

**ICLR 2022 Workshop on the Elements of Reasoning:
Objects, Structure, and Causality (OSC)**

Causality in Computer Vision

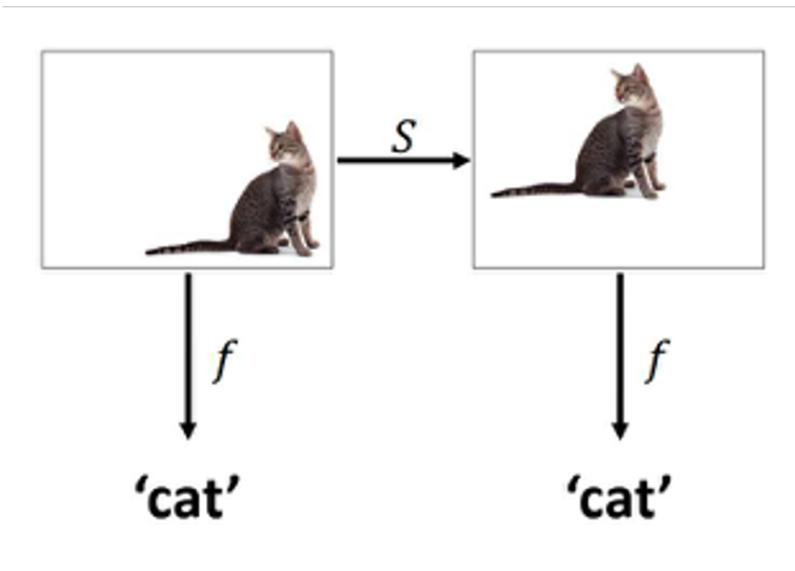
Invariant Learning with Insufficient Data

Qianru SUN

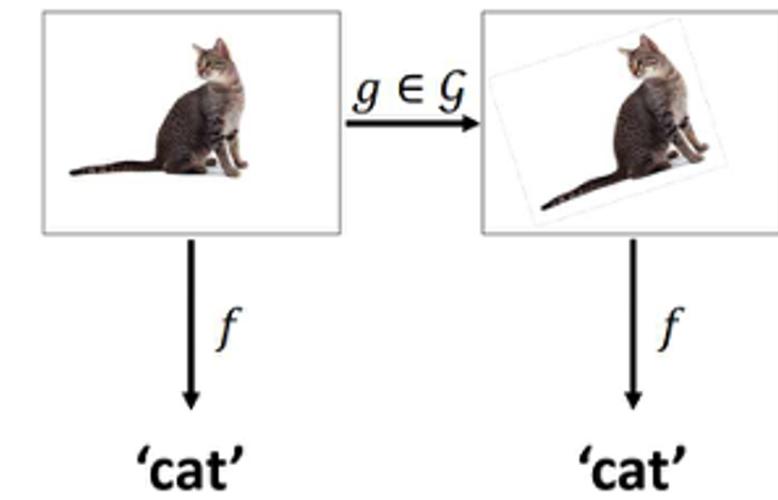
29 Apr 2022

Invariant Learning?

Invariant to Shift

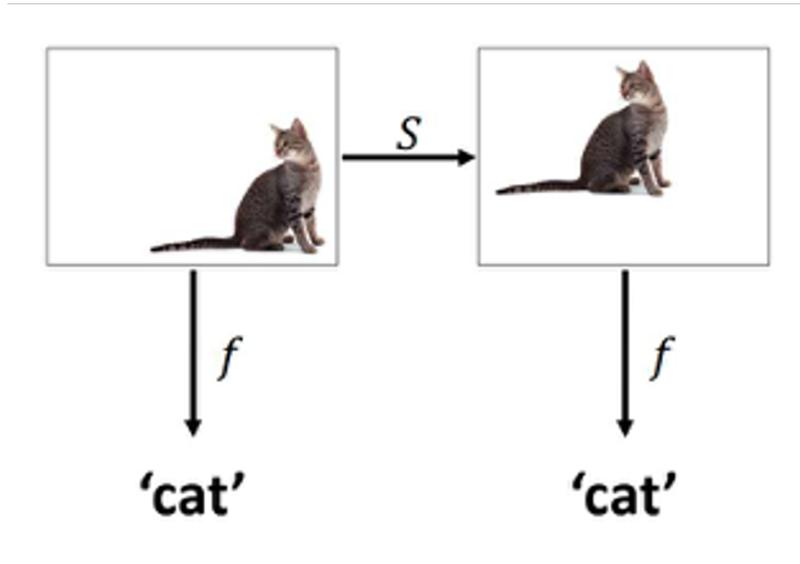


Invariant to Rotation



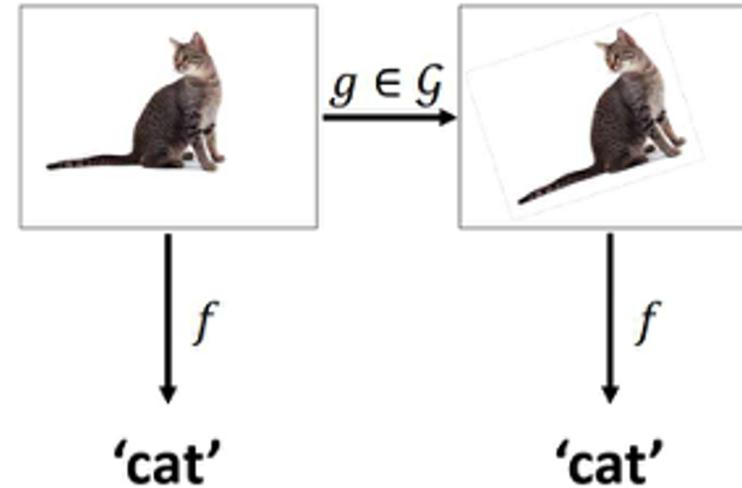
Invariant Learning with Sufficient Data?

