

香港中文大學  
The Chinese University of Hong Kong

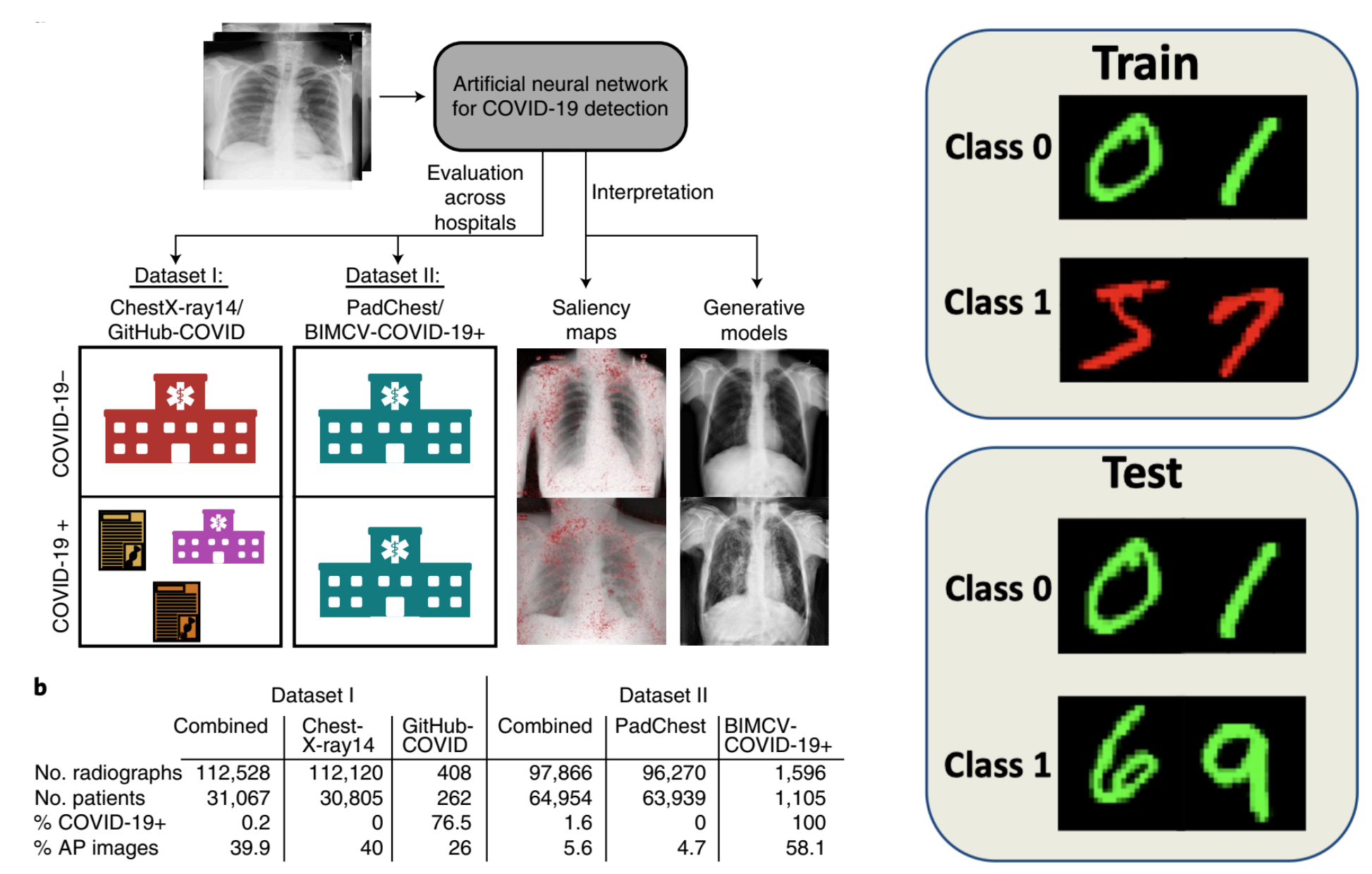
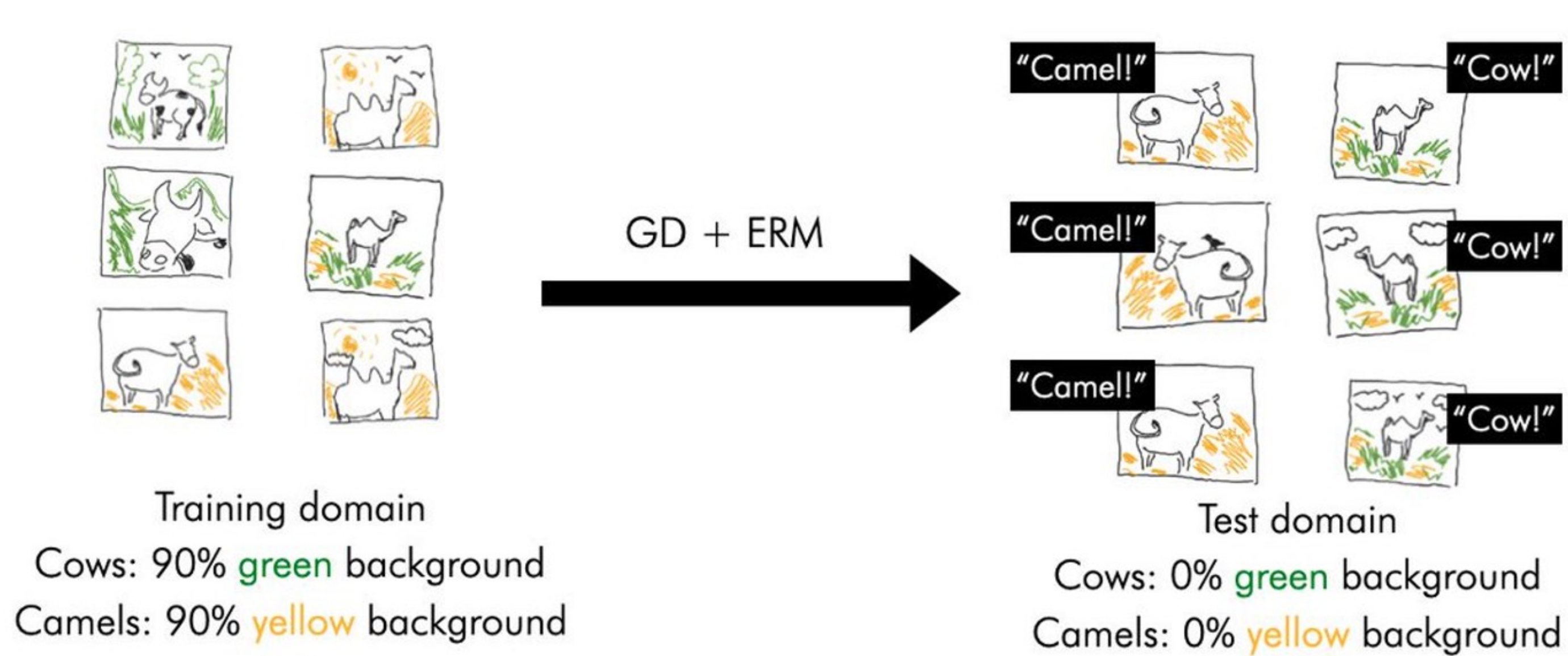


# Pareto Invariant Risk Minimization: Towards Mitigating the Optimization Dilemma in OOD Generalization

Yongqiang Chen  
CUHK & Tencent AI Lab

*with Kaiwen Zhou, Yatao Bian, Binghui Xie,  
Bingzhe Wu, Peilin Zhao, Bo Han, James Cheng and others.*

# Out-of-Distribution generalization



( Beery et al., 2018; Arjovsky et al., 2019; DeGrave et al. 2021; Ahuja et al., 2021; Zhang et al., 2022)

Models trained with Empirical Risk Minimization (ERM) are often:

- prone to **spurious correlations**
- can hardly generalize to **OOD** data

# Previous works focus on OOD objectives

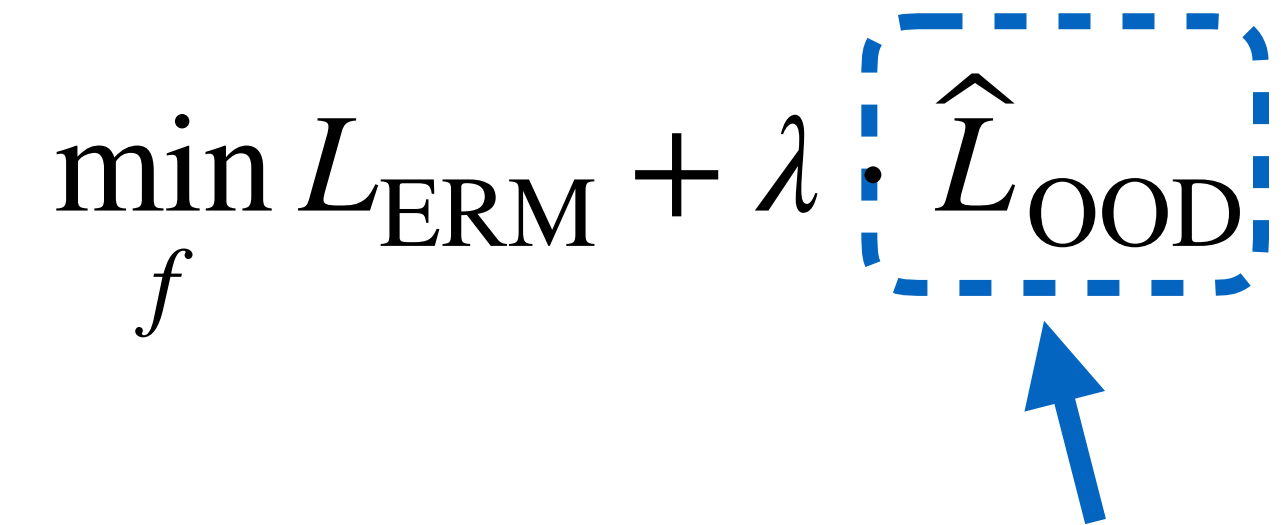
Previous works mostly focus on developing better **optimization objectives**:

$$\min_f L_{\text{ERM}} + \lambda \hat{L}_{\text{OOD}}$$

Regularization via some OOD objective

# The Optimization Dilemma in OOD Generalization

Previous works mostly focus on developing better *optimization objectives*:

$$\min_f L_{\text{ERM}} + \lambda \cdot \hat{L}_{\text{OOD}}$$


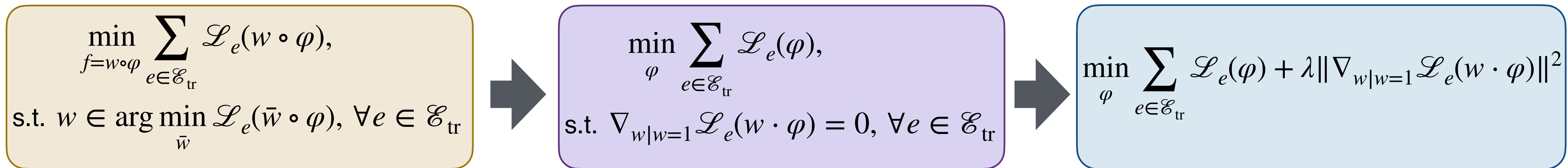
Regularization via some *relaxed* OOD objective

# The Optimization Dilemma in OOD Generalization

Previous works mostly focus on developing better **optimization objectives**:

$$\min_f L_{\text{ERM}} + \lambda \hat{L}_{\text{OOD}}$$

Regularization via some **relaxed** OOD objective



IRM  
😊

Linearized IRM with  $w \in \mathbb{R}^d$

IRM<sub>ℒ</sub>  
😓

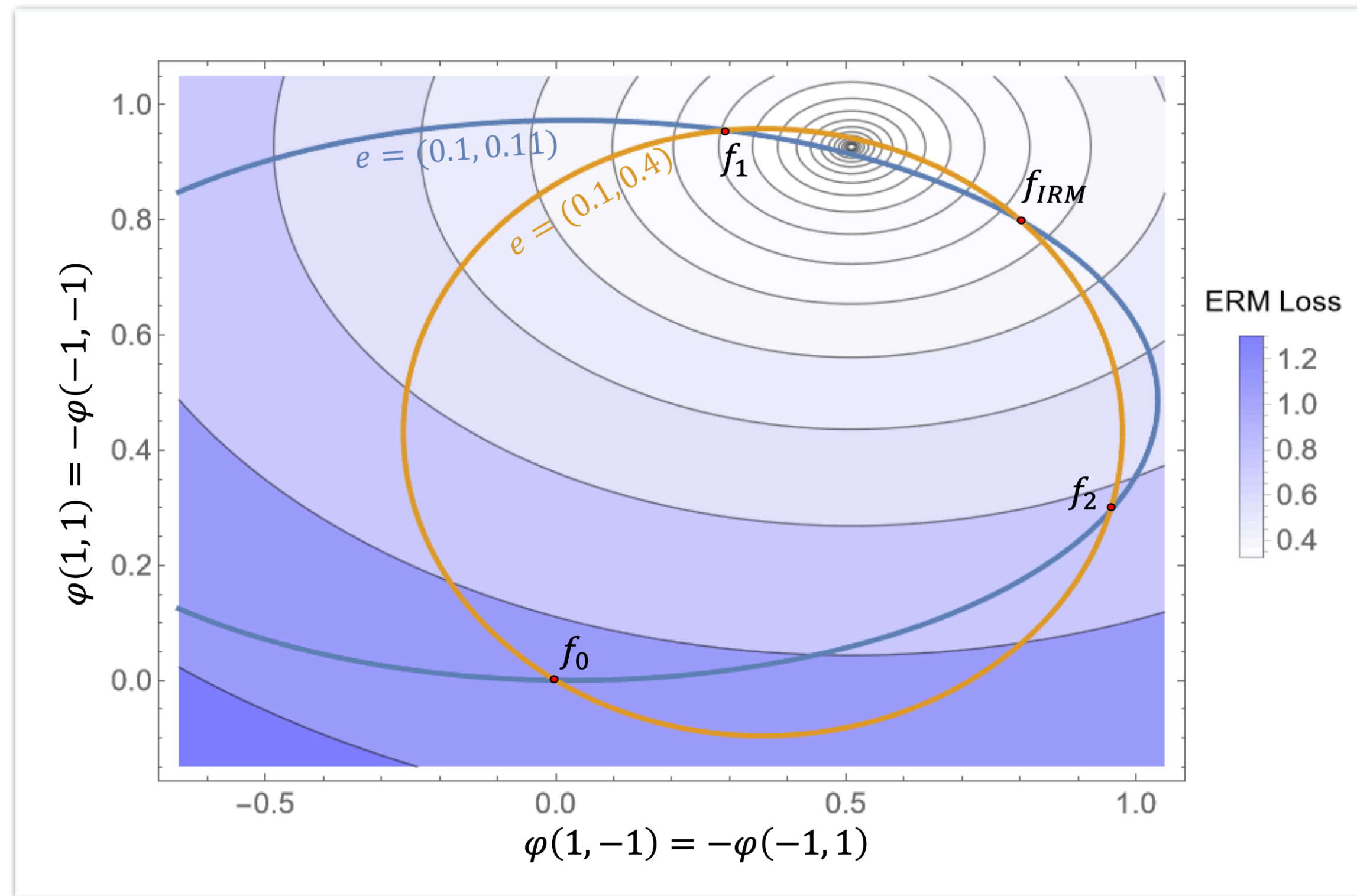
Soften the constraints

IRMv1  
😓 😓

# The Optimization Dilemma in OOD Generalization



The practical variants of IRM can have very different behaviors from the original IRM.



