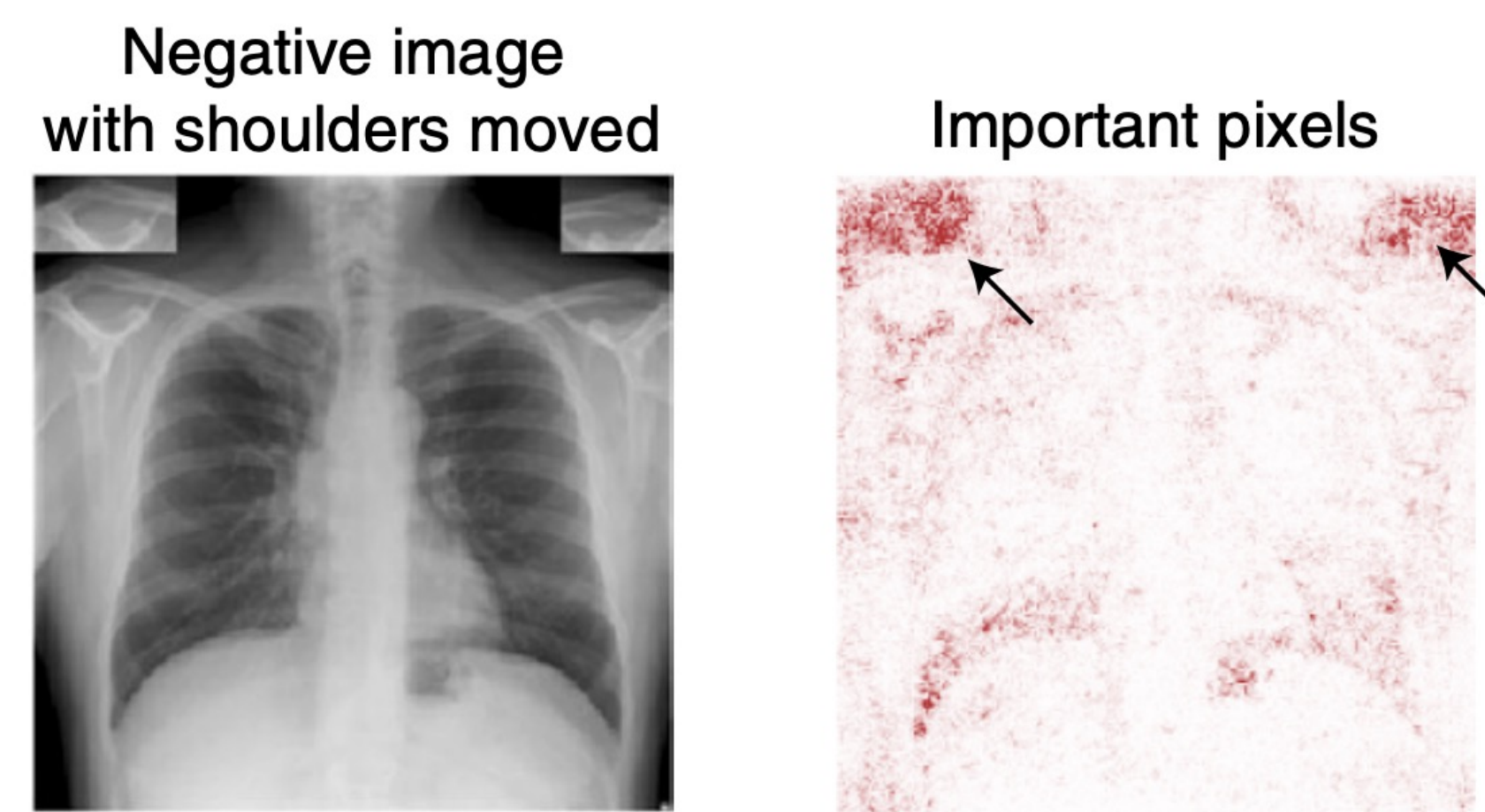




1 Simplicity Bias and Domain Shifts

DNNs learn simple/biased shortcuts rather than complex/stable features.

