

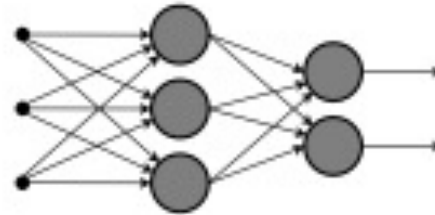
Fishr

Invariant Gradient Variances for Out-of-Distribution Generalization

Alexandre Ramé (PhD)
Corentin Dancette (PhD)
Matthieu Cord (Professor)



➤ DNNs to detect Covid from medical scans ...



Positive
Negative

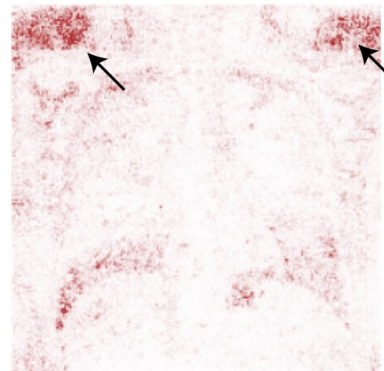
➤ ... but DNNs memorized biased shortcuts

- age: children vs. adults
- position: standing up vs. lying down

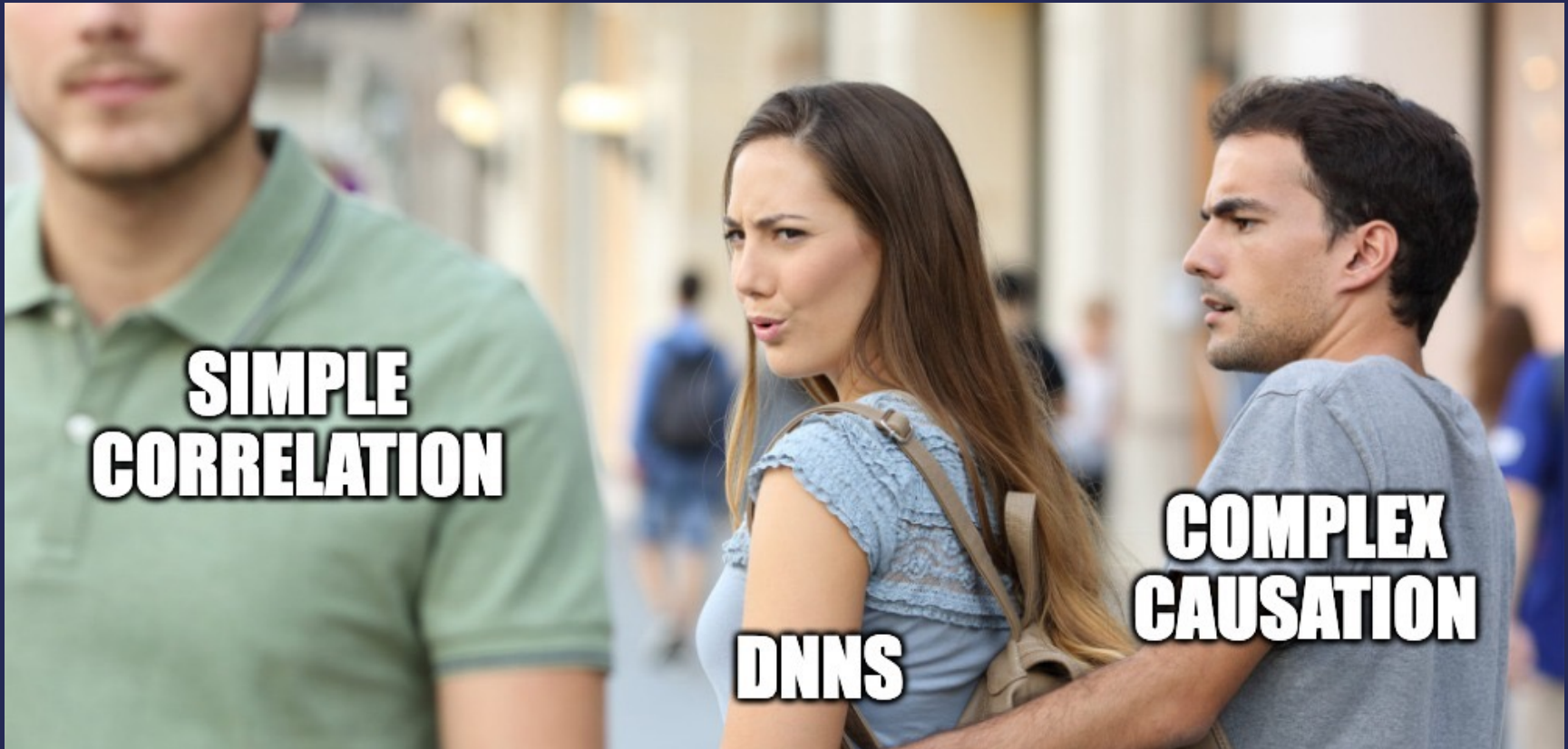
Negative image
with shoulders moved



Important pixels



(rather than analyzing lung fields)