The ellipsoids are the solutions satisfying the **invariant constraints** in  $\text{IRM}_{\mathcal{S}}$

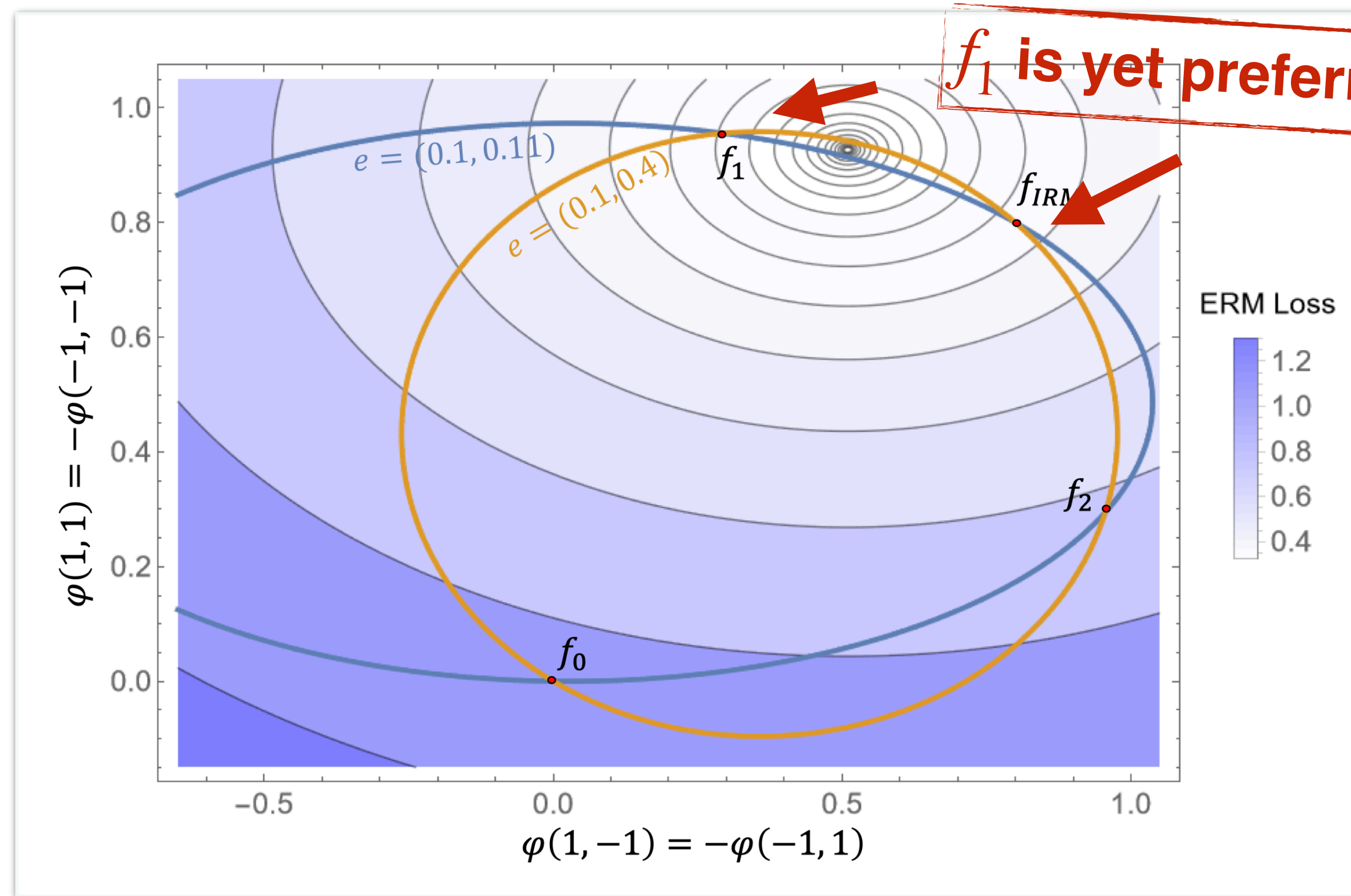
$$\nabla_{w|w=1} \mathcal{L}_e(w \cdot \varphi) = 0, \forall e \in \mathcal{E}_{\text{tr}}$$

Illustration of IRMv1 failures

# The Optimization Dilemma in OOD Generalization



The practical variants of IRM can have very different behaviors from the original IRM.



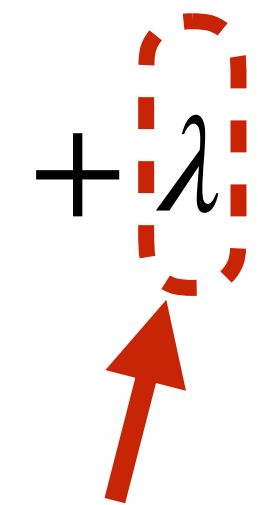
The ellipsoids are the solutions satisfying the **invariant constraints** in  $IRM_{\mathcal{S}}$

$$\nabla_{w|w=1} \mathcal{L}_e(w \cdot \varphi) = 0, \forall e \in \mathcal{E}_{\text{tr}}$$

Illustration of IRMv1 failures

# The Optimization Dilemma in OOD Generalization

Previous works mostly focus on developing better *optimization objectives*:

$$\min_f L_{\text{ERM}} + \lambda \cdot \hat{L}_{\text{OOD}}$$


$\lambda$  is *hard to tune*



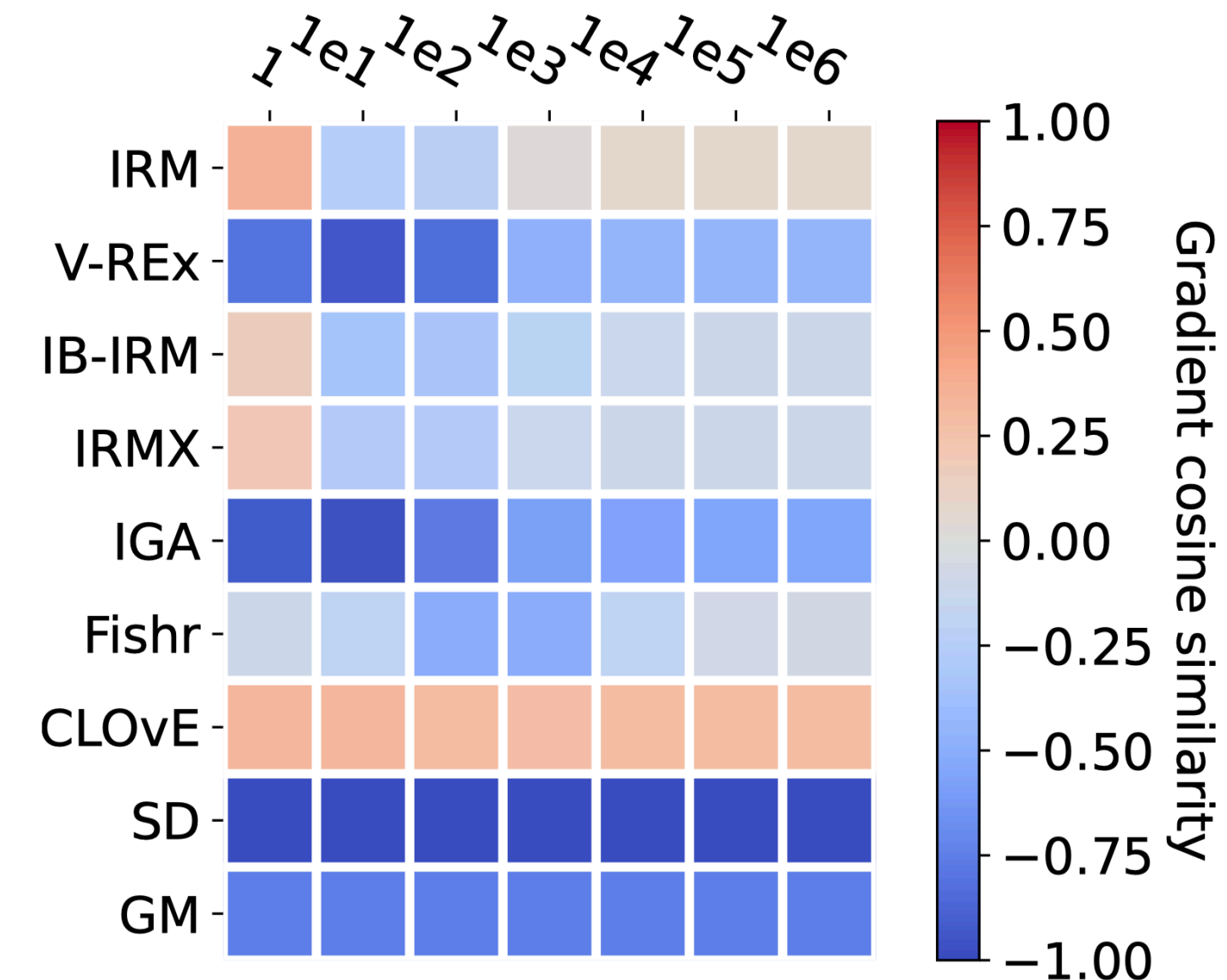
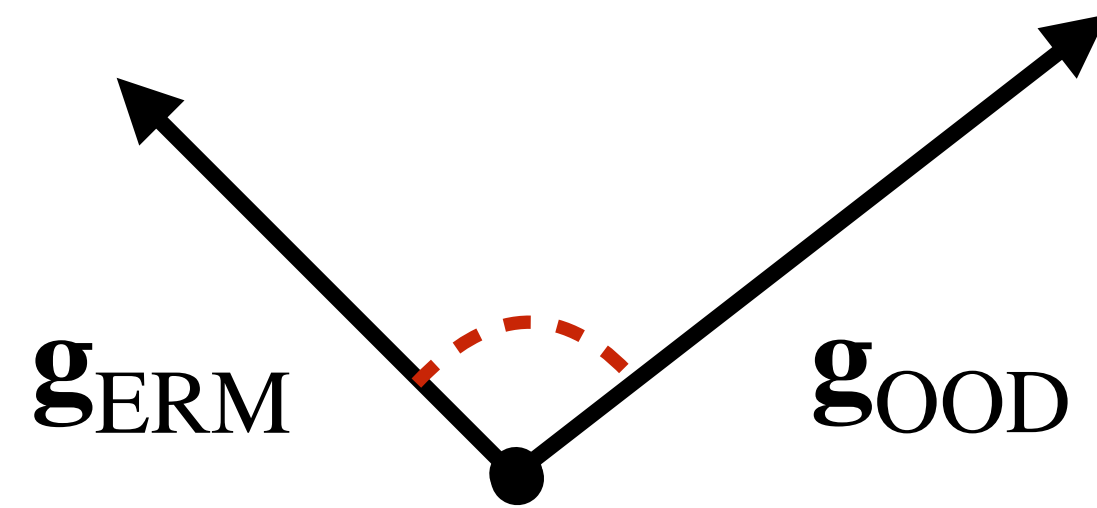
# The Optimization Dilemma in OOD Generalization

Previous works mostly focus on developing better *optimization objectives*:

$$\min_f L_{\text{ERM}} + \lambda \cdot \hat{L}_{\text{OOD}}$$

$\lambda$  is *hard to tune*

**Gradient Conflicts** generically exist between ERM and OOD objectives:

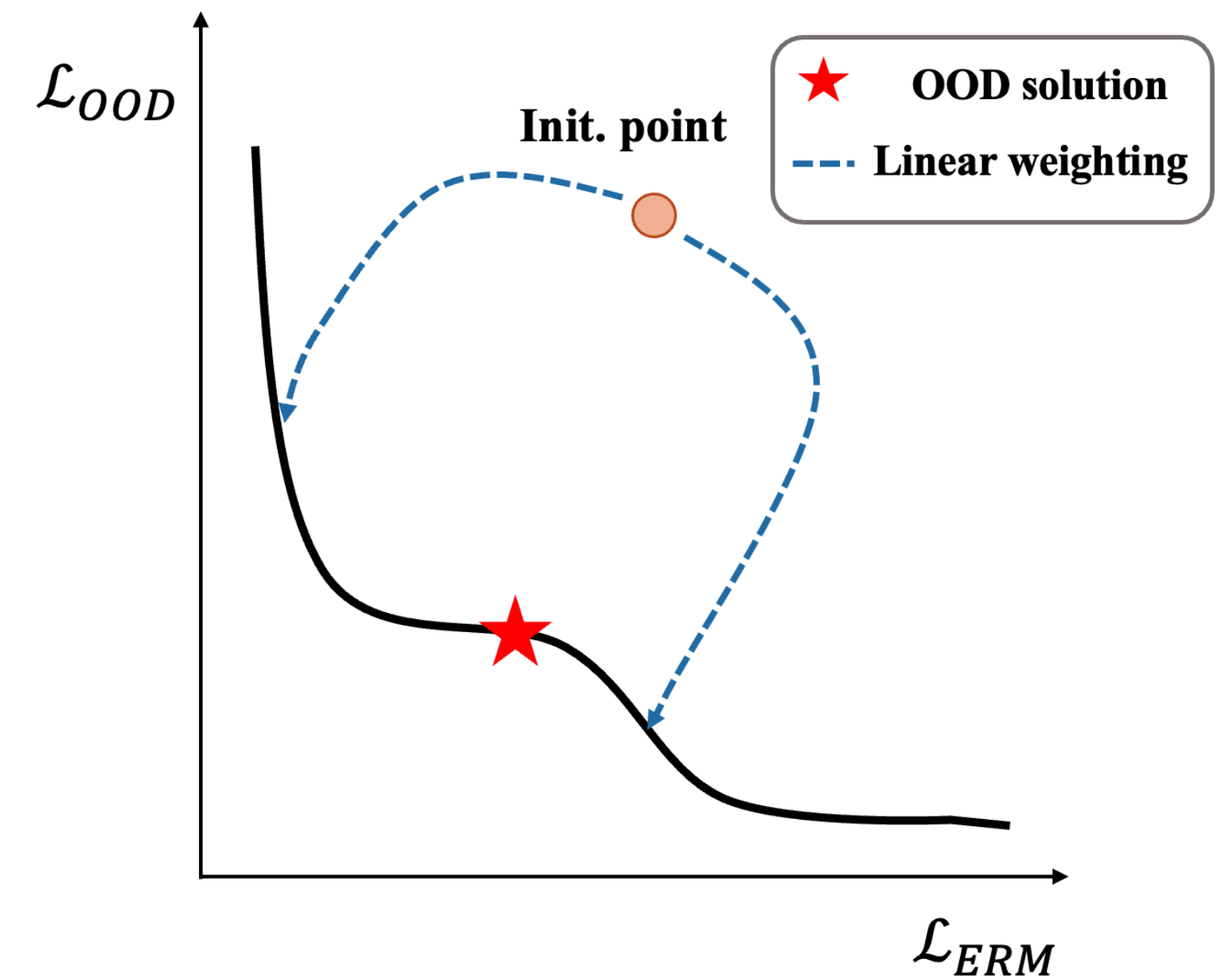


# The Optimization Dilemma in OOD Generalization

The typically used linear weighting scheme cannot reach ***non-convex part of pareto front solutions***