Invariant to Shift



Range of Shift covers every position

Invariant to Rotation

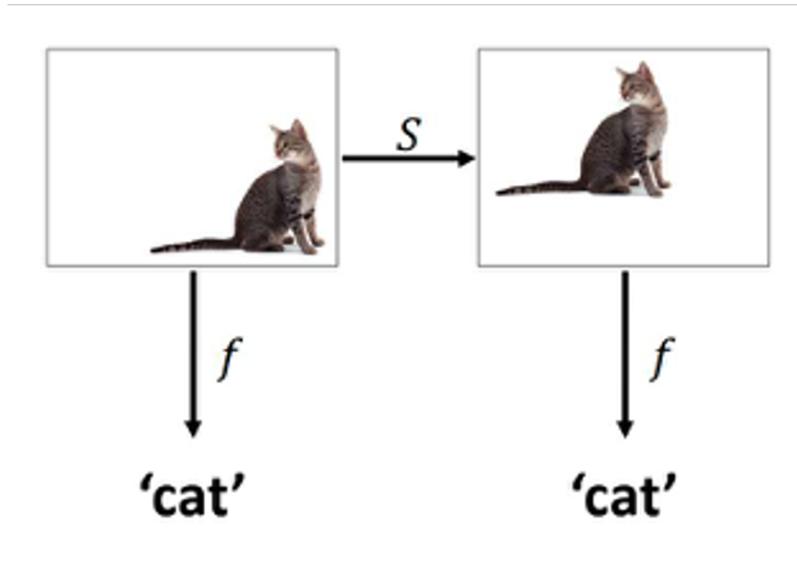


Range of Rotation covers every view angle

Invariant Learning with **Sufficient Data**?

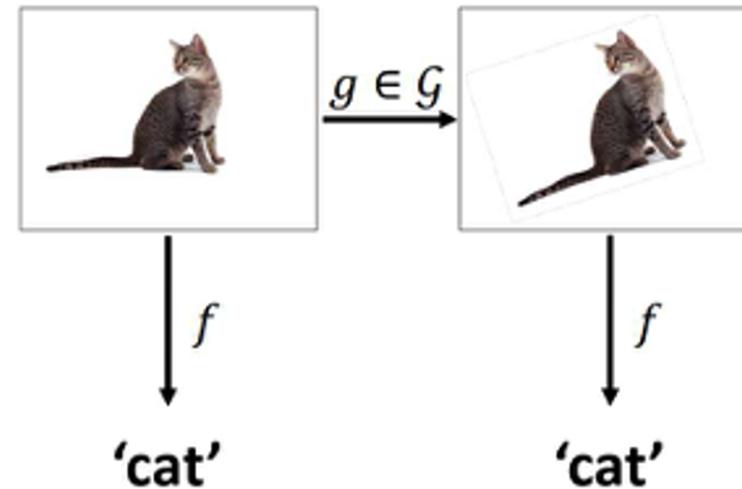
Key Property: Disentangling Class, Position, and View Angle