⇒ Simplicity bias deteriorates out-of-distribution generalization

Framework: Training with Multiple Domains

Invariance paradigm: the causal mechanism is invariant across domains



ERM (empirical risk minimization) and invariance approaches

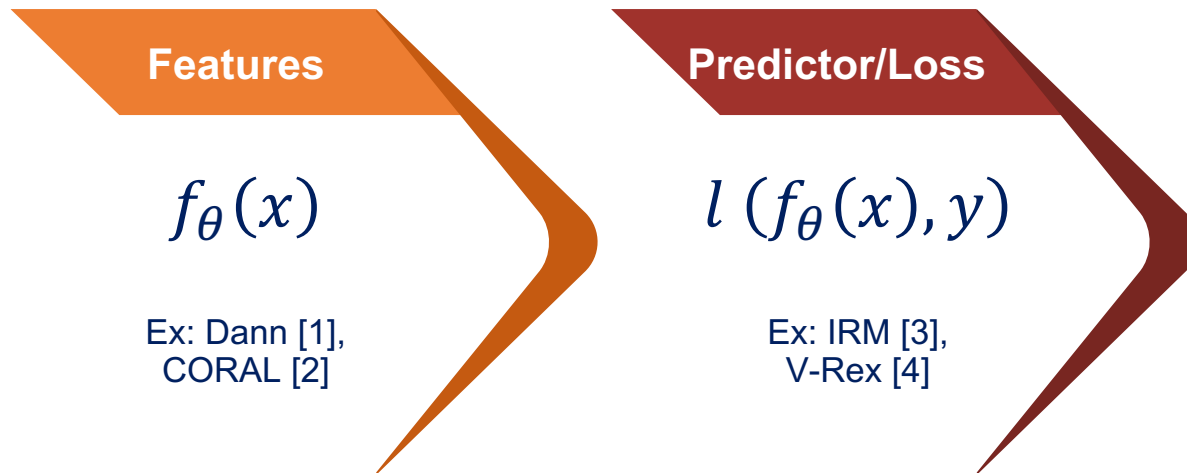
$$\mathcal{L}_{ERM} = R_A + R_B \quad \leftarrow \text{Sum of domain-level risks}$$

As most works, we add an invariance regularization on top of ERM:

$$\mathcal{L}_{invariance} = \mathcal{L}_{ERM} + \underbrace{\lambda \times distance(\phi_A, \phi_B)}_{\text{Invariance regularization}}$$

Diagram annotations:
- "Hyperparameter" points to λ
- "Domain-level statistics" points to ϕ_A and ϕ_B

➤ Invariance in features or losses



[1] Domain-Adversarial Training of Neural Networks. Ganin *et al.*, JMLR 2016

[2] Deep coral: Correlation alignment for deep domain adaptation. Sun and Saenko, ECCV 2016

[3] Invariant risk minimization. Arjovsky *et al.*, 2019

[4] Out-of-distribution generalization via risk extrapolation. Krueger *et al.*, ICML 2021

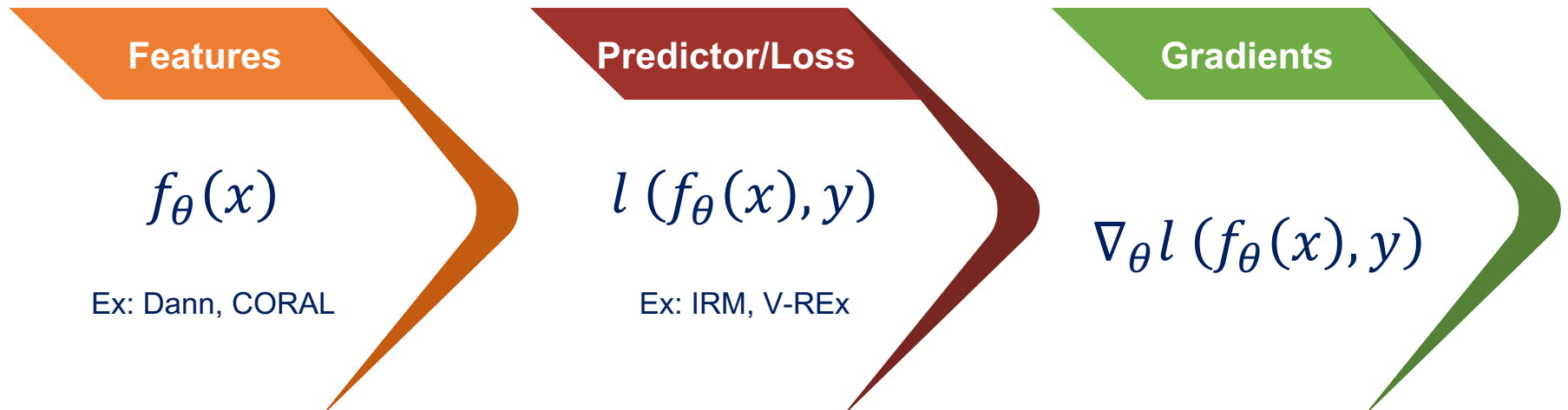


DomainBed

Dataset	Domains					
Colored MNIST	+90%	+80%	-90%			
<i>(degree of correlation between color and label)</i>						
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		
<i>(camera trap location)</i>						
DomainNet	Clipart	Infographic	Painting	QuickDraw	Photo	Sketch