$$\min_f L_{\text{ERM}} + \lambda \cdot \hat{L}_{\text{OOD}}$$

The linear weight scheme

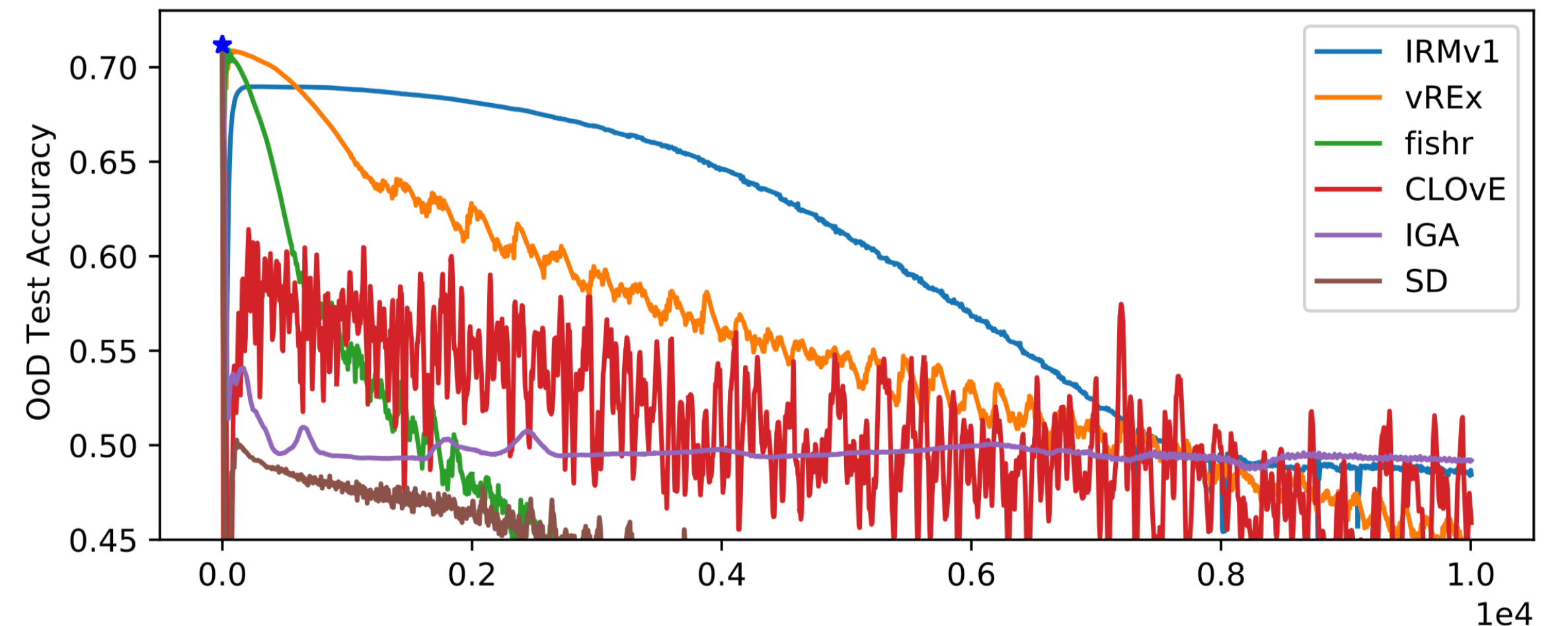
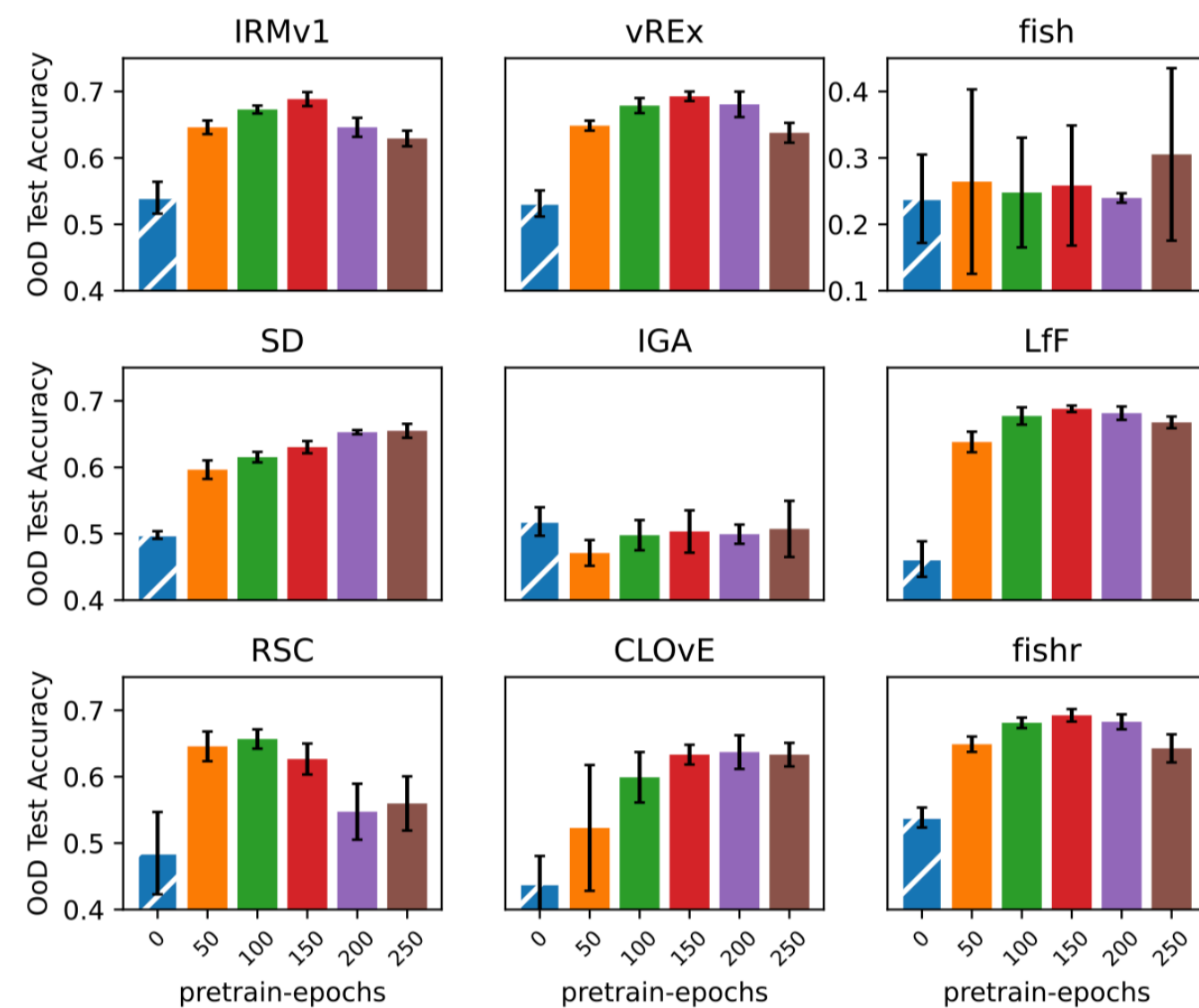
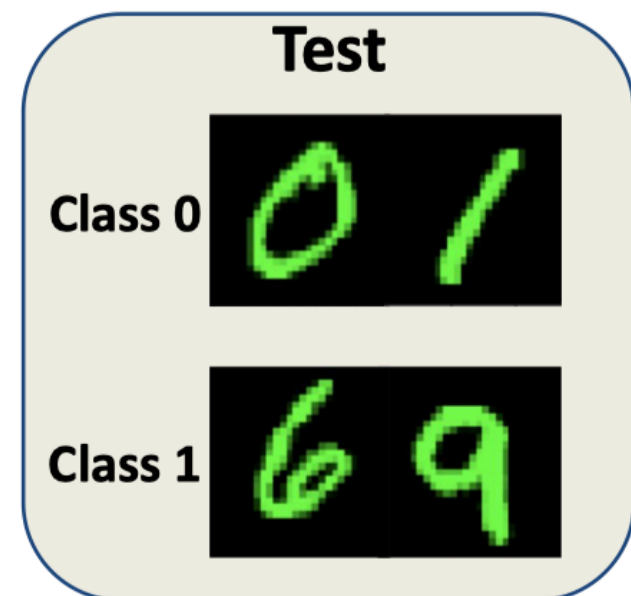
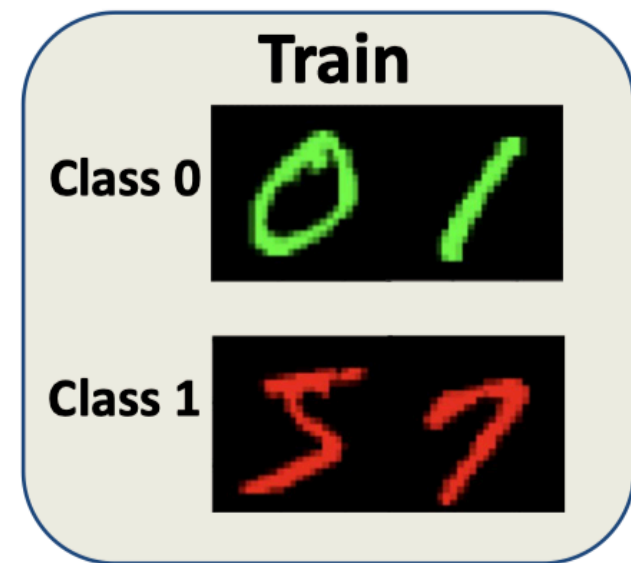


# The Optimization Dilemma in OOD Generalization

Even the desired solution is reachable, the scheme requires **exhaustive hyperparameter tuning**:

$$\min_f L_{\text{ERM}} + \lambda \cdot \hat{L}_{\text{OOD}}$$

$\lambda$  is **too strong** to learn the correlation;  $\lambda$  is **too weak** to keep the invariance



# The Optimization Dilemma in OOD Generalization

The usual optimization formula of OOD objectives in practice:

$$\min_f L_{\text{ERM}} + \lambda \cdot \hat{L}_{\text{OOD}}$$

$\lambda$  is **hard to tune** Regularization via some **relaxed** OOD objective

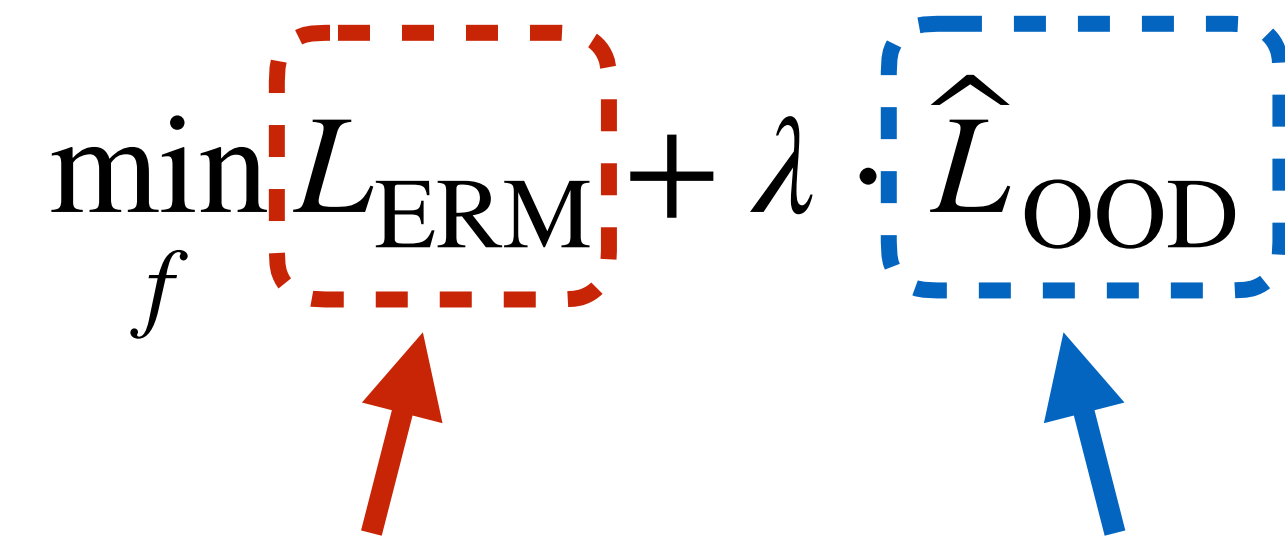
- $\hat{L}_{\text{OOD}}$  usually has **a large gap** from the original one;
- $\lambda$  is **hard to tune**, i.e.,
  - ▶ Not all potentially optimal solutions are reachable;
  - ▶ Even reachable, it still requires exhaustive tuning efforts to find a proper  $\lambda$ ;

*As the traditional optimization scheme fails*

***How to obtain a desired OOD solution  
under the ERM and OOD conflicts?***

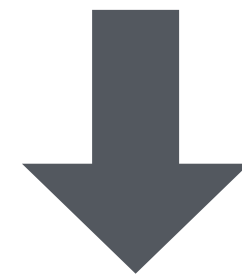
# From a Multi-Objective Optimization perspective...