Ex: to detect Covid, DNNs analyze body/shoulder positions rather than lung fields.



⇒ lack of robustness, failures under shifts.
Main challenge: we want to detect causation rather than correlation.

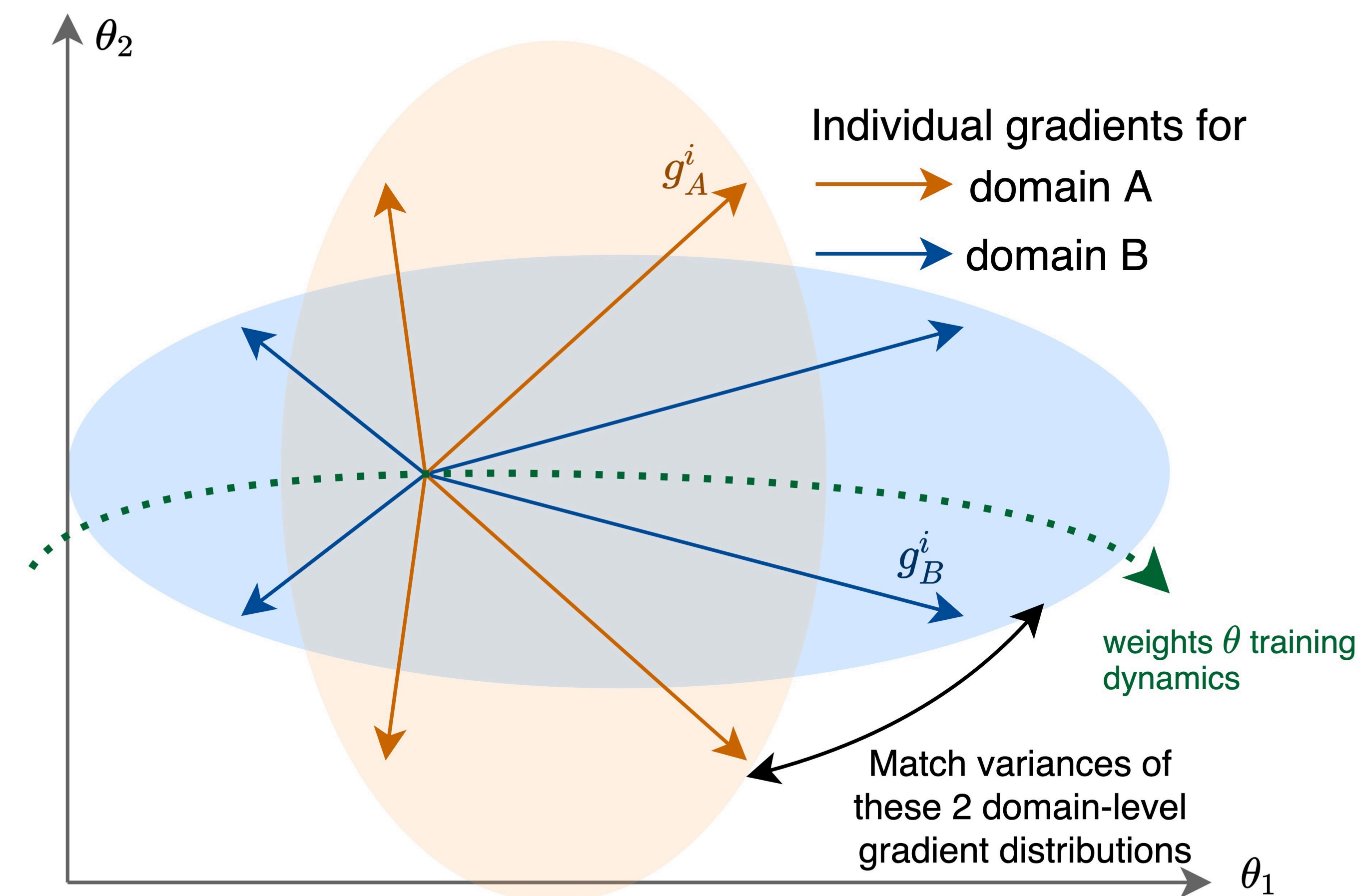
2 Invariance Paradigm

Assumption: the causal mechanism is invariant across the various training domains.

