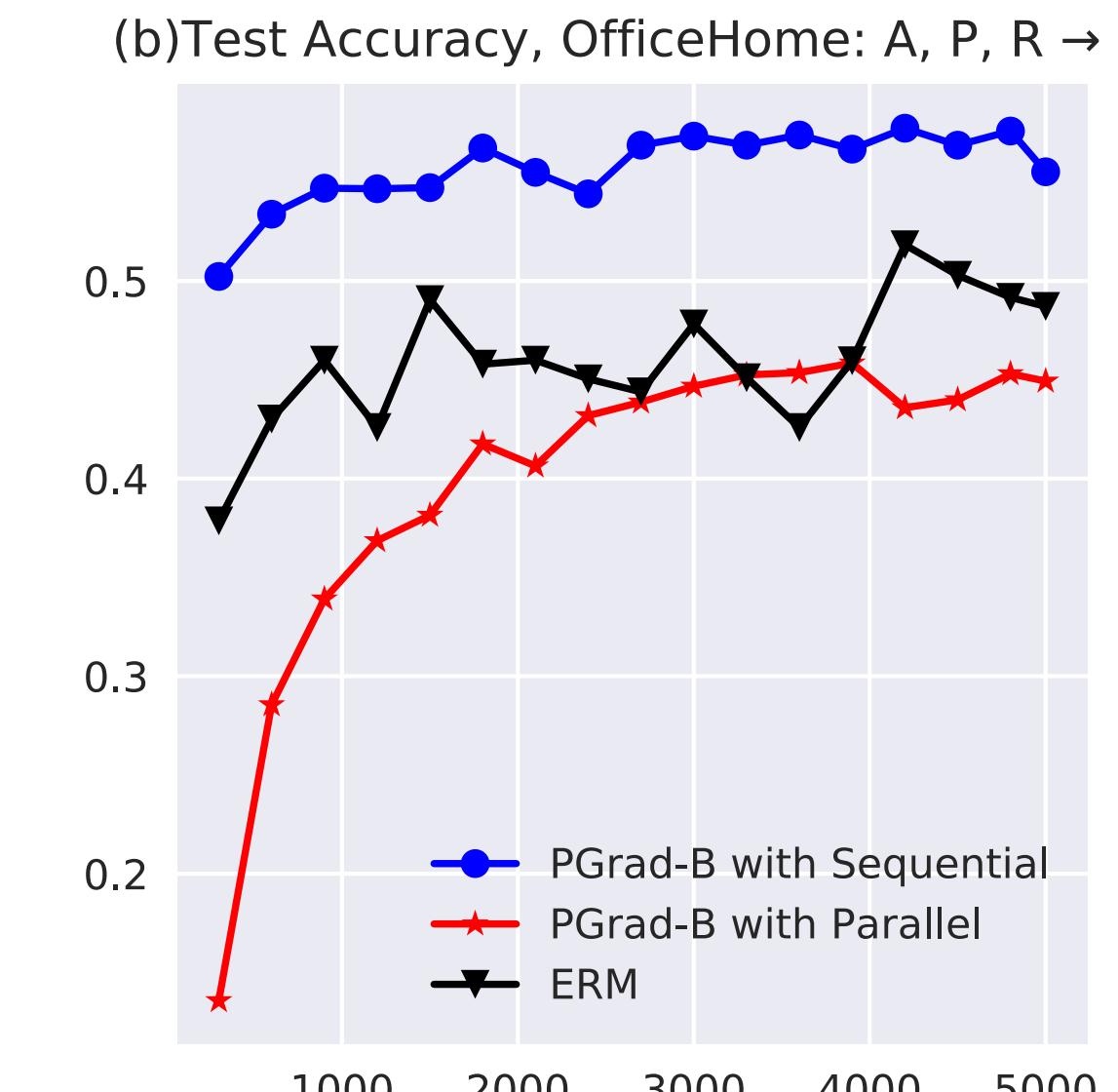
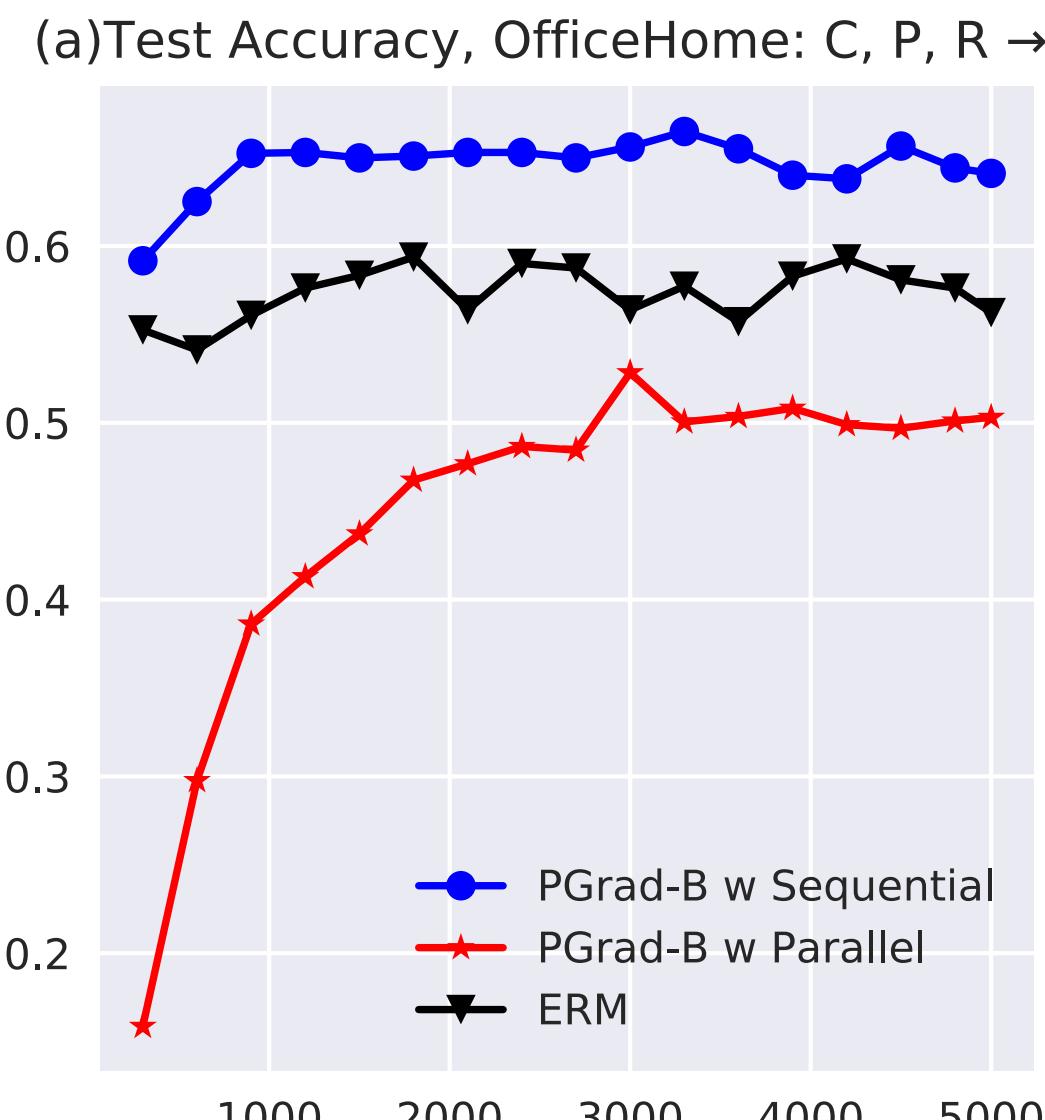
# Comparing Sequential VS. Parallel Training