The optimization of IRM essentially handles the *trade-off* between

$$\min_f L_{\text{ERM}} + \lambda \cdot \hat{L}_{\text{OOD}}$$


Capturing the statistical correlations

Enforcing the invariance of learned correlations



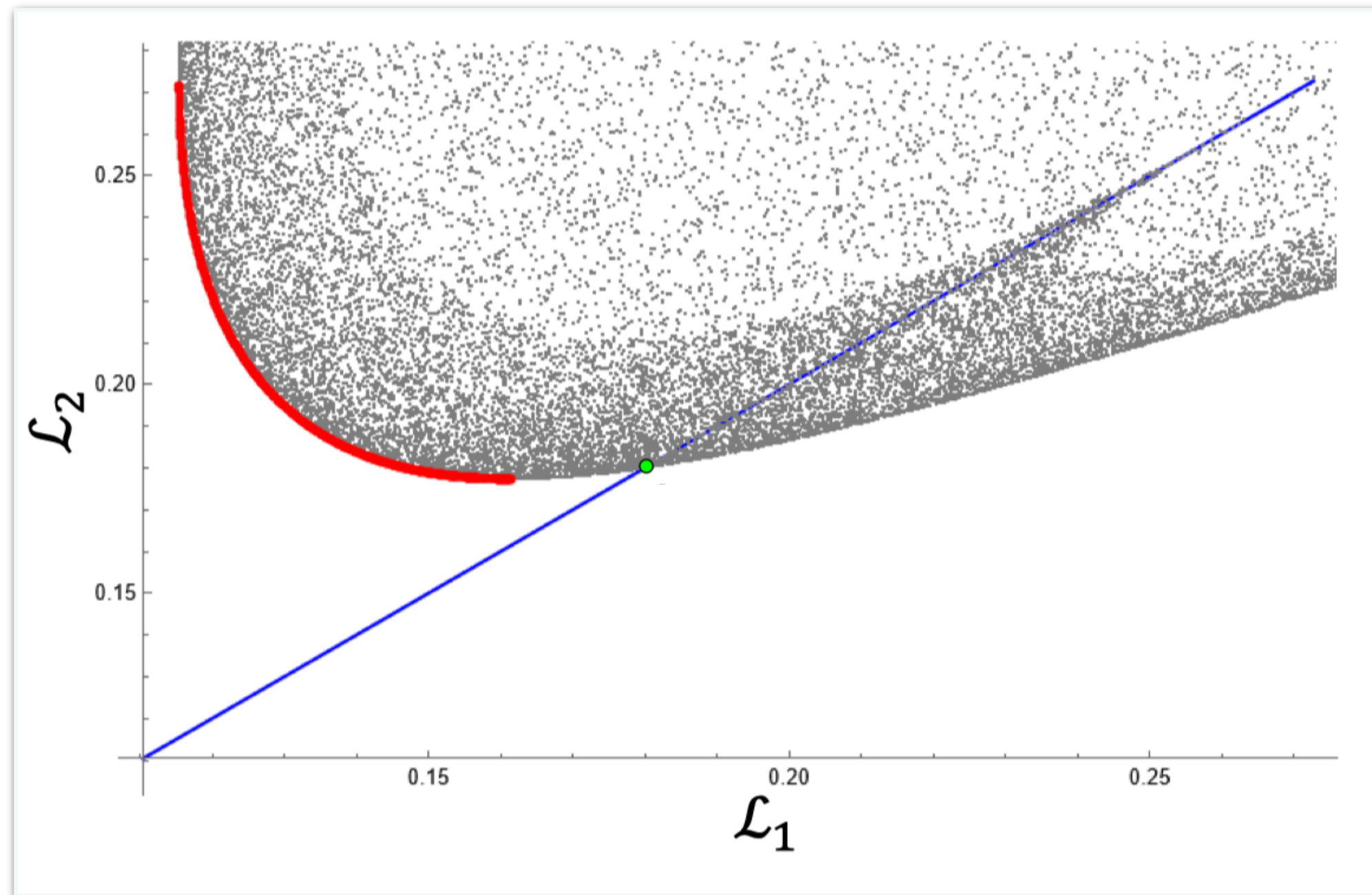
**Oh, it's a Multi-Objective Optimization (MOO)!**

$$\min_f \{L_{\text{ERM}}, \hat{L}_{\text{OOD}}\}^T$$

# From a Multi-Objective Optimization perspective...

Assume we have the Multi-Objective Optimization (MOO) problem with 2 objectives:

$$\min_{f=w \cdot \varphi} \{L_1, L_2\}^T$$



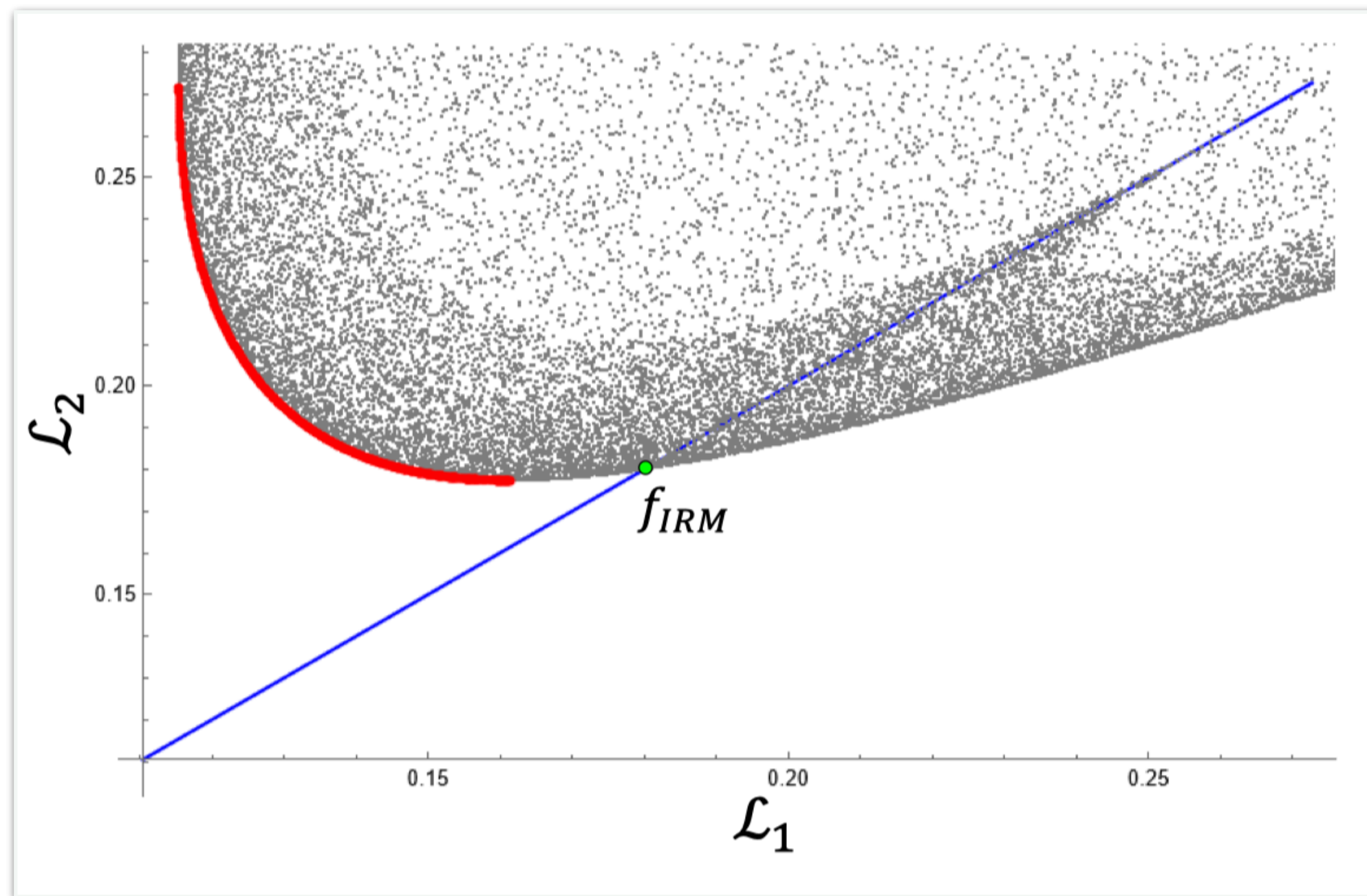
Simulated Pareto front

- A solution  $f$  (with  $\{L_1, L_2\}^T$ ) **dominates**  $\bar{f}$  (with  $\{\bar{L}_1, \bar{L}_2\}^T$ ) if both  $L_1 \leq \bar{L}_1$  and  $L_2 \leq \bar{L}_2$ ;
- **Pareto optimal solutions** are the set of solutions dominated by none;
- Their images form the **Pareto front**;

# From a Multi-Objective Optimization perspective...

Assume we have 2 training environments, a natural MOO formulation of IRMv1 is:

$$\min_{f=w \cdot \varphi} \{L_1, L_2, L_{\text{IRM}}\}^T$$



Simulated Pareto front

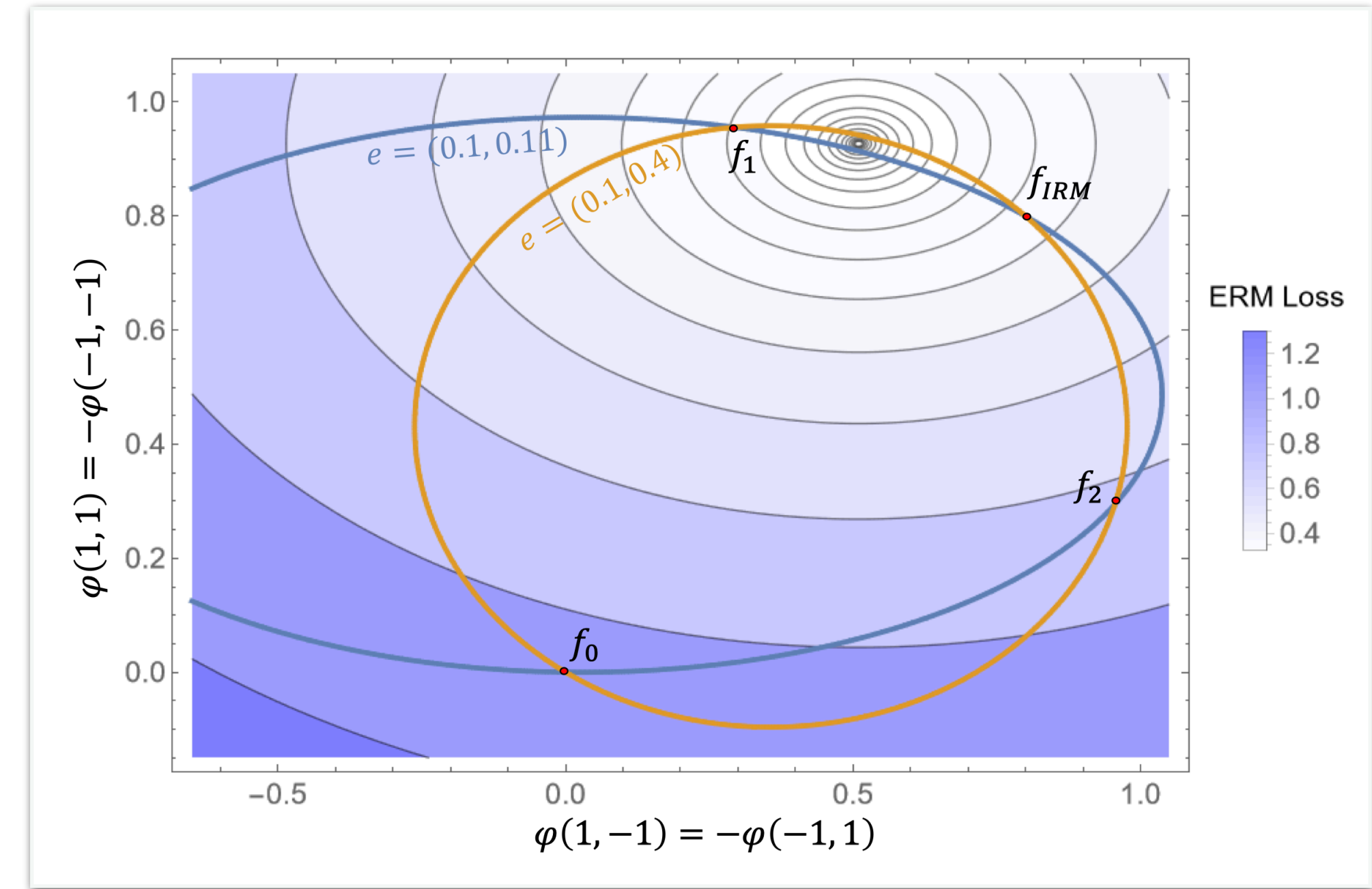


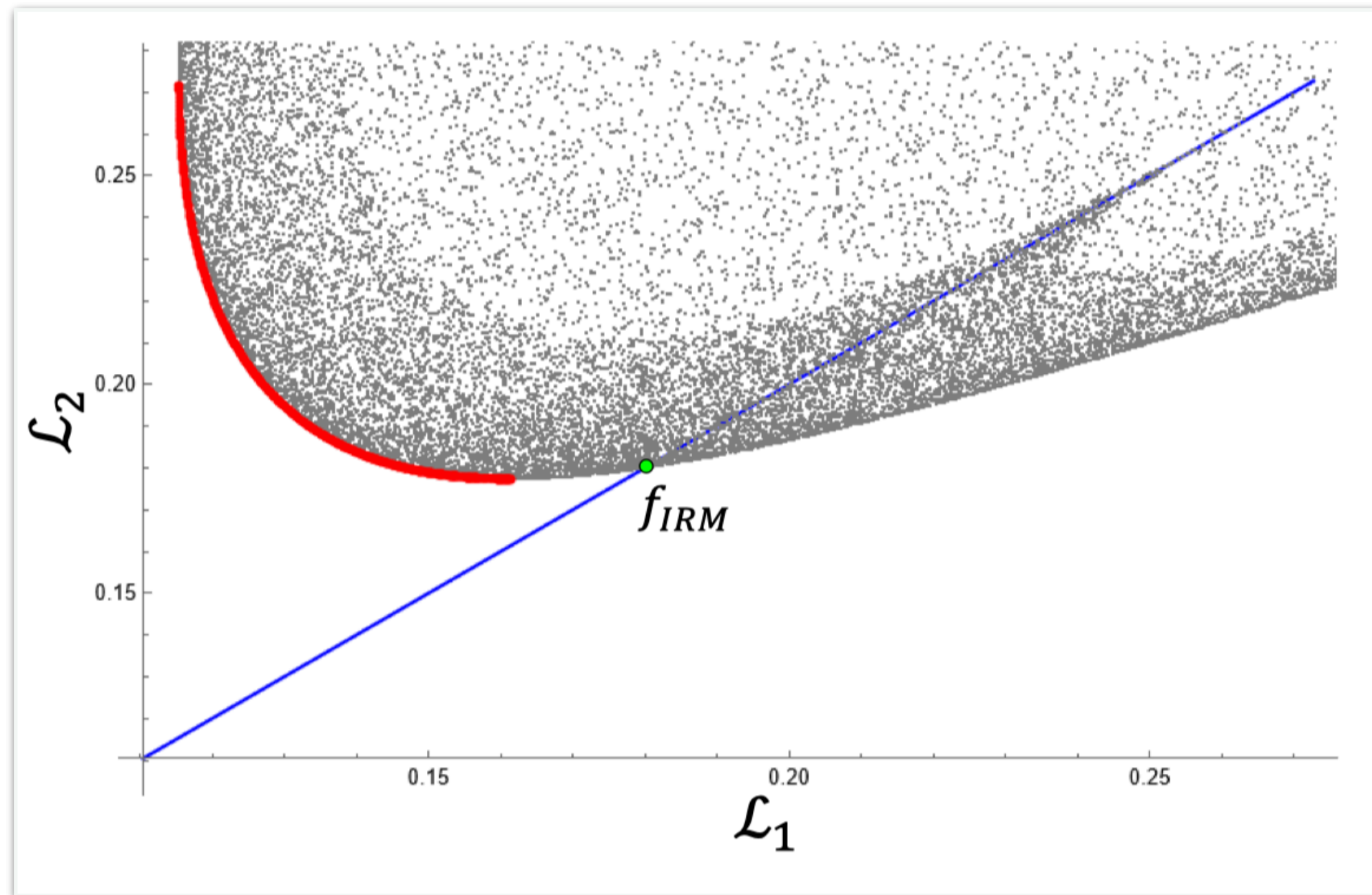
Illustration of IRMv1 failures



# From a Multi-Objective Optimization perspective...

The failures of practical IRM variants is because of using **bad objectives!**

$$\min_{f=w \cdot \varphi} \{L_1, L_2, L_{\text{IRM}}\}^T$$



Simulated Pareto front

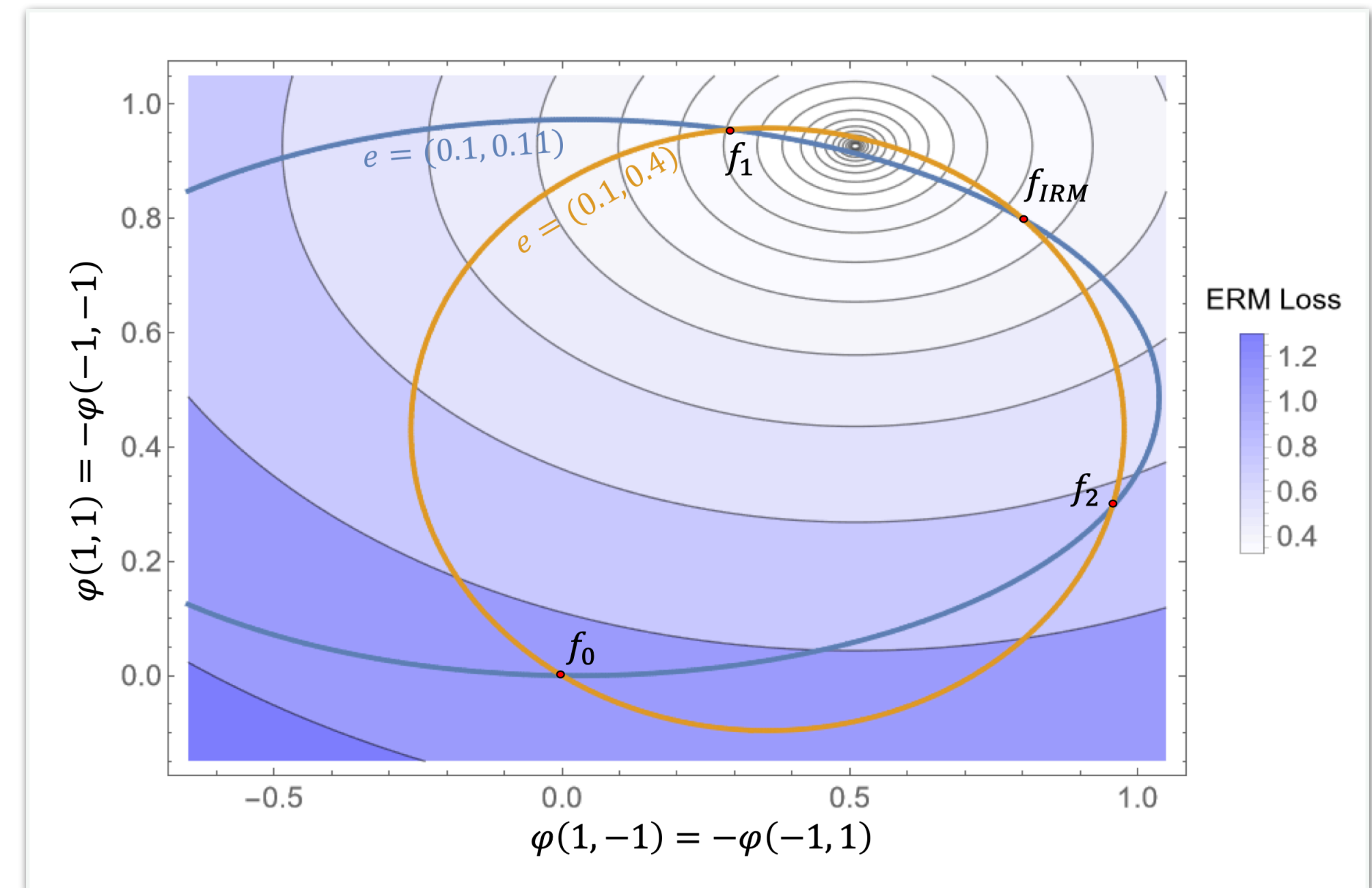


Illustration of IRMv1 failures

# Robustify MOO objectives

IRM can extrapolate **stationary points** of **negative** combinations of training environments:

$$\left\{ \sum_{e \in \mathcal{E}_{\text{tr}}} \lambda_e \mathcal{D}_e \mid \sum_{e \in \mathcal{E}_{\text{tr}}} \lambda_e = 1, \lambda_e \geq 0, \forall e \right\} \rightarrow \left\{ \sum_{e \in \mathcal{E}_{\text{tr}}} \lambda_e \mathcal{D}_e \mid \sum_{e \in \mathcal{E}_{\text{tr}}} \lambda_e = 1, \lambda_e \leq 0, \forall e \right\}$$

Invariance buys extrapolation powers

Queries with decreasing popularity  
e.g. "ICLR schedule"

Queries with decreasing popularity  
e.g. "Easter bunny"

Queries with constant popularity  
e.g. "Orange juice"

An invariant regression on the training environments is optimal far beyond their convex hull.

# Robustify MOO objectives

We can introduce **additional** guidance that **directly** enforces extrapolation at certain region.

$$\left\{ \sum_{e \in \mathcal{E}_{\text{tr}}} \lambda_e \mathcal{D}_e \mid \sum_{e \in \mathcal{E}_{\text{tr}}} \lambda_e = 1, \lambda_e \geq 0, \forall e \right\} \rightarrow \left\{ \sum_{e \in \mathcal{E}_{\text{tr}}} \lambda_e \mathcal{D}_e \mid \sum_{e \in \mathcal{E}_{\text{tr}}} \lambda_e = 1, \lambda_e \leq 0, \forall e \right\} \rightarrow \left\{ \sum_{e \in \mathcal{E}_{\text{tr}}} \lambda_e \mathcal{D}_e \mid \sum_{e \in \mathcal{E}_{\text{tr}}} \lambda_e = 1, \lambda_e \leq -\beta, \forall e \right\}$$

Invariance buys extrapolation powers

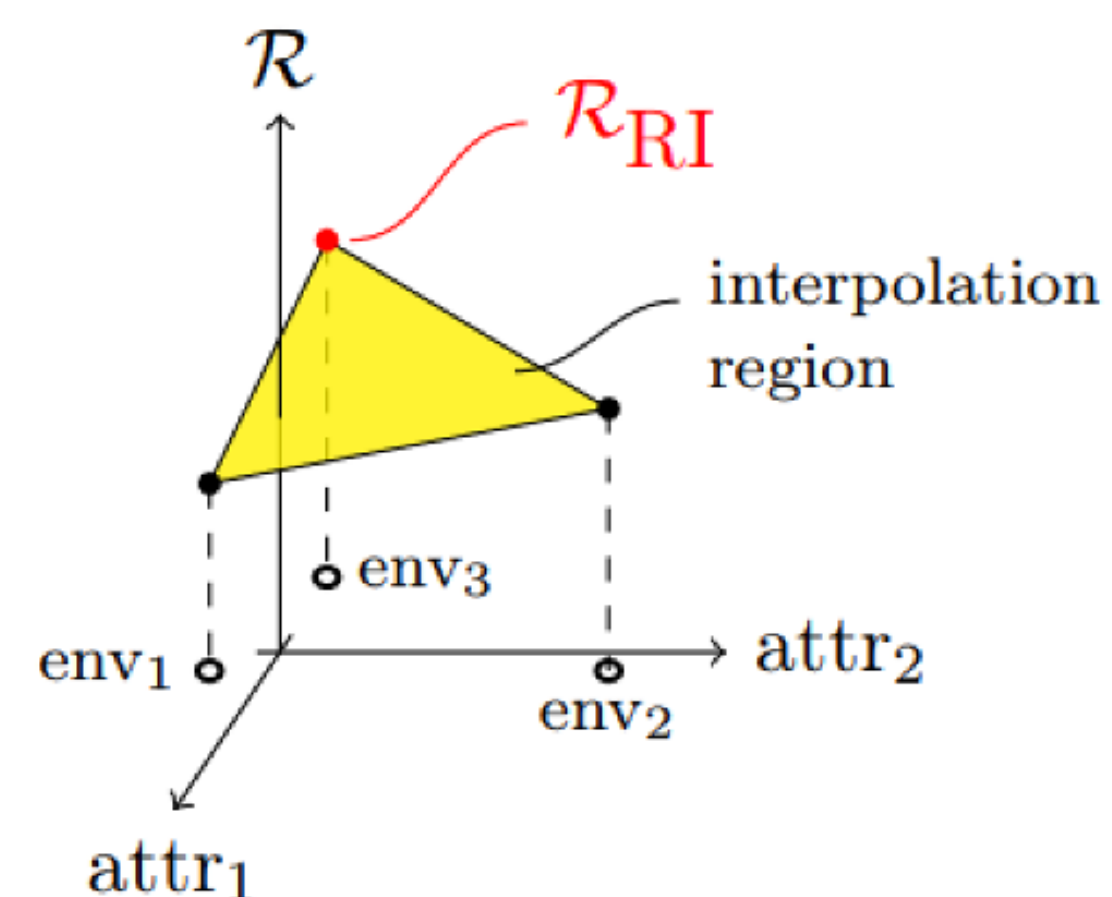
Queries with decreasing popularity  
e.g. "ICLR schedule"

Queries with decreasing popularity  
e.g. "Easter bunny"

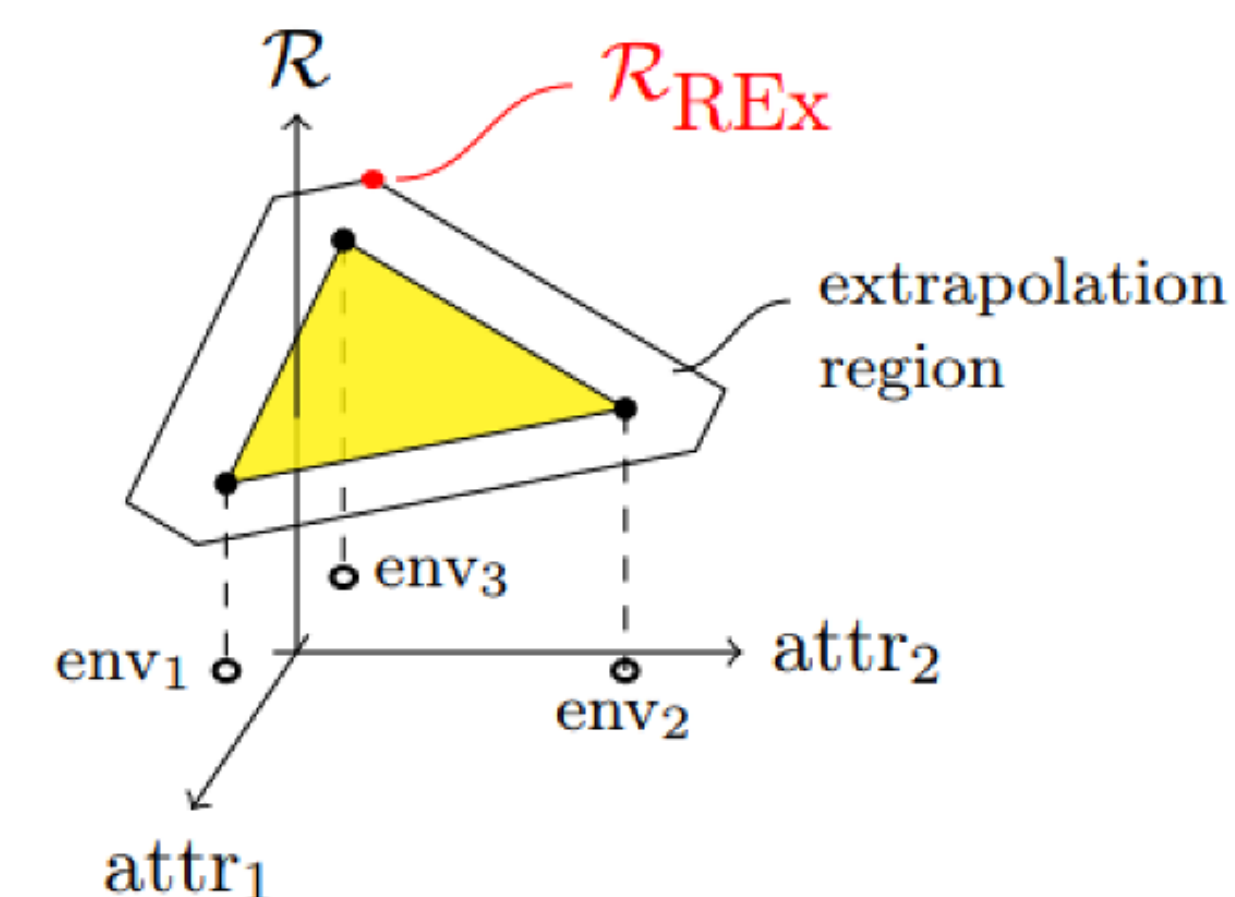
Queries with constant popularity  
e.g. "Orange juice"

An invariant regression on the training environments is optimal far beyond their convex hull.