No traditional methods outperform ERM in DomainBed

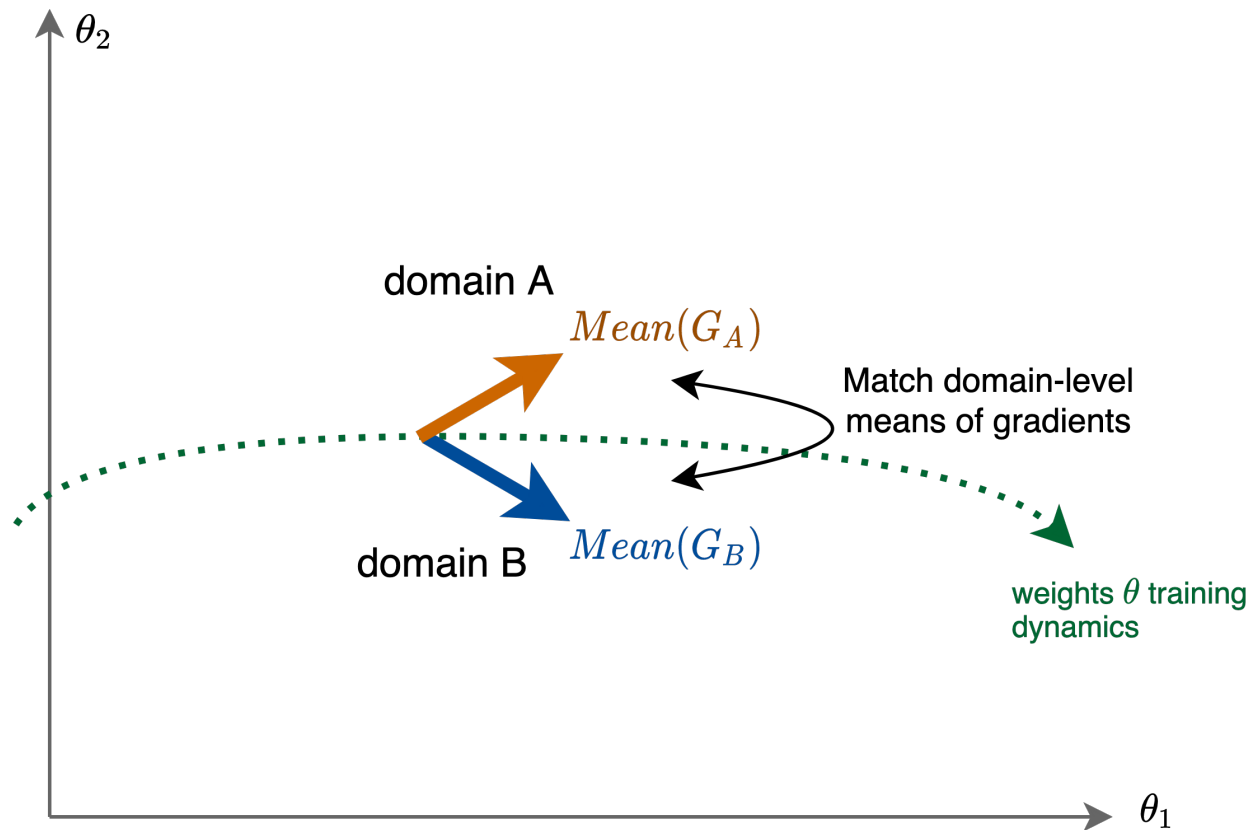
➤ Invariance in gradients !



$$G_e = \left[\nabla_{\theta} l(f_{\theta}(x_e^i), y_e^i) \right]_{i=1}^{n_e} \text{ for domain } e \in \{A, B\}$$