(a): Sequential training will reinforce robust direction

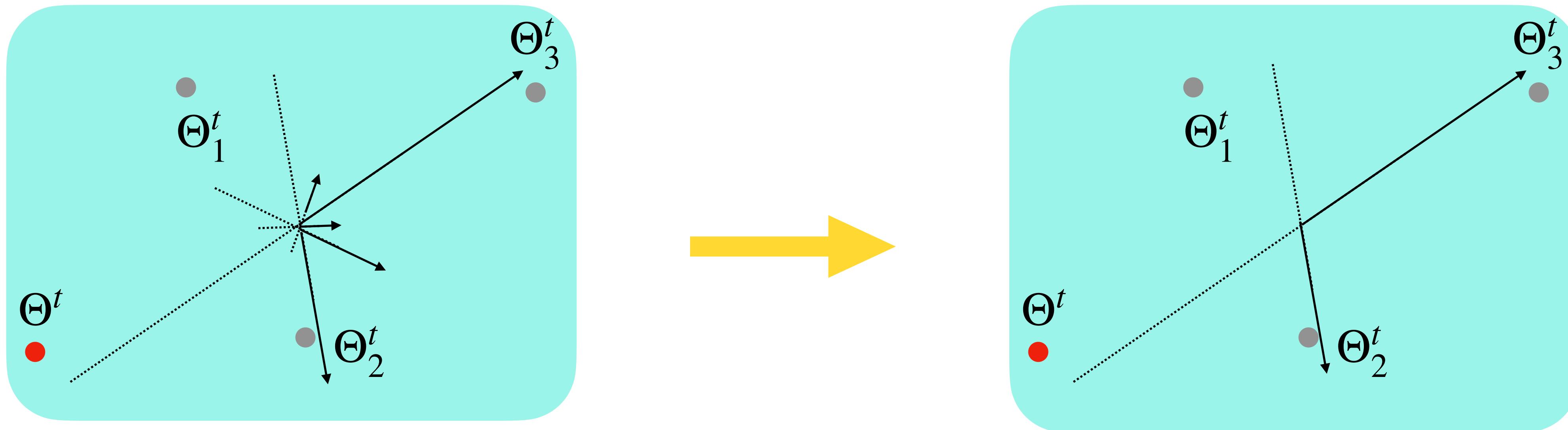


(b): Parallel training will suppress robust direction



# Analyzing PGrad Hyperparameter-k

- Suppressing insignificant principal directions



# How Many Bottom Directions to Suppress?

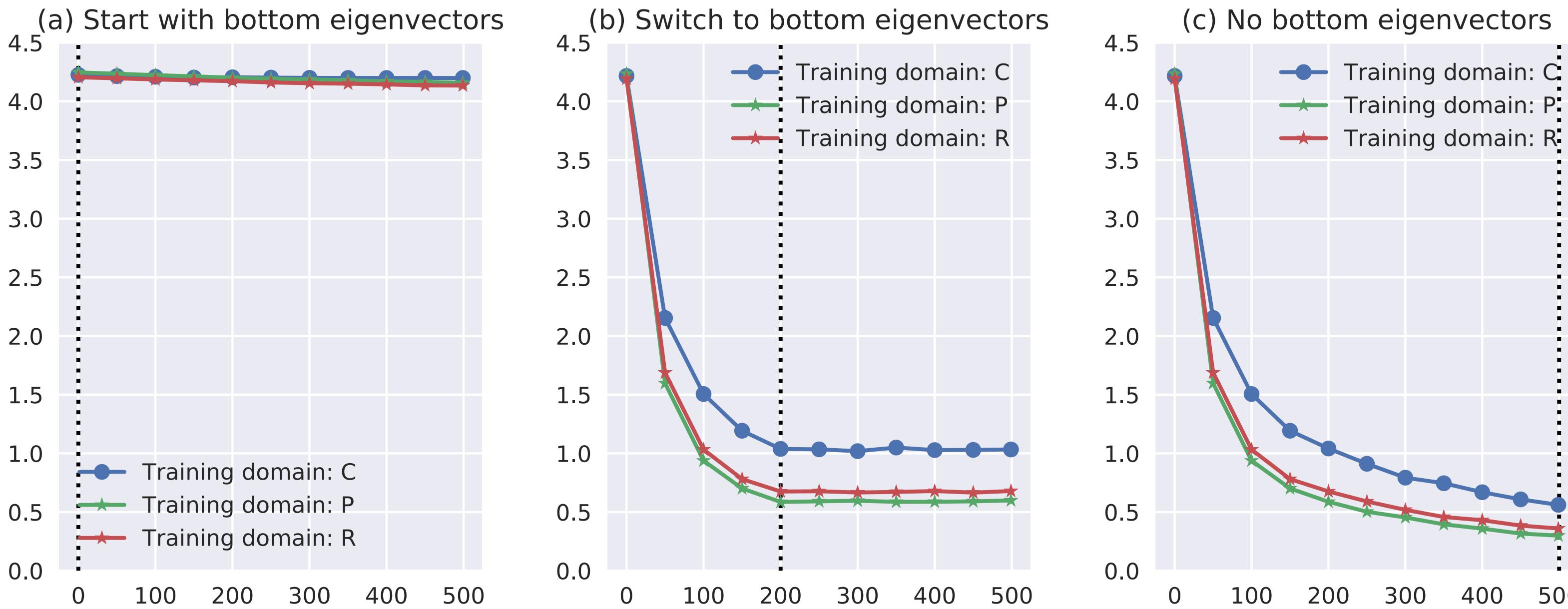


Table 3: Analysis the effect of varying  $k$ . The experiments are performed on PACS dataset. We highlight **first** and **second** best results.

Method	Algorithms	P	A	C	S	Avg
PGrad	$k = 0$	$98.0 \pm 0.2$	$87.3 \pm 0.2$	$76.8 \pm 0.4$	$73.4 \pm 1.3$	83.9
	$k = 2$	$97.8 \pm 0.0$	$87.5 \pm 0.3$	$78.2 \pm 0.8$	$74.0 \pm 1.5$	84.4
	$k = 3$	$97.8 \pm 0.0$	$87.8 \pm 0.4$	$78.4 \pm 0.6$	$77.2 \pm 1.1$	85.3
	$k = 4$	$97.4 \pm 0.1$	$87.6 \pm 0.3$	$79.1 \pm 1.0$	$76.3 \pm 1.2$	85.1
PGrad-B	$k = 0$	$97.5 \pm 0.1$	<u><b><math>89.1 \pm 0.8</math></b></u>	<u><b><math>80.3 \pm 0.6</math></b></u>	$77.5 \pm 0.4$	86.1