Invariant to Shift



Range of Shift covers **every position**

Invariant to Rotation

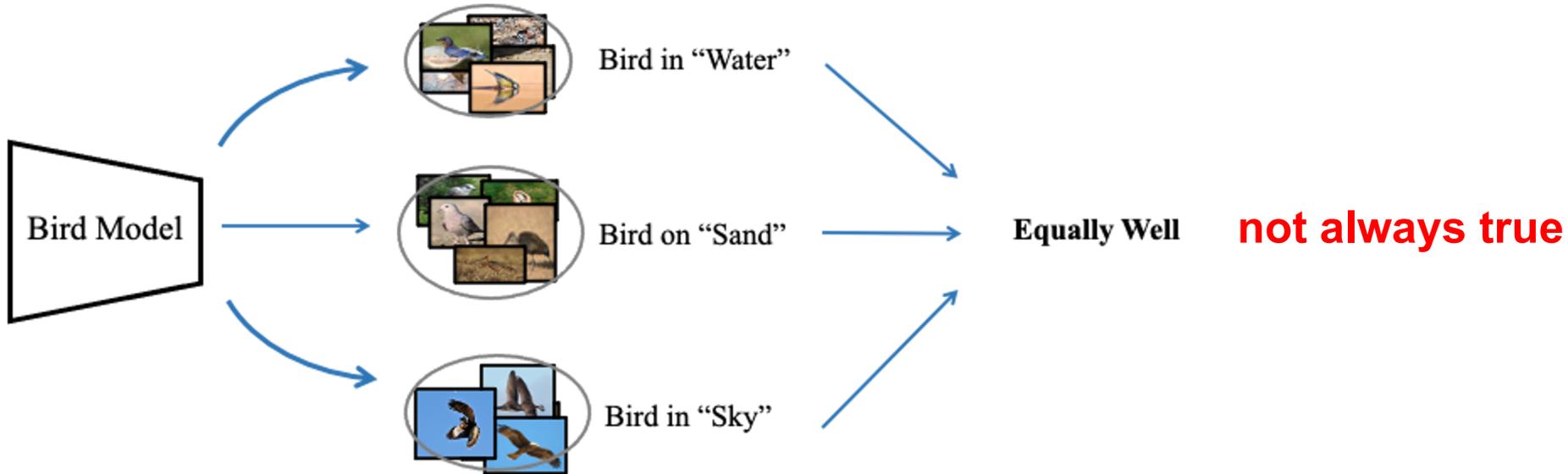


Range of Rotation covers **every view angle**

Invariant Learning with **Sufficient Data**?

Key Property: Disentangling Class, Position, and View Angle

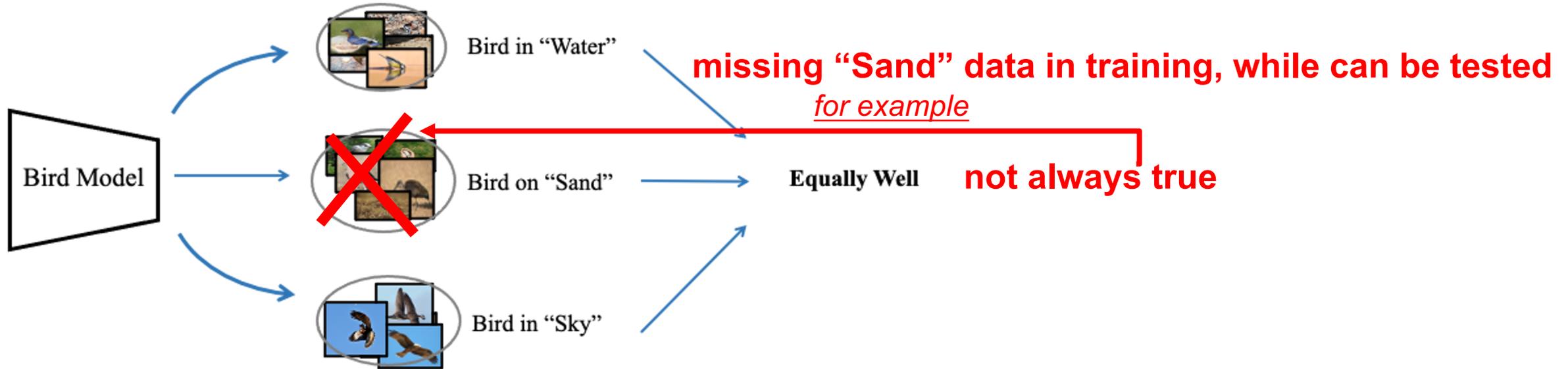
This is, however, not easy to achieve in real-world data ...



Invariant Learning with **Sufficient Data**?

Key Property: Disentangling Class, Position, and View Angle

This is, however, not easy to achieve in real images of complex foregrounds and backgrounds ...

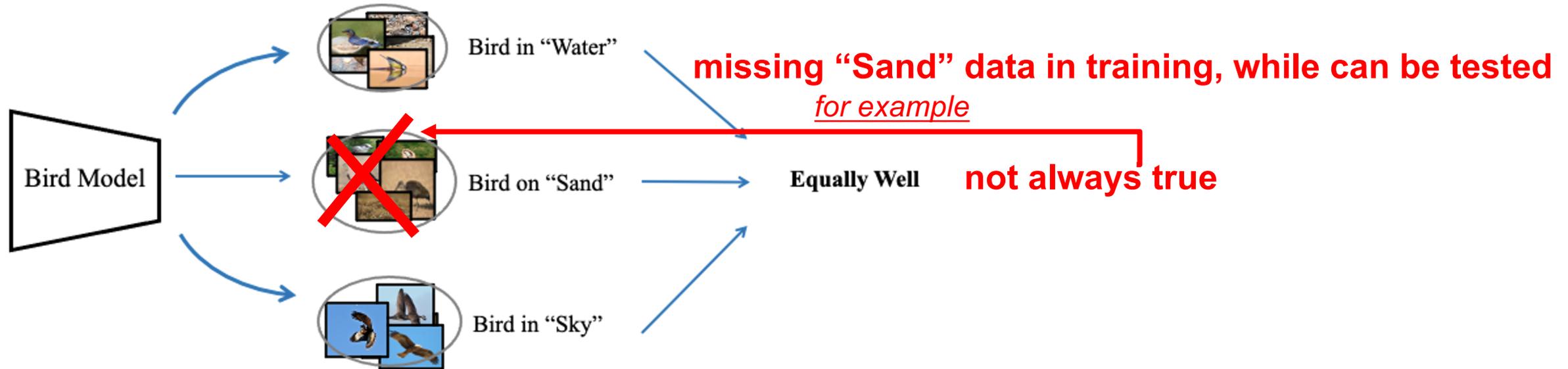


Training data is often **Insufficient**.