➤ Matching domain-level gradient means

$$\text{Regularization: } \| \text{Mean}(G_A) - \text{Mean}(G_B) \|_2$$

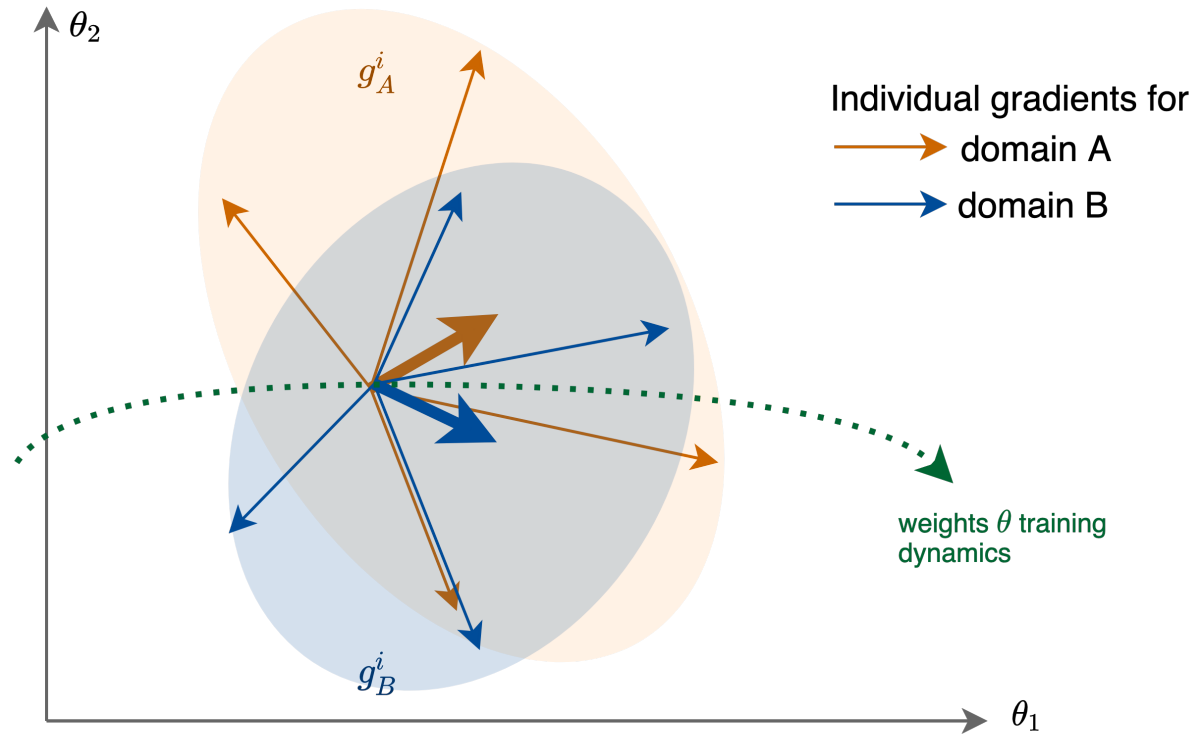


[1] Out-of-distribution generalization with maximal invariant predictor. Koyama and Yamaguchi, 2020

[2] Fish: Gradient matching for domain generalization. Shi *et al.*, ICLR 2022



Gradient distributions richer than gradient means



[1] Gradient diversity: a key ingredient for scalable distributed learning. Yin *et al.*, AISTATS 2018

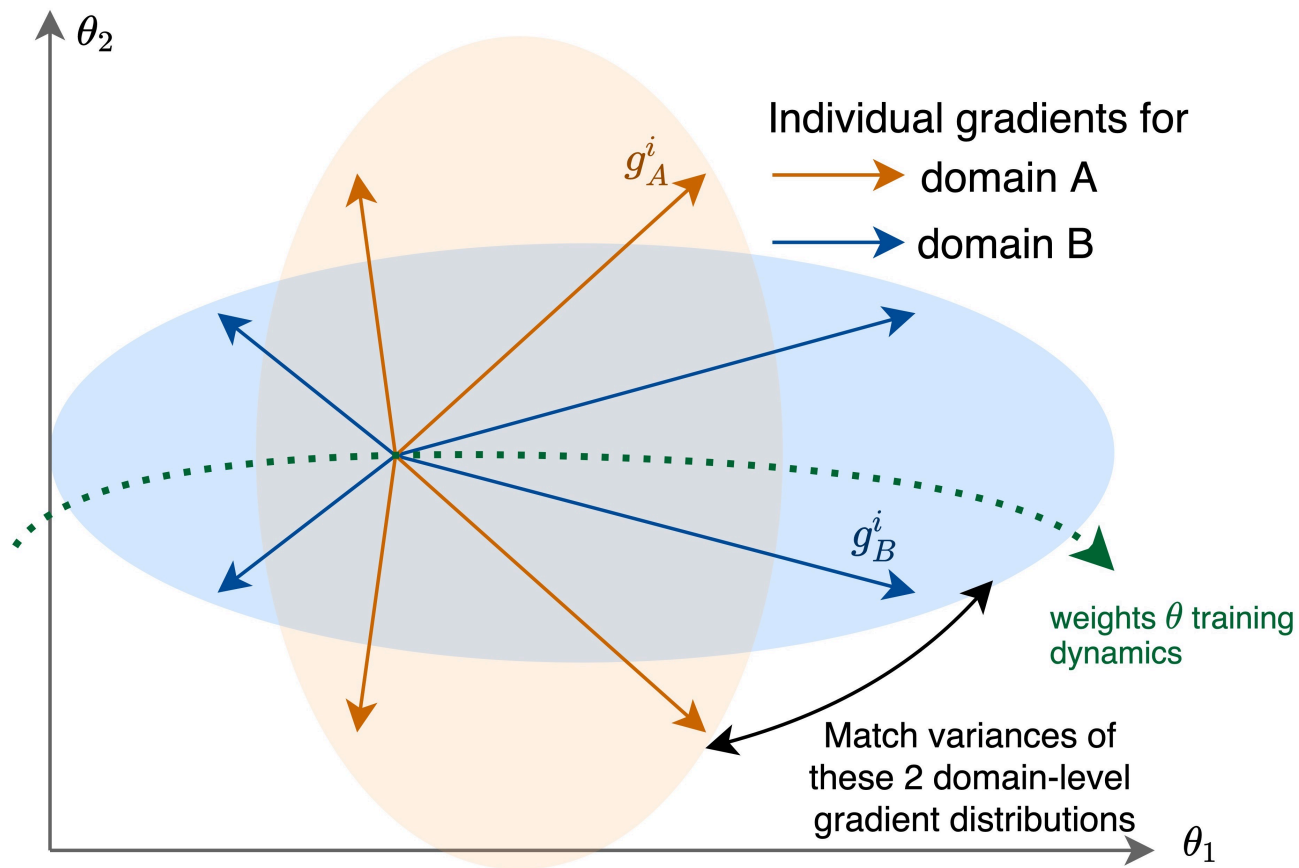
[2] The impact of neural network overparameterization on gradient confusion and stochastic gradient descent. Sankararaman *et al.*, ICML 2020