	$k = 2$	$97.7 \pm 0.2$	$88.5 \pm 1.0$	$79.9 \pm 1.1$	<u><b><math>79.2 \pm 0.7</math></b></u>	86.4
	$k = 4$	<u><b><math>98.0 \pm 0.2</math></b></u>	<b><math>89.9 \pm 0.2</math></b>	$80.0 \pm 0.6$	<b><math>80.1 \pm 0.9</math></b>	<b>87.0</b>
	$k = 7$	$97.6 \pm 0.3$	$88.2 \pm 0.8$	<b><math>81.1 \pm 1.3</math></b>	$79.0 \pm 1.5$	<u><b>86.5</b></u>

# Analyzing Training Time Cost

Model's time efficiency is defined as the required time (seconds) for each update

Table 12: Training time of the PGrad-B evaluated on both PACS and OfficeHome dataset.

Method	PACS				OFFICEHOME			
	P	A	C	S	A	C	P	R
DANN	1.18	1.23	1.10	1.17	1.39	1.47	1.25	1.26
MLDG	1.03	1.02	1.00	1.01	0.49	0.46	0.43	0.41
Fish	1.63	1.66	1.61	1.58	1.47	1.57	1.45	1.58
PGrad-B	1.62	1.67	1.63	1.70	1.70	1.73	1.75	1.71

Thanks for attending our presentation!