Invariant Learning with **Insufficient Data**

Key Property: Disentangling Class, Position, and View Angle

This is, however, not easy to achieve in real images of complex foregrounds and backgrounds ...



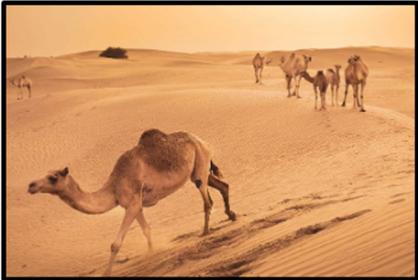
Training data is often **Insufficient**.

Invariant Learning with **Insufficient Data (OOD Data)**

When Model Fails in OOD Generalization

For example, the model's prediction is misled by biased backgrounds:

Training



Camel + Sand



Cow + Grass

Testing



Camel + Grass

Prediction: Cow



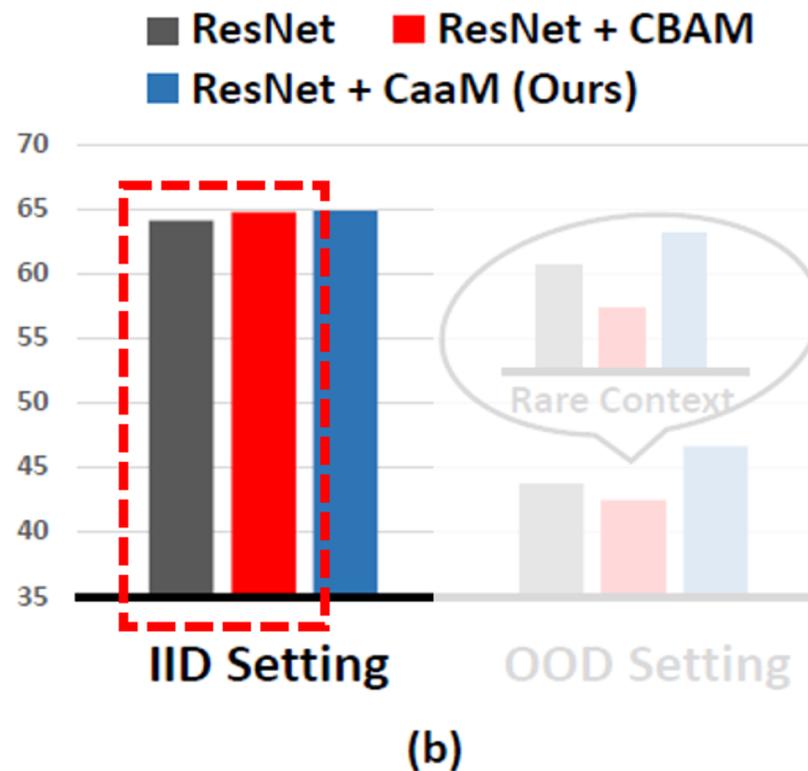
Cow + Sand

Prediction: Camel



When Model Fails in OOD Generalization

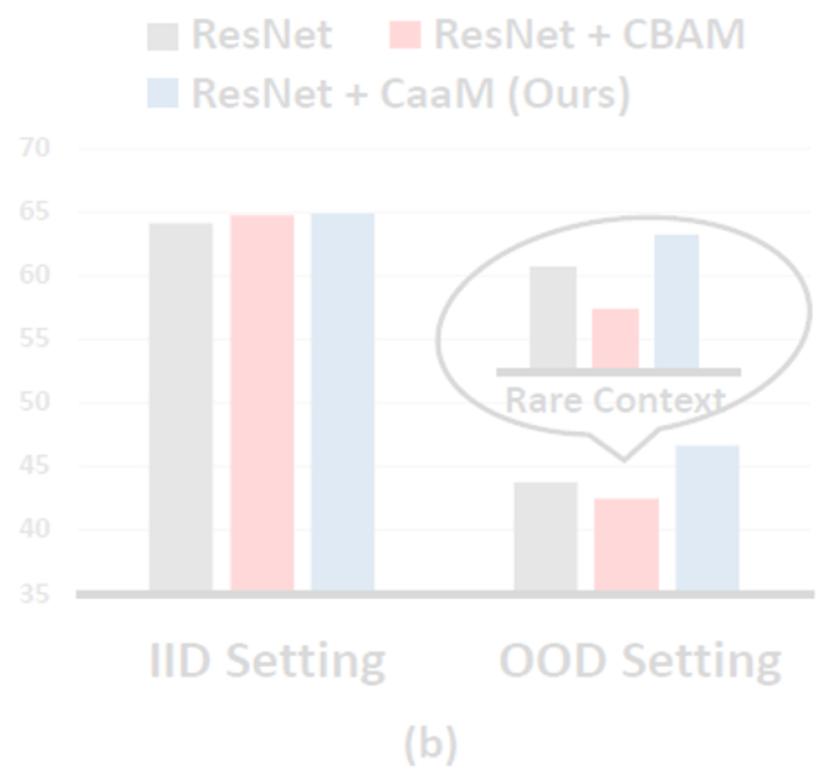
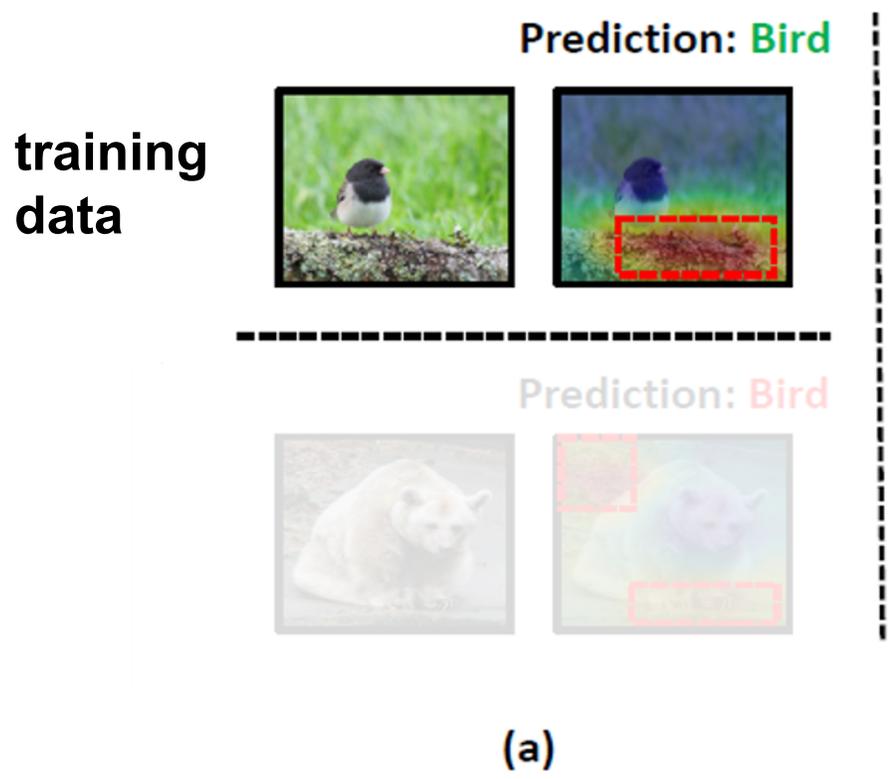
Another example is the model attention is biased to backgrounds:



“Attention is all you need”

When Model Fails in OOD Generalization

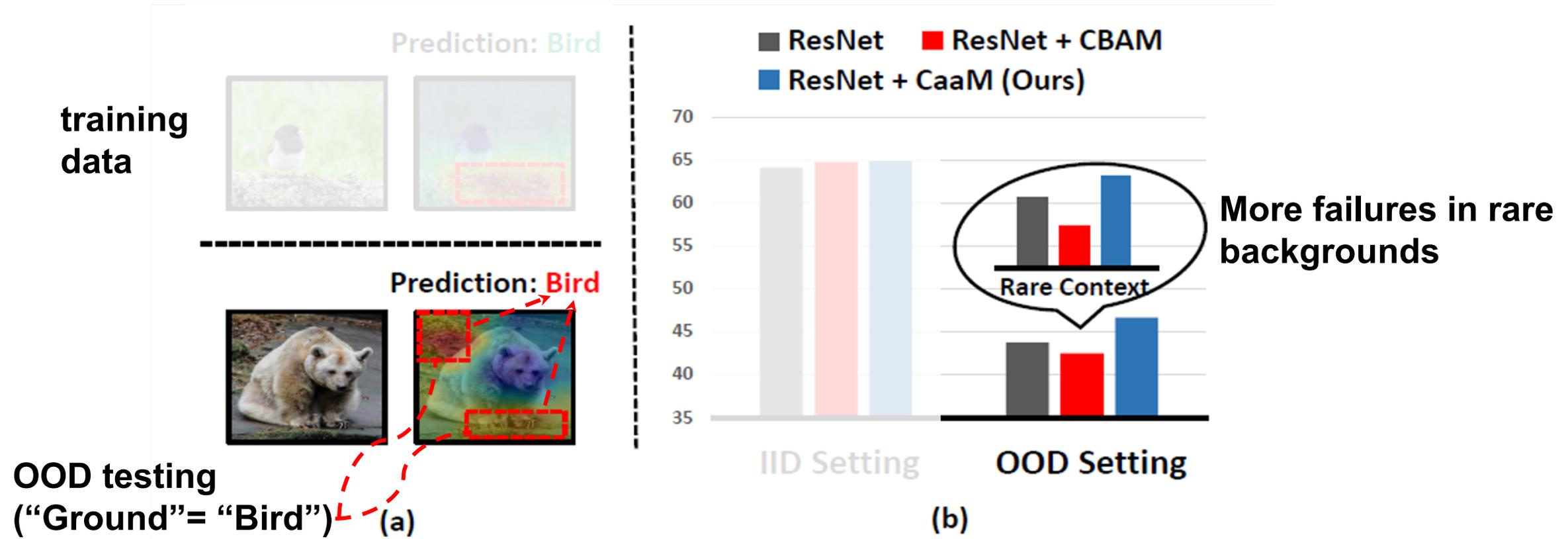
Another example is the model attention is biased to backgrounds:



“Attention is all you need”?