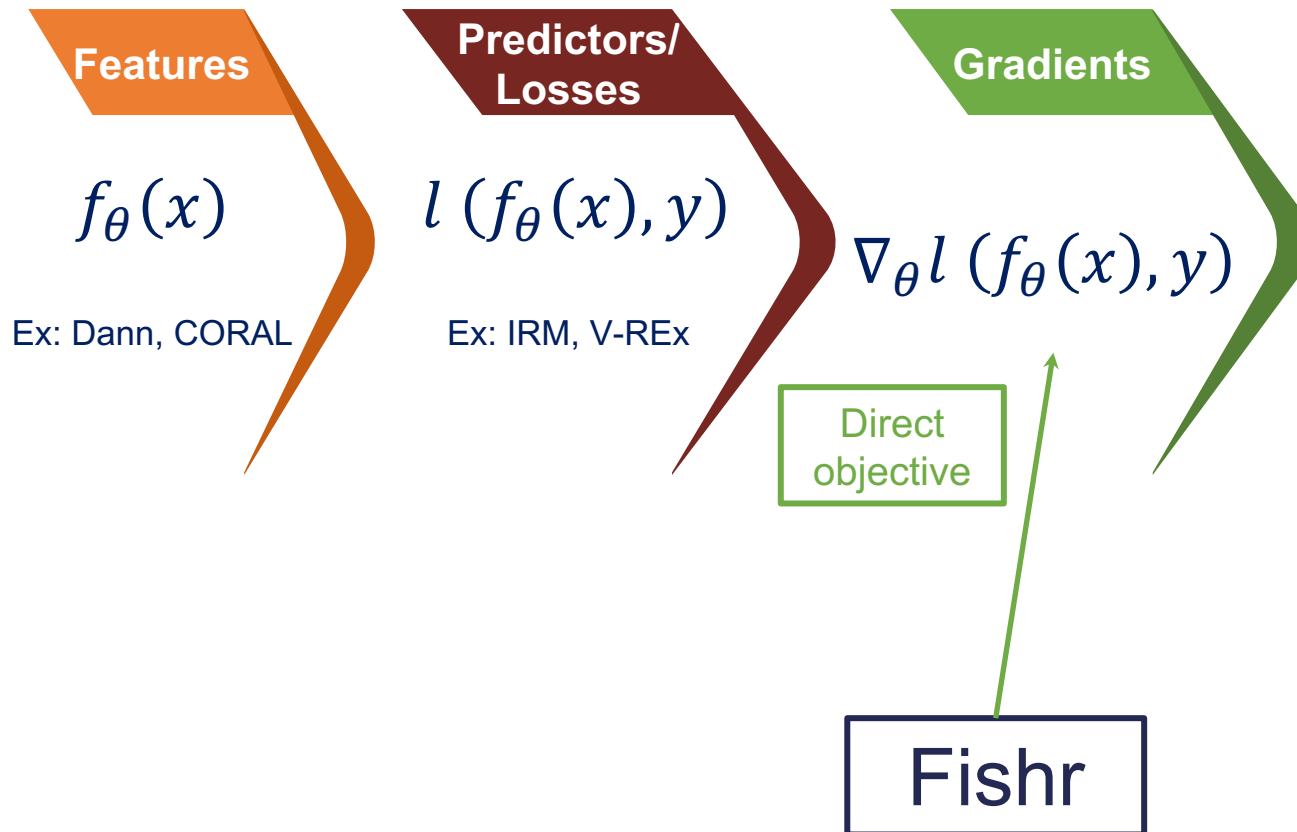
Fishr: invariant gradient variances

$$\text{Regularization: } \| \text{Var}(G_A) - \text{Var}(G_B) \| \frac{2}{2}$$