Risk Interpolation



Risk Extrapolation

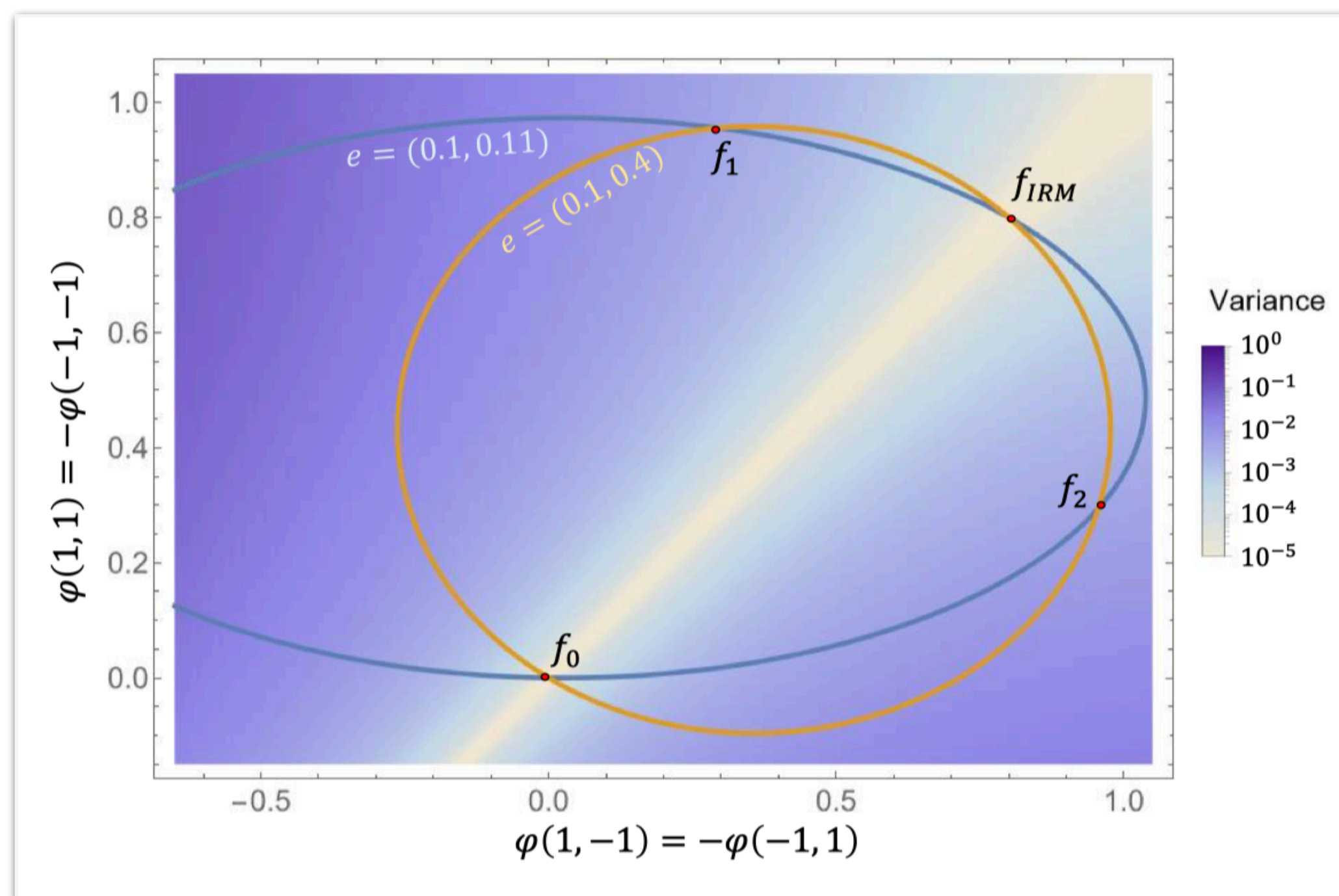


This brings us a new MOO objectives, IRMX:  $\min_{f=w \cdot \varphi} \{L_1, L_2, L_{\text{IRM}}, L_{\text{REx}}\}^T$

# PAIR: PAreto Invariant Risk minimization



A PAIRed journey into the adventure of extrapolation:  $\min_{f=w \cdot \varphi} \{L_{\text{ERM}}, L_{\text{IRM}}, L_{\text{REX}}\}^T$



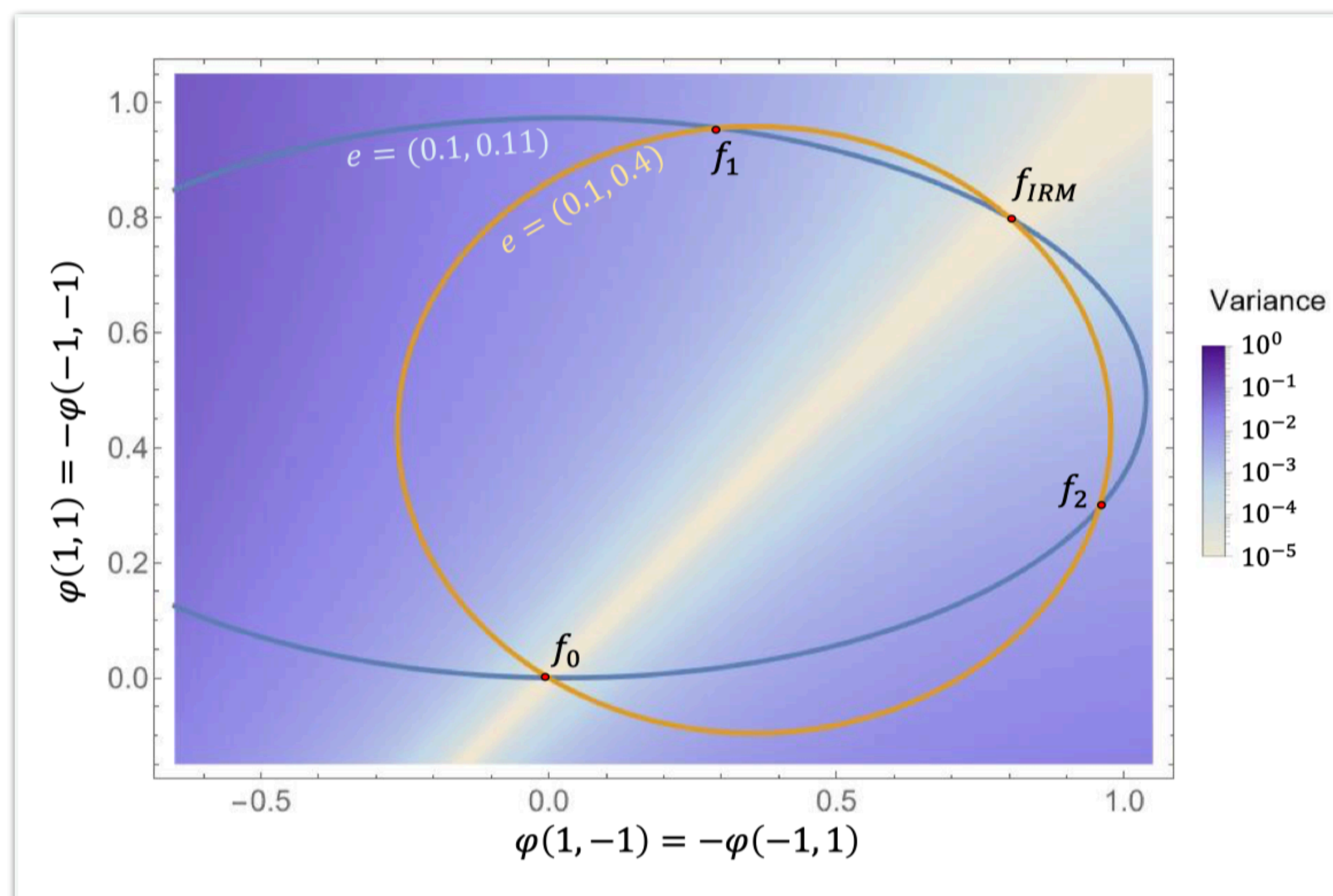
## Theoretical results (Informal):

IRMx solves the IRMv1 failures under any environment settings in (*Kamath et al., 2021*).

# PAIR: PAreto Invariant Risk minimization



A PAIRed journey into the adventure of extrapolation:  $\min_{f=w \cdot \varphi} \{L_{\text{ERM}}, L_{\text{IRM}}, L_{\text{REX}}\}^T$



**IRMX raises more challenges in hp. tuning!**

## Theoretical results (Informal):

IRMX solves the IRMv1 failures under any environment settings in (Kamath et al., 2021).

# PAIR: PAreto Invariant Risk minimization

---

IRMX raises more challenges in the optimization:

$$\min_{f=w \cdot \varphi} \{L_{\text{ERM}}, L_{\text{IRM}}, L_{\text{REX}}\}^T$$

- The Pareto front becomes **more complicated**:

# PAIR: PAreto Invariant Risk minimization

IRMX raises more challenges in the optimization:

$$\min_{f=w \cdot \varphi} \{L_{\text{ERM}}, L_{\text{IRM}}, L_{\text{REX}}\}^T$$

- The Pareto front becomes **more complicated**:
  - ✓ The optimizer needs to be able to reach **any** Pareto optimal solutions!

e.g., MGDA algorithms (*Désidéri, 2012*)

# PAIR: PAreto Invariant Risk minimization

IRMX raises more challenges in the optimization:

$$\min_{f=w \cdot \varphi} \{L_{\text{ERM}}, L_{\text{IRM}}, L_{\text{REX}}\}^T$$

- The Pareto front becomes **more complicated**:
  - ✓ The optimizer needs to be able to reach **any** Pareto optimal solutions!
- There can be **multiple** Pareto optimal solutions:



# PAIR: PAreto Invariant Risk minimization

IRMX raises more challenges in the optimization:

$$\min_{f=w \cdot \varphi} \{L_{\text{ERM}}, L_{\text{IRM}}, L_{\text{REX}}\}^T$$

- The Pareto front becomes **more complicated**:
  - ✓ The optimizer needs to be able to reach **any** Pareto optimal solutions!
- There can be **multiple** Pareto optimal solutions:
  - ✓ A **preference** of each objective is required!

## Exact Pareto Optimality:

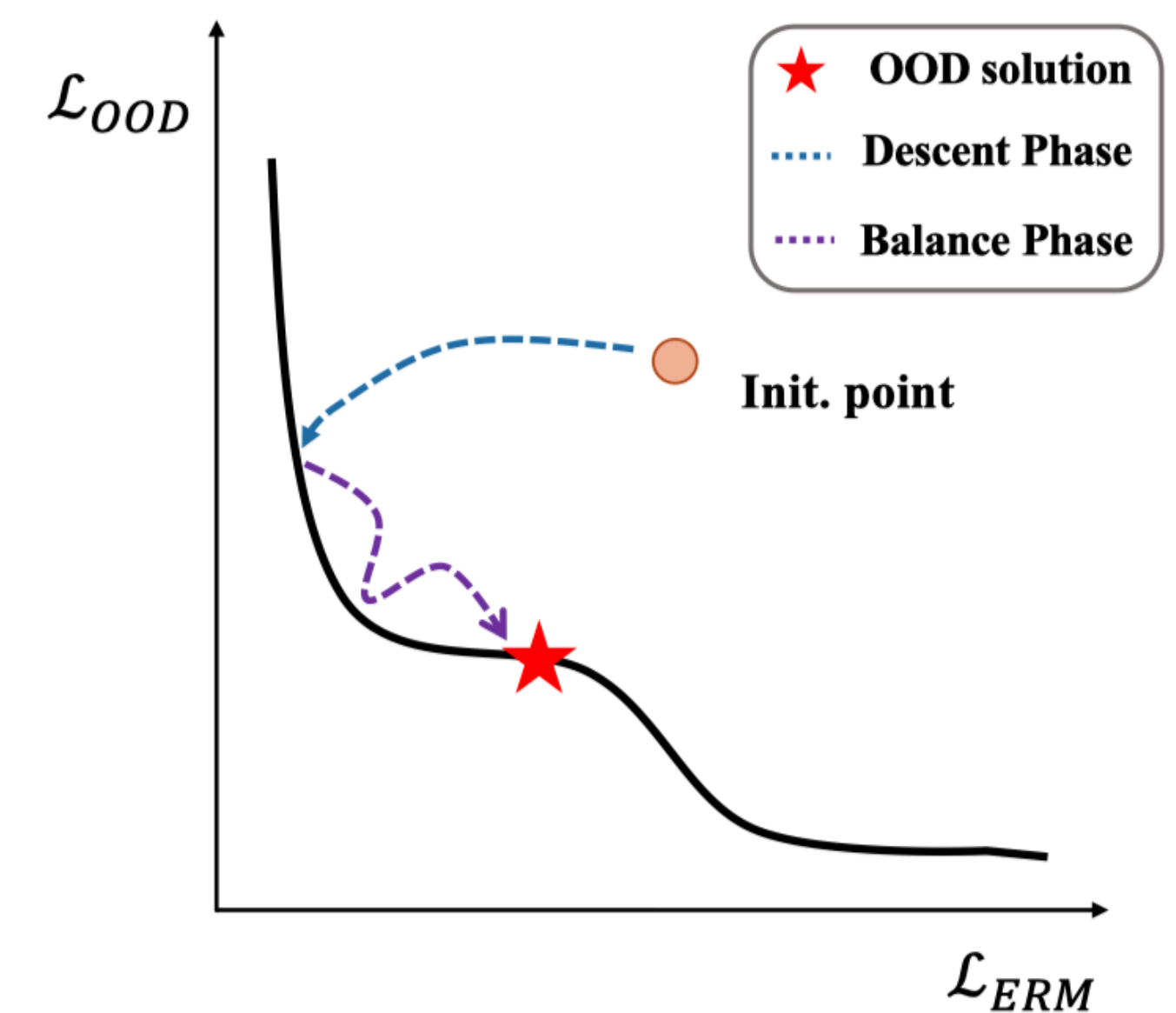
Given a preference  $\mathbf{p} = \{p_{\text{ERM}}, p_{\text{IRM}}, p_{\text{REX}}\}^T$  for each objective, a solution  $\hat{\mathbf{L}} = \{\hat{L}_{\text{ERM}}, \hat{L}_{\text{IRM}}, \hat{L}_{\text{REX}}\}^T$  satisfies Exact Pareto Optimality iff.  $p_{\text{ERM}} \hat{L}_{\text{ERM}} = p_{\text{IRM}} \hat{L}_{\text{IRM}} = p_{\text{REX}} \hat{L}_{\text{REX}}$ .