When Model Fails in OOD Generalization

Another example is the model attention is biased to backgrounds:



“Attention is all you need”?

What is the reason?

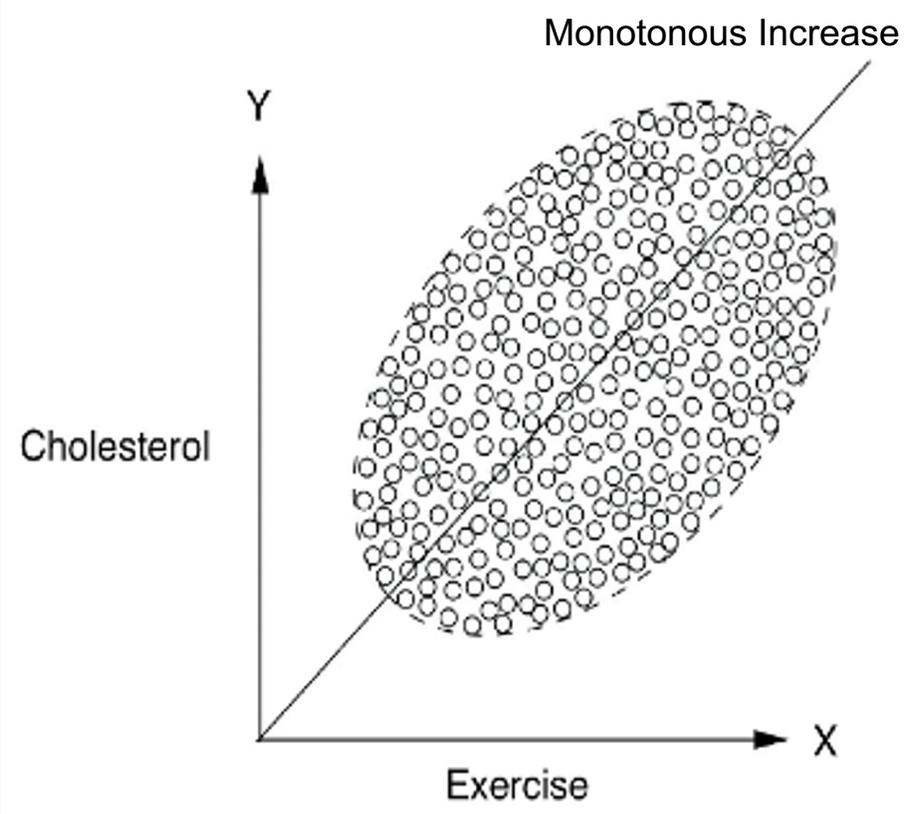
Models learn only correlations $P(Y | X)$

How to solve the issue?

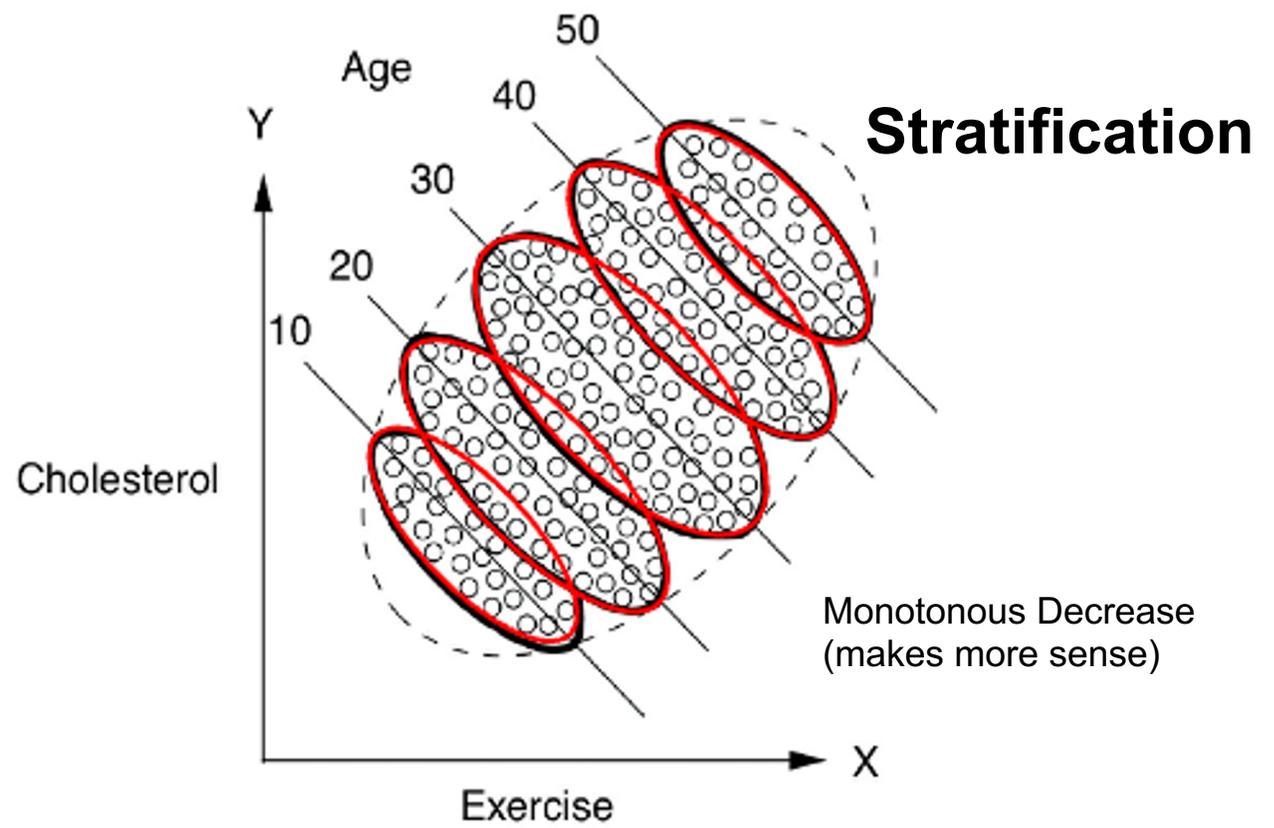
Causal Intervention (from correlation to causality):