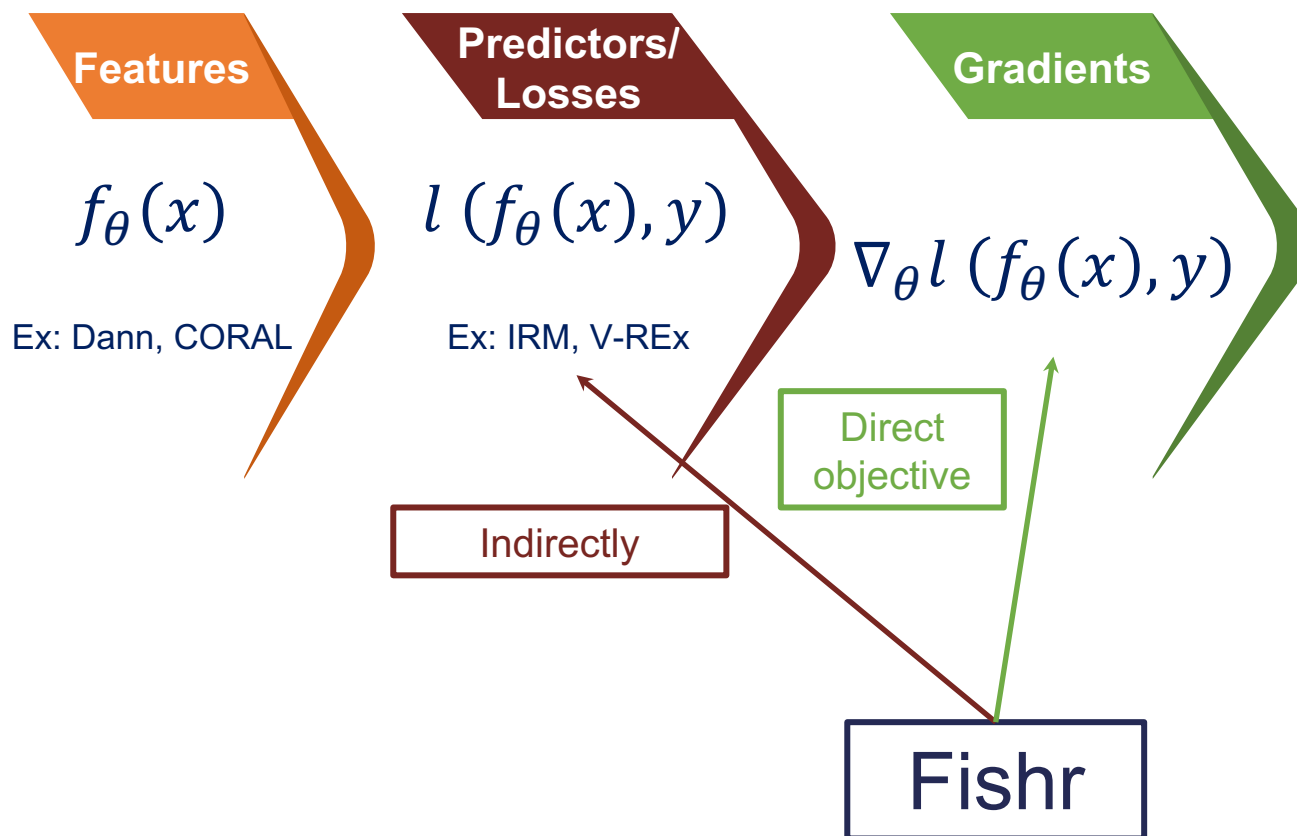
$$\text{where for } e \in \{A, B\}, G_e = [\nabla_{\theta} l(f_{\theta}(x_e^i), y_e^i)]_{i=1}^{n_e}$$