# PAIR: PAreto Invariant Risk minimization

IRMX raises more challenges in the optimization:

$$\min_{f=w \cdot \varphi} \{L_{\text{ERM}}, L_{\text{IRM}}, L_{\text{REX}}\}^T$$

- The Pareto front becomes **more complicated**:
  - ✓ The optimizer needs to be able to reach **any** Pareto optimal solutions!
- There can be **multiple** Pareto optimal solutions:
  - ✓ A **preference** of each objective is required! **PAIR-o** as the OOD optimizer;



*Exact Pareto optimal search*

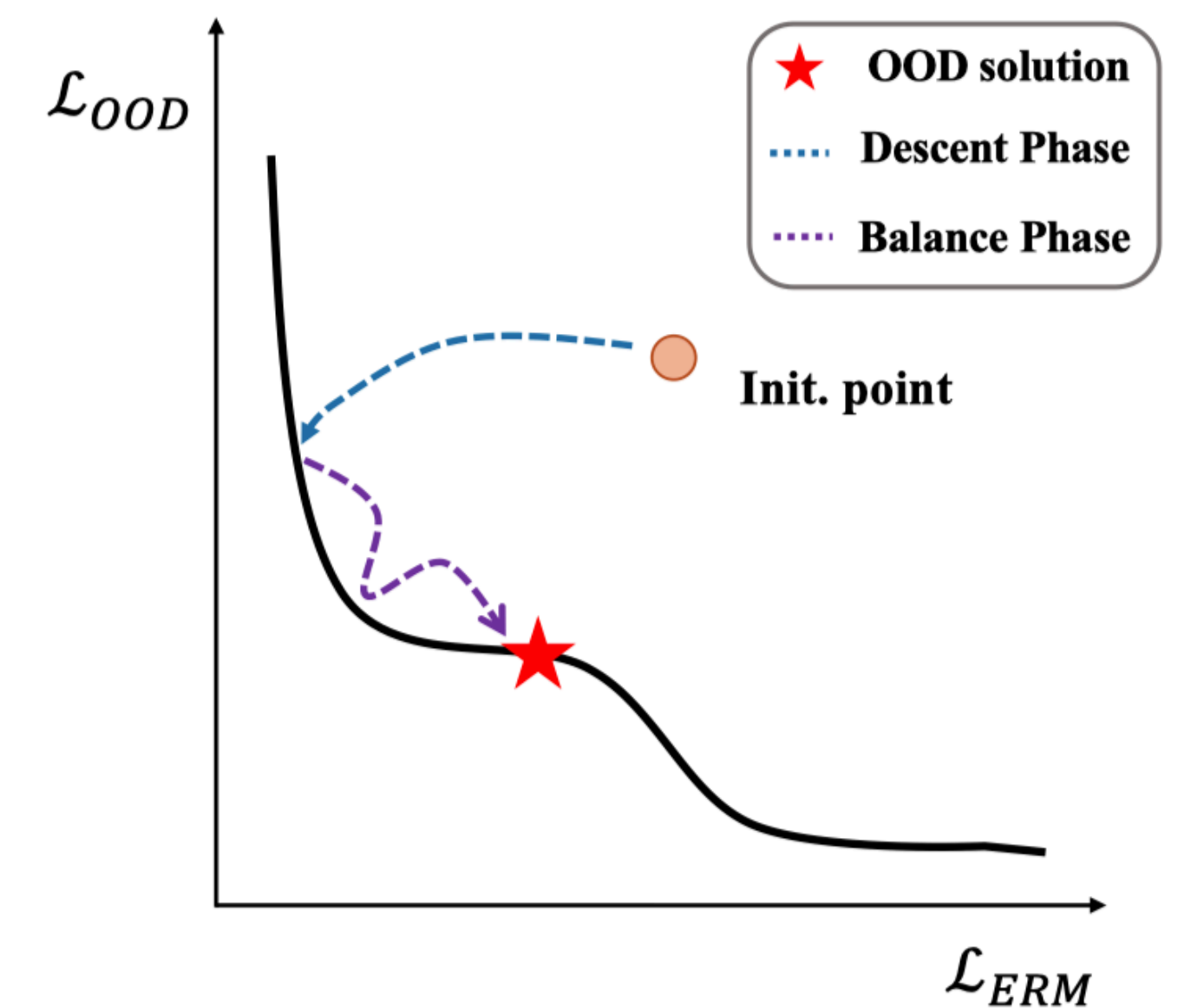
## Theoretical results (Informal):

Under mild assumptions, let  $f_{\text{OOD}}$  be the desired OOD solution w.r.t. an underlying preference  $\mathbf{p}_{\text{OOD}}$ , PAIR-o converges and approximates to  $f_{\text{OOD}}$  for any approximated  $\hat{\mathbf{p}}_{\text{OOD}}$ .

# PAIR: PAreto Invariant Risk minimization

IRMX raises more challenges in the optimization:

$$\min_{f=w \cdot \varphi} \{L_{\text{ERM}}, L_{\text{IRM}}, L_{\text{REX}}\}^T$$



*Exact Pareto optimal search*

- The Pareto front becomes **more complicated**:
  - ✓ The optimizer needs to be able to reach **any** Pareto optimal solutions!
- There can be **multiple** Pareto optimal solutions:
  - ✓ A **preference** of each objective is required! **PAIR-o** as the OOD optimizer;
  - ✓ It also motivates a new model selection criteria, by selecting models that maximally satisfy the Exact Pareto Optimality! **PAIR-s** as the OOD model selector;

# Causal Invariance Recovery Tests

## Regression target:

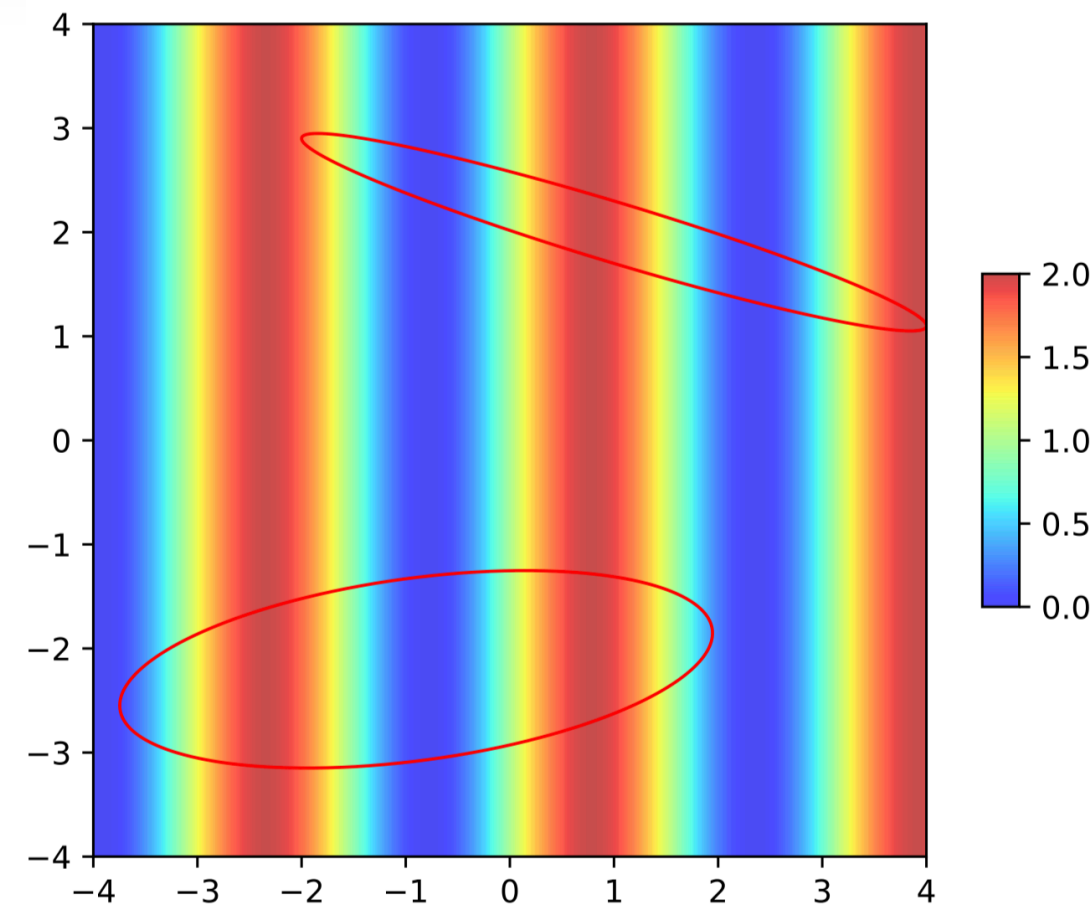
$Y = \sin(X_1) + 1$ , only depends on the x-axis;

## Training envs:

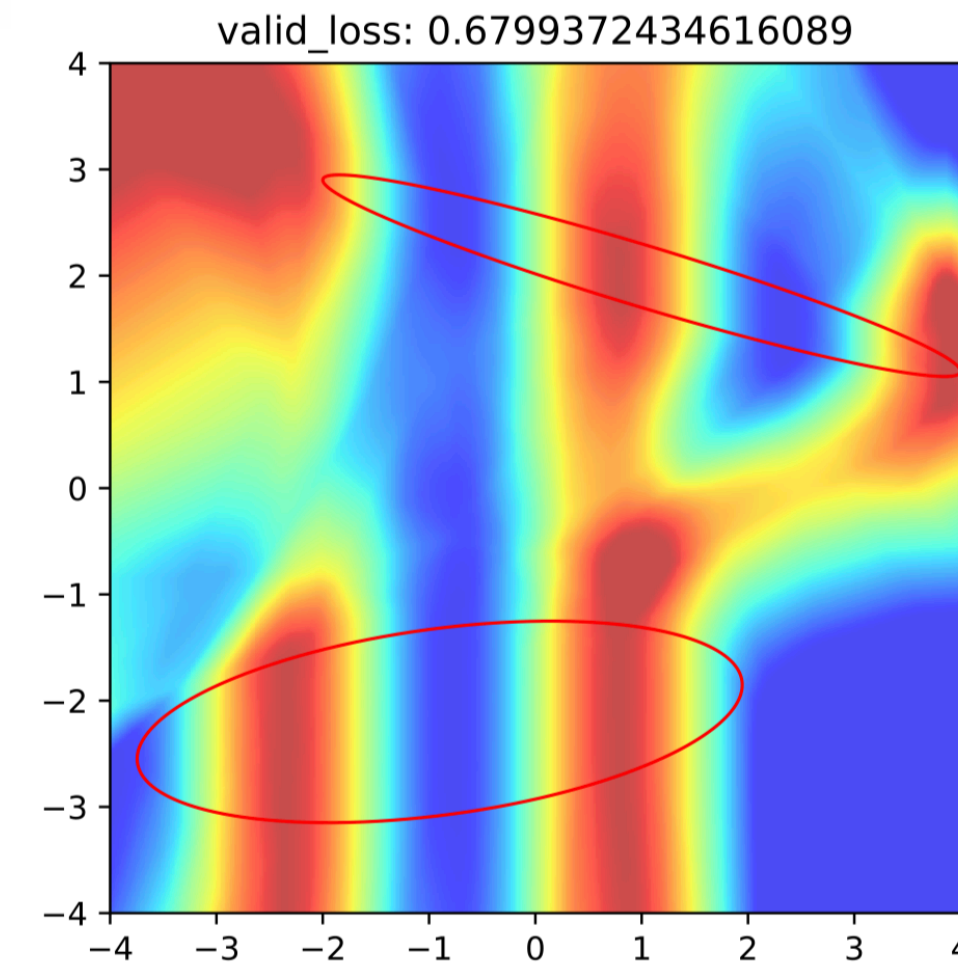
Two elliptical regions (Gaussian distributions) marked in red;

## Invariance:

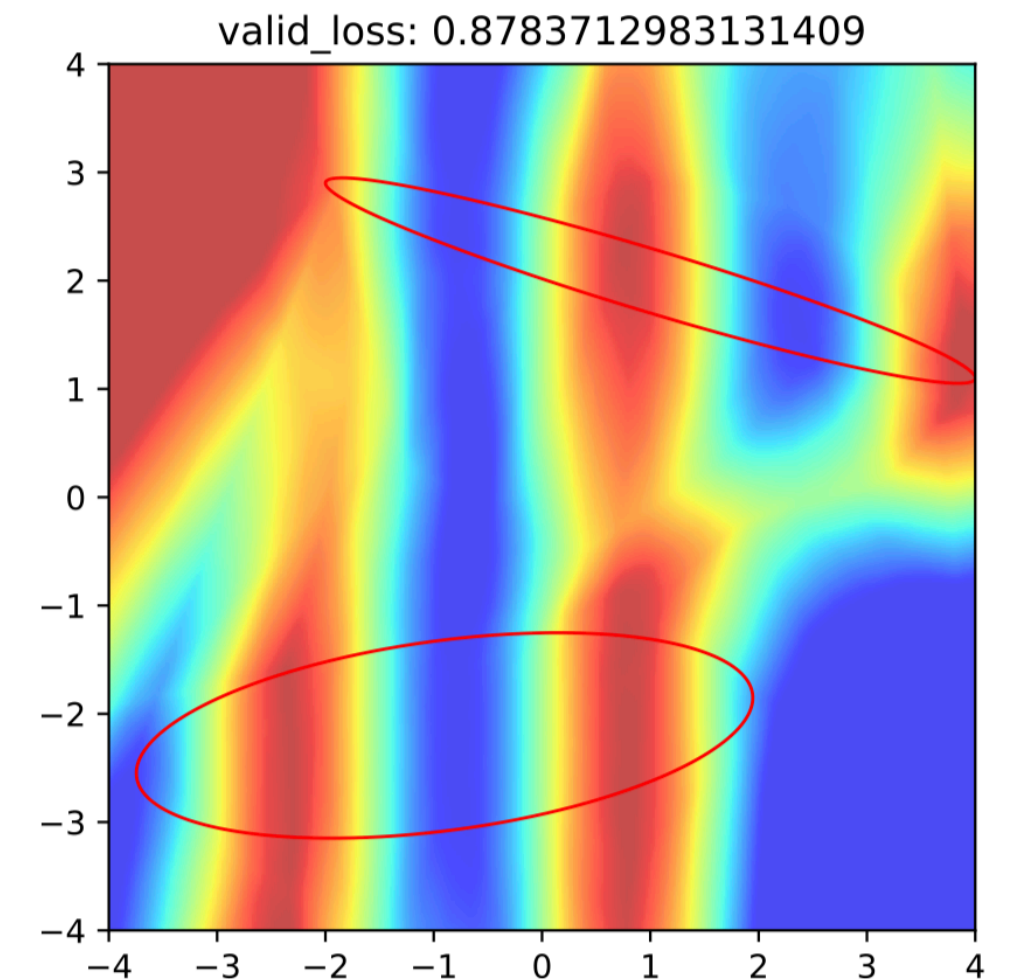
The **overlapped** x-axis region, i.e.,  $[-2, 2]$ .