$$P(Y | X) \rightarrow P(Y | \text{do}(X))$$

Causal Intervention: $P(Y | X) \rightarrow P(Y|do(X))$

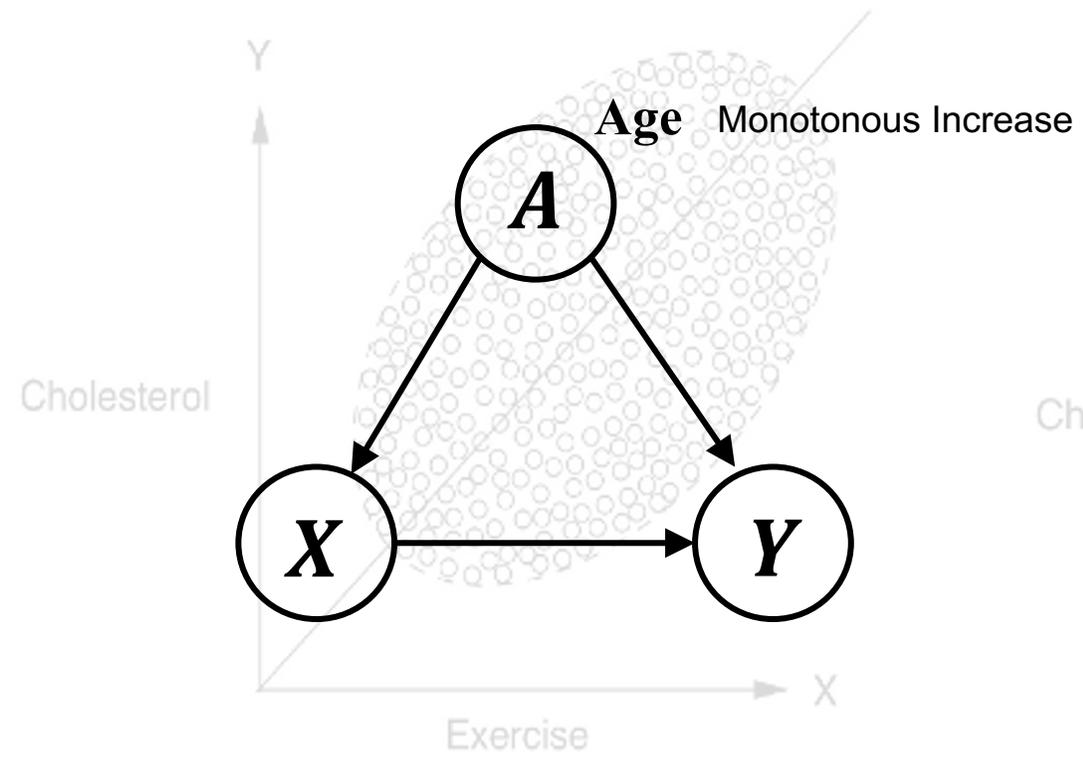


$P(Y | X)$

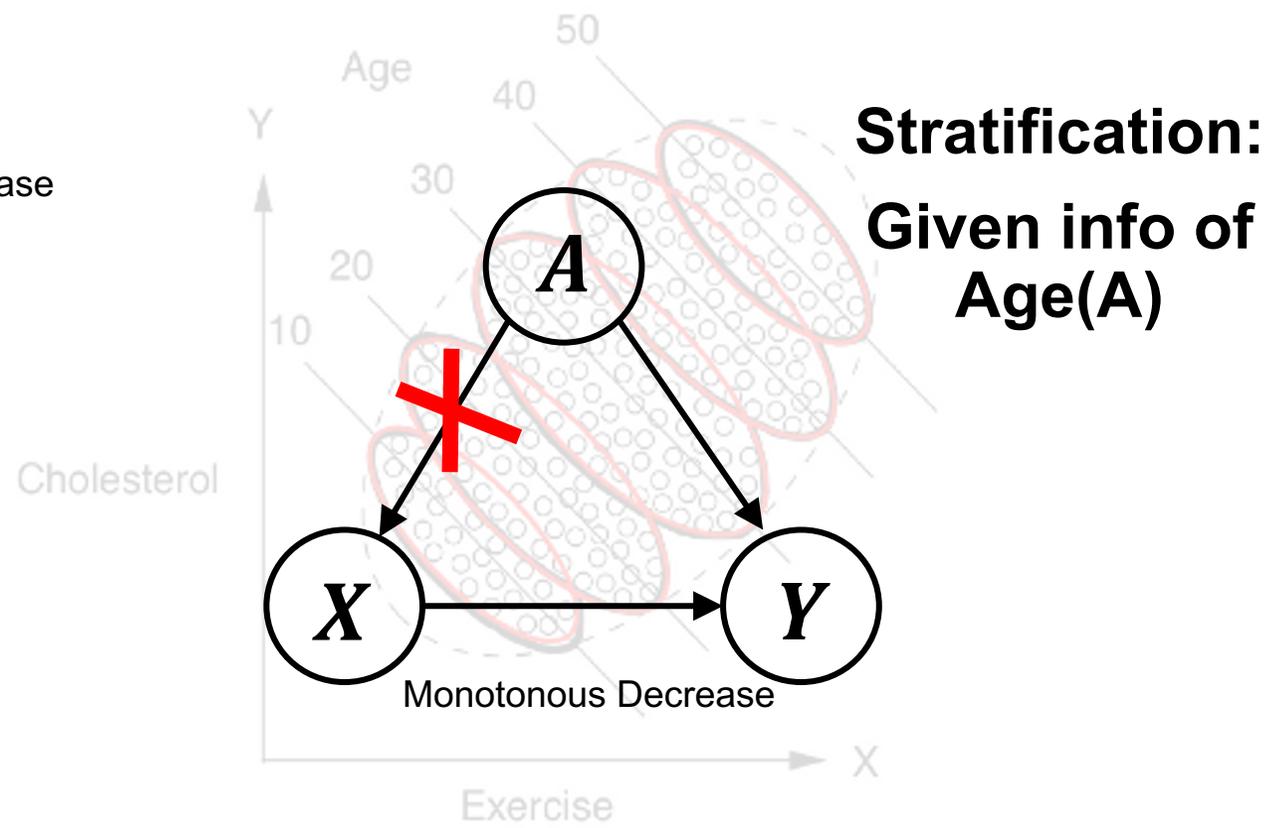


$P(Y|do(X))$ ✓

Causal Intervention: $P(Y | X) \rightarrow P(Y|do(X))$



$P(Y | X)$

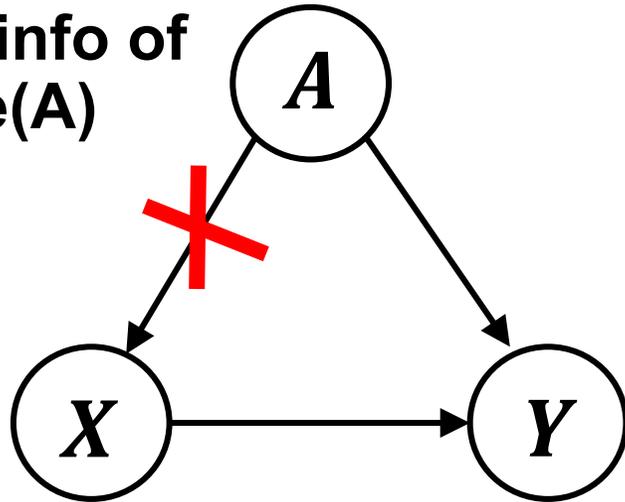


$P(Y|do(X))$