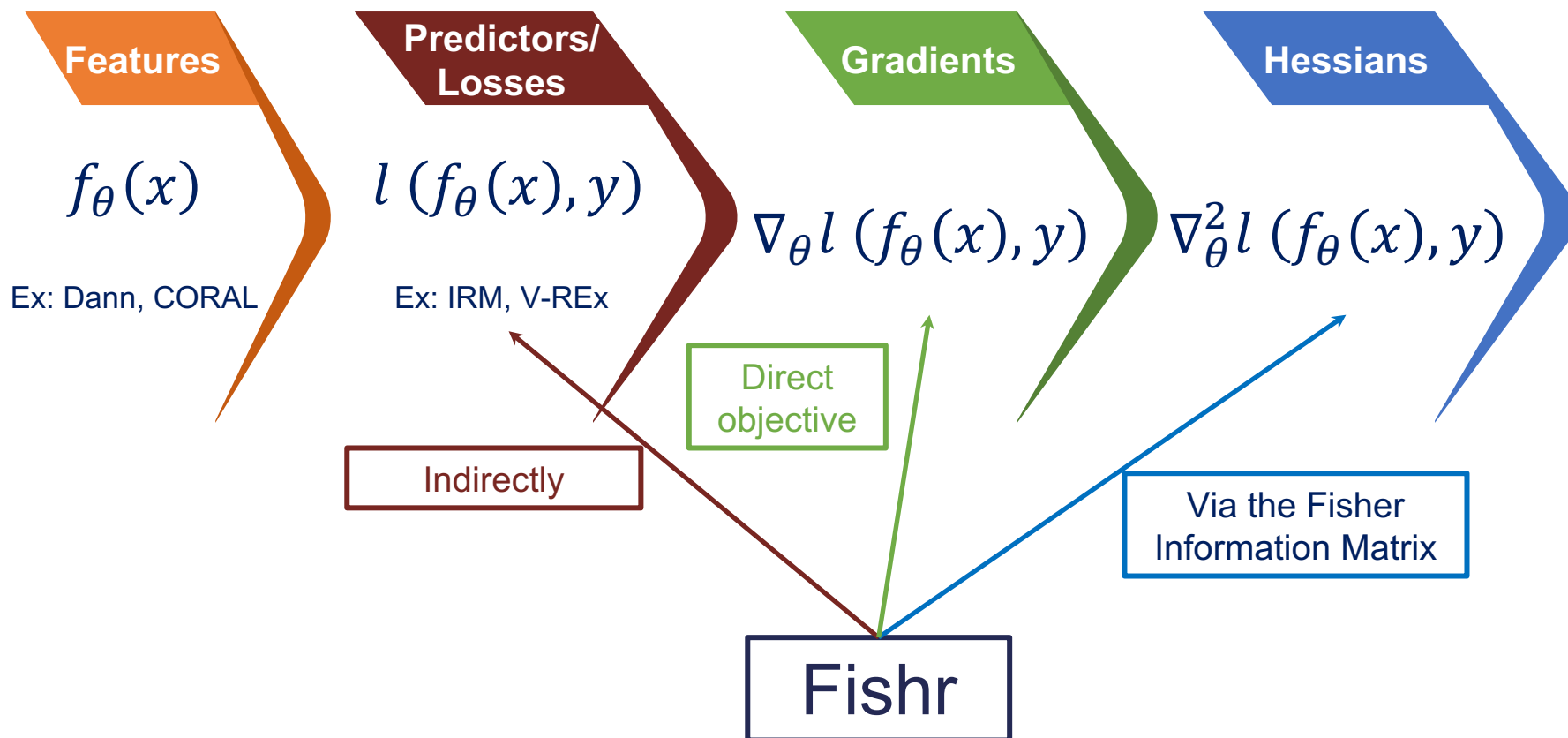
➤ Invariant gradients



➤ Invariant gradients thus invariant losses

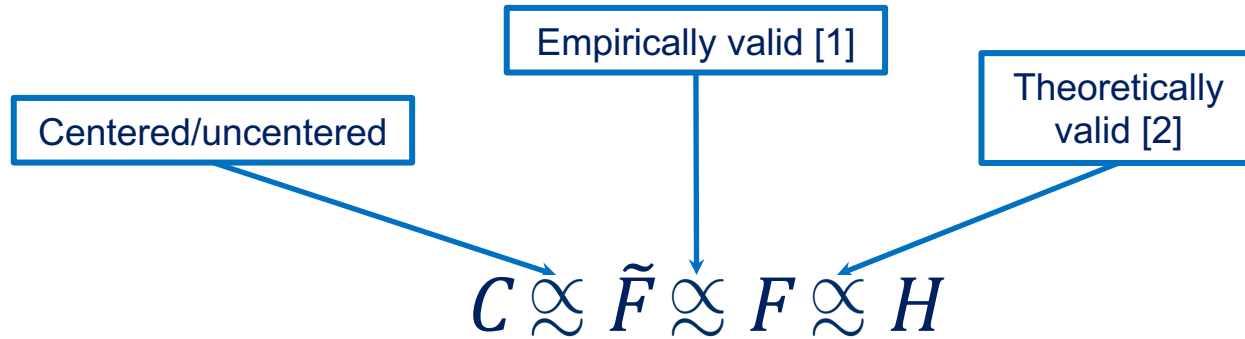


➤ Invariant gradients thus invariant losses ... and invariant Hessians!