How to enforce *invariance* across domains?

- In *features* extracted by the encoder $f_\phi(x)$. Ex: CORAL(AAAI16), DANN(JMLR16).
- In *predictors* with invariant risks $l(f_\theta(x), y)$. Ex: IRM, V-Rex(ICML21).
- In *gradients* of the loss w.r.t. the weights of the network $\nabla_\theta l(f_\theta(x), y)$. Ex: Fish(ICLR22).

3 Fishr: Invariant Gradient Variances ... and also invariant risks and Hessians



$$\mathcal{L}_{Fishr} = \mathcal{L}_{ERM} + \lambda \| \text{Var}(G_A) - \text{Var}(G_B) \|_2^2,$$

λ controls the strength of our regularization matching the variance of gradients $G_e = [\nabla_\theta l(f_\theta(x_e^i), y_e^i)]_{i=1}^{n_e}$ across domains $e \in \{A, B\}$.

Theoretical analysis

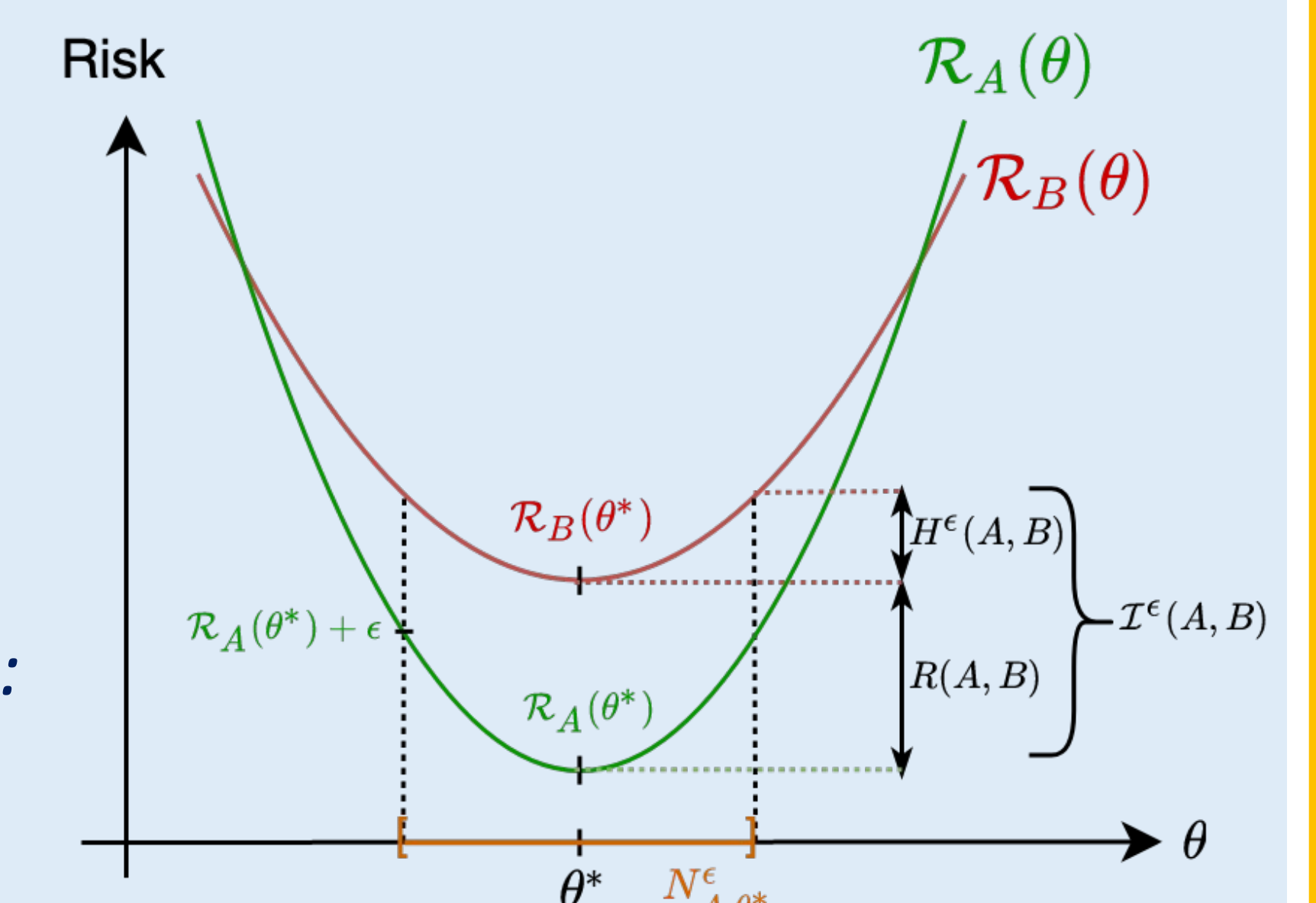
Gradient variances as a proxy to match:

- risks: $\mathcal{R}_e(\theta) = \sum_i l(f_\theta(x_e^i), y_e^i) / n_e$,
- Hessians: $\mathcal{H}_e(\theta) = \sum_i \nabla_\theta^2 l(f_\theta(x_e^i), y_e^i) / n_e$ (proof via the Fisher Information Matrix).

We prove Fishr reduces *domain inconsistencies*:

$$\mathcal{J}^\epsilon(A, B) = \max_{\theta \in N_{A, \theta^*}^\epsilon} |R_B(\theta) - R_A(\theta^*)|$$

in a ϵ neighbourhood N_{A, θ^*}^ϵ around weights θ^* .



4 State-of-the-art Performances on DomainBed Benchmark (at almost no computational overhead)

Dataset	Domains
Colored MNIST	+90% +80% -90% (degree of correlation between color and label)
Rotated MNIST	0° 15° 30° 45° 60° 75°
VLCS	Caltech101 LabelMe SUN09 VOC2007
PACS	Art Cartoon Photo Sketch
Office-Home	Art Clipart Product Photo
Terra Incognita	L100 L38 L43 L46 (camera trap location)
DomainNet	Clipart Infographic Painting QuickDraw Photo Sketch

Algo.	Invariance	Acc. ↑								Rank ↓	
		cMNIST	rMNIST	VLCS	PACS	OHome	TerraI	DNet	Avg	Avg	
ERM	✗	57.8	97.8	77.6	86.7	66.4	53.0	41.3	68.7	9.1	
CORAL	Features	58.6	98.0	77.7	<u>87.1</u>	68.4	52.8	<u>41.8</u>	<u>69.2</u>	<u>4.6</u>	
DANN		57.0	<u>97.9</u>	79.7	85.2	65.3	50.6	38.3	67.7	11.9	
IRM	Predictors	<u>67.7</u>	97.5	76.9	84.5	63.0	50.5	28.0	66.9	14.7	
V-REx		67.0	<u>97.9</u>	78.1	87.2	65.7	51.4	30.1	68.2	7.7	
Fish	Gradients	61.8	<u>97.9</u>	77.8	85.8	66.0	50.8	43.4	69.1	8.4	
Fishr		68.8	97.8	<u>78.2</u>	86.9	<u>68.2</u>	53.6	<u>41.8</u>	70.8	3.9	