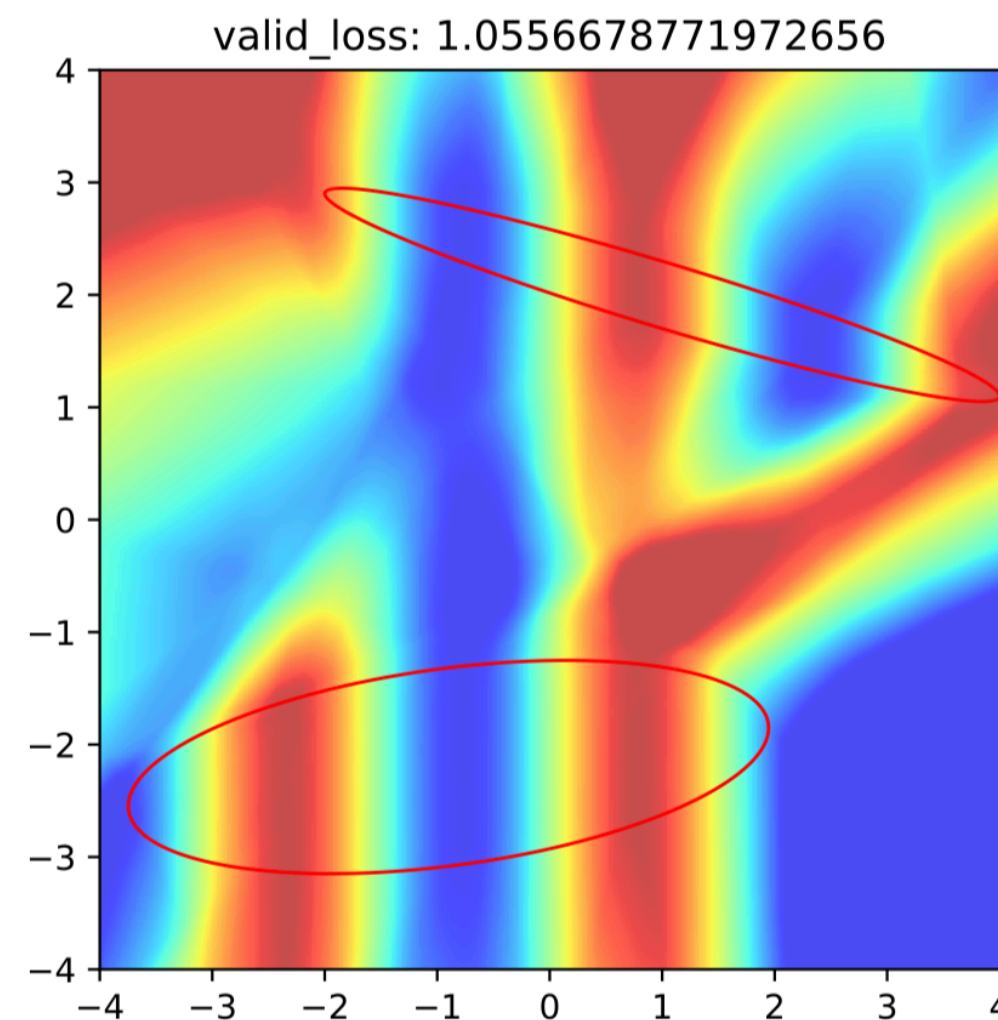
Ground Truth



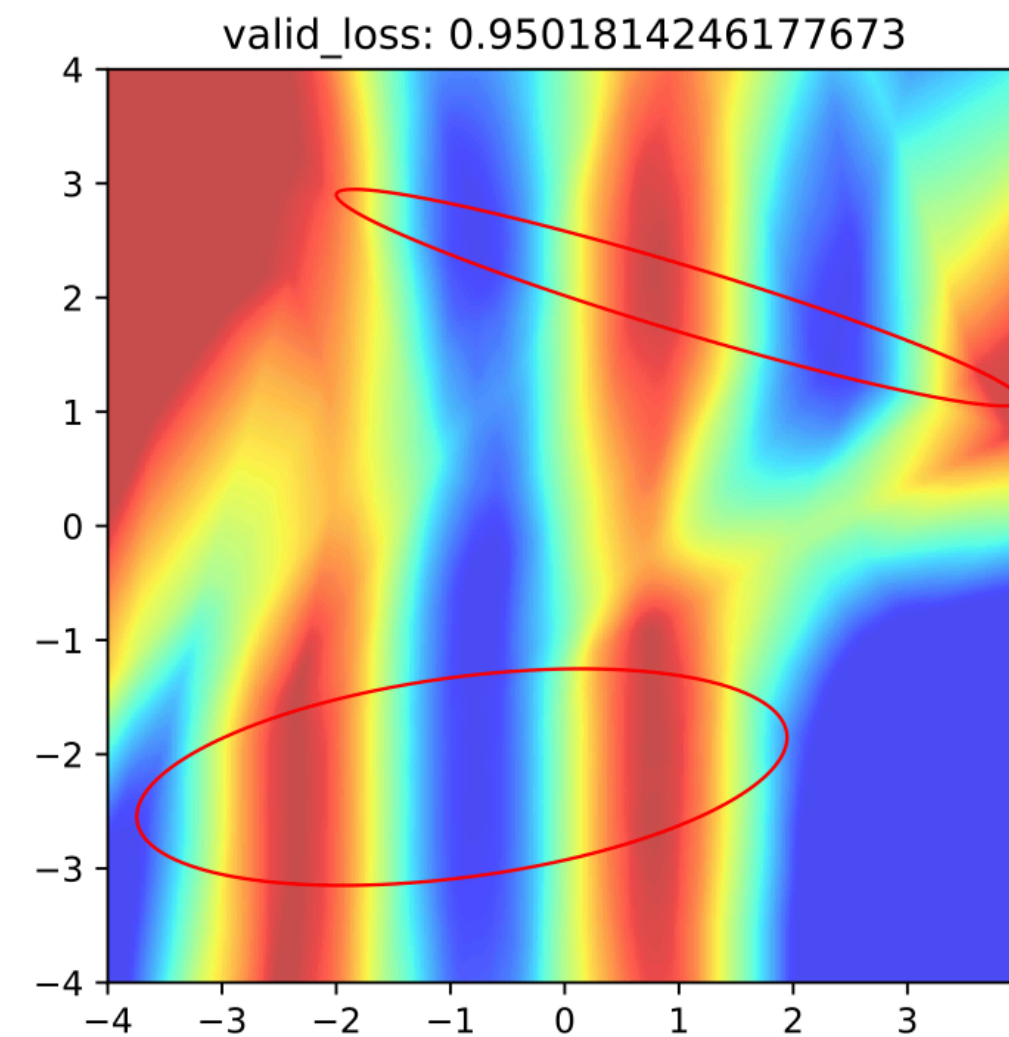
ERM



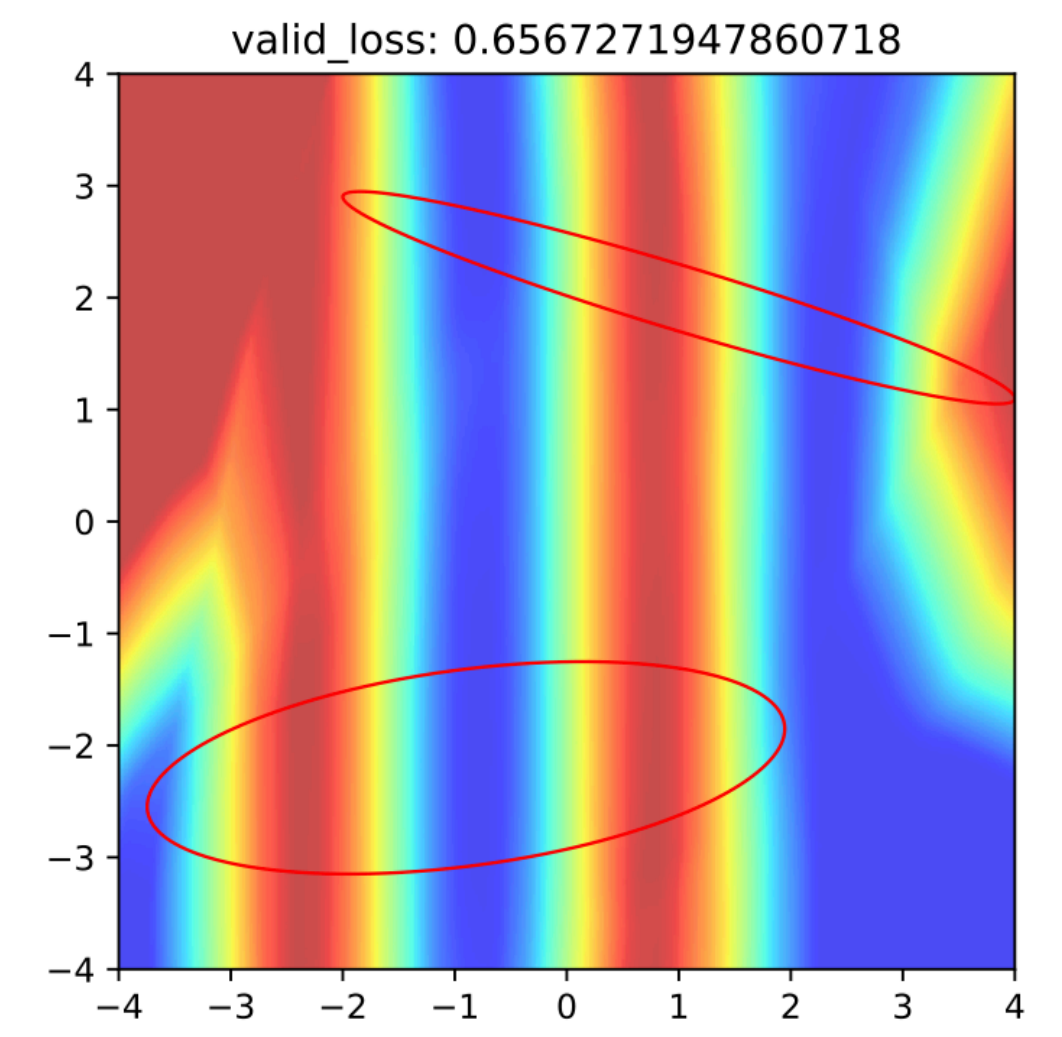
IRMv1



VREx



IRMx



PAIR

# PAIR as the optimizer

Table 2: OOD generalization performances on WILDS benchmark.

	CAMELYON17	CIVILCOMMENTS	FMoW	IWILDCAM	POVERTYMAP	RXR1	AVG. RANK( $\downarrow$ ) <sup>†</sup>
	Avg. acc. (%)	Worst acc. (%)	Worst acc. (%)	Macro F1	Worst Pearson r	Avg. acc. (%)	
ERM	70.3 ( $\pm 6.4$ )	56.0 ( $\pm 3.6$ )	32.3 ( $\pm 1.25$ )	30.8 ( $\pm 1.3$ )	0.45 ( $\pm 0.06$ )	29.9 ( $\pm 0.4$ )	4.50
CORAL	59.5 ( $\pm 7.7$ )	65.6 ( $\pm 1.3$ )	31.7 ( $\pm 1.24$ )	<b>32.7</b> ( $\pm 0.2$ )	0.44 ( $\pm 0.07$ )	28.4 ( $\pm 0.3$ )	5.50
GroupDRO	68.4 ( $\pm 7.3$ )	70.0 ( $\pm 2.0$ )	30.8 ( $\pm 0.81$ )	23.8 ( $\pm 2.0$ )	0.39 ( $\pm 0.06$ )	23.0 ( $\pm 0.3$ )	6.83
IRMv1	64.2 ( $\pm 8.1$ )	66.3 ( $\pm 2.1$ )	30.0 ( $\pm 1.37$ )	15.1 ( $\pm 4.9$ )	0.43 ( $\pm 0.07$ )	8.2 ( $\pm 0.8$ )	7.67
V-REx	71.5 ( $\pm 8.3$ )	64.9 ( $\pm 1.2$ )	27.2 ( $\pm 0.78$ )	27.6 ( $\pm 0.7$ )	0.40 ( $\pm 0.06$ )	7.5 ( $\pm 0.8$ )	7.00
Fish	74.3 ( $\pm 7.7$ )	73.9 ( $\pm 0.2$ )	34.6 ( $\pm 0.51$ )	24.8 ( $\pm 0.7$ )	0.43 ( $\pm 0.05$ )	10.1 ( $\pm 1.5$ )	4.33
LISA	<b>74.7</b> ( $\pm 6.1$ )	70.8 ( $\pm 1.0$ )	33.5 ( $\pm 0.70$ )	24.0 ( $\pm 0.5$ )	<b>0.48</b> ( $\pm 0.07$ )	<b>31.9</b> ( $\pm 0.8$ )	2.67
IRMX	67.0 ( $\pm 6.6$ )	74.3 ( $\pm 0.8$ )	33.7 ( $\pm 0.78$ )	26.6 ( $\pm 0.9$ )	0.45 ( $\pm 0.04$ )	28.7 ( $\pm 0.2$ )	4.00
<b>PAIR-o</b>	74.0 ( $\pm 7.0$ )	<b>75.2</b> ( $\pm 0.7$ )	<b>35.5</b> ( $\pm 1.13$ )	27.9 ( $\pm 0.7$ )	0.47 ( $\pm 0.06$ )	28.8 ( $\pm 0.1$ )	<b>2.17</b>

<sup>†</sup>Averaged rank is reported because of the dataset heterogeneity. A lower rank is better.

PAIR re-empowers IRMv1 and achieves new state-of-the-arts across **6 challenging realistic datasets**.

# PAIR as the model selector

Table 3: OOD generalization performances using DOMAINBED evaluation protocol.

	PAIR-s	COLOREDMNIST <sup>†</sup>				PACS <sup>‡</sup>					TERRAINCOGNITA <sup>†</sup>				
		+90%	+80%	10%	$\Delta$ wr.	A	C	P	S	$\Delta$ wr.	L100	L38	L43	L46	$\Delta$ wr.
ERM		71.0	<b>73.4</b>	10.0		87.2	79.5	95.5	76.9		46.7	<b>41.8</b>	57.4	39.7	
DANN		71.0	<b>73.4</b>	10.0		86.5	79.9	97.1	75.3		46.1	41.2	56.7	35.6	
DANN	✓	71.6	73.3	10.9	+0.9	87.0	81.4	96.8	77.5	+2.2	43.1	41.1	55.2	38.7	+3.1
GroupDRO		72.6	73.1	9.9		87.7	82.1	98.0	79.6		48.4	40.3	57.9	40.0	
GroupDRO	✓	<b>72.7</b>	73.2	13.0	+3.1	86.7	<b>83.2</b>	<b>97.8</b>	81.4	+1.8	48.4	40.3	57.9	40.0	+0.0
IRMv1		72.3	72.6	9.9		82.3	80.8	95.8	78.9		48.4	35.6	55.4	40.1	
IRMv1	✓	67.4	64.8	<b>24.2</b>	+14.3	85.3	81.7	97.4	79.7	+0.8	40.4	38.3	48.8	37.0	+1.4
Fishr		72.2	73.1	9.9		<b>88.4</b>	82.2	97.7	81.6		49.2	40.6	57.9	40.4	
Fishr	✓	69.1	70.9	22.6	+12.7	87.4	82.6	97.5	<b>82.2</b>	+0.6	<b>51.0</b>	40.7	<b>58.2</b>	<b>40.8</b>	+0.3

<sup>†</sup>Using the training domain validation accuracy. <sup>‡</sup>Using the test domain validation accuracy.

PAIR-s substantially improves the worst environment performance of all representative OOD methods up to **10%**.

# Summary

We provided a new understanding of the optimization dilemma in OOD generalization from the Multi-Objective Optimization perspective.

We attributed the failures of OOD optimization to the compromised robustness of relaxed OOD objectives and the unreliable optimization scheme.

We highlighted the importance of trading-off the ERM and OOD objectives and proposed a new optimization scheme PAIR to mitigate the dilemma.



Paper



Code

## Thank you!

Contact: [yqchen@cse.cuhk.edu.hk](mailto:yqchen@cse.cuhk.edu.hk)