Gradient covariance approximates the Hessian (via the Fisher Information Matrix)



Name	Statistics	Formula
C	Gradient Covariance	$Cov(G)$
\tilde{F}	Empirical Fisher Information Matrix	$\sum_{i=1}^n \nabla_{\theta} \log(p_{\theta}(y^i x^i)) \nabla_{\theta} \log(p_{\theta}(y^i x^i))^T$
F	True Fisher Information Matrix	$\sum_{i=1}^n \mathbb{E}_{\hat{y} \sim P_{\theta}(\cdot x^i)} [\nabla_{\theta} \log(p_{\theta}(\hat{y} x^i)) \nabla_{\theta} \log(p_{\theta}(\hat{y} x^i))^T]$
H	Hessian	$\sum_{i=1}^n \nabla_{\theta}^2 l(f_{\theta}(x^i), y^i)$

[1] On the interplay between noise and curvature and its effect on optimization and generalization. Thomas *et al.*, AISTATS 2020

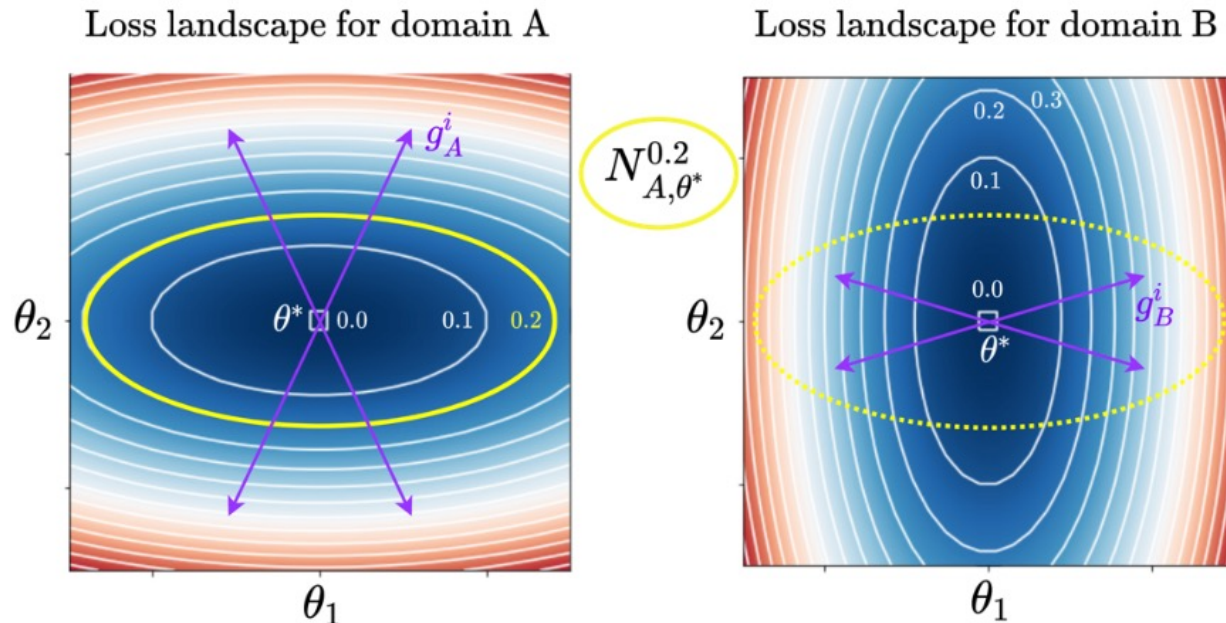
[2] New insights and perspectives on the natural gradient method. Martens, 2014

➤ Fishr matches domain-level loss landscapes

With Fishr at convergence at θ^* , $R_A^{\theta^*} \approx R_B^{\theta^*}$, $G_A^{\theta^*} \approx G_B^{\theta^*}$, $H_A^{\theta^*} \approx H_B^{\theta^*}$

Via a 2nd-order Taylor expansion, \forall weights θ close to θ^* :

$$\begin{aligned} R_A^\theta &\approx R_A^{\theta^*} + (\theta - \theta^*)G_A^{\theta^*} + (\theta - \theta^*)H_A^{\theta^*}(\theta - \theta^*) \\ &\approx R_B^{\theta^*} + (\theta - \theta^*)G_B^{\theta^*} + (\theta - \theta^*)H_B^{\theta^*}(\theta - \theta^*) \\ &\approx R_B^\theta \end{aligned}$$



➤ DomainBed

Reference benchmark for OOD generalization, imposing the *code, datasets, training procedures, hyperparameter search, model selection* etc.

Algo.	Invariance	Acc. ↑								Rank ↓
		cMNIST	rMNIST	VLCS	PACS	OHome	TerraI	DNet	Avg	Avg
ERM	✗	57.8	97.8	77.6	86.7	66.4	<u>53.0</u>	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



Fishr Contributions

❖ Theoretically

- Invariant gradient variances ...
- but also invariant losses and Hessians to align landscapes

❖ Empirically

- Simple and scalable
- State of the Art on DomainBed for OOD generalization

Code available: <https://github.com/alexrame/fishr>

Merci !