Causal Intervention: $P(Y | X) \rightarrow P(Y|do(X))$

Stratification:

Given info of Age(A)



$$P(Y|do(x)) = \sum_{\boxed{a}} P(Y|X, A = a) \underline{P(A = a)}$$

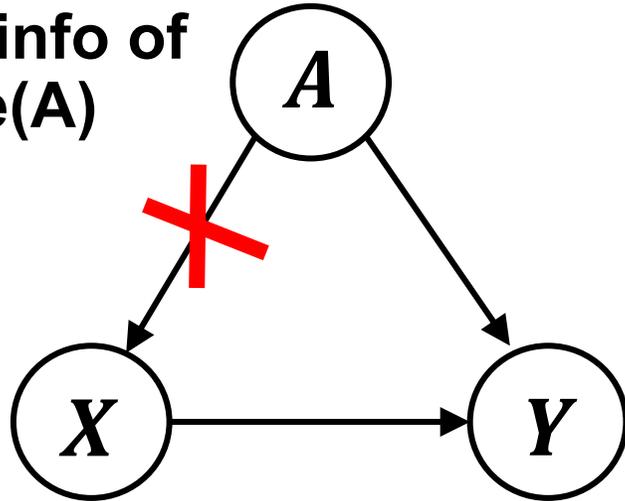
Stratify the Confounder

Weighted by a sample-agnostic prior

Causal Intervention: $P(Y | X) \rightarrow P(Y|do(X))$

Stratification:

Given info of Age(A)



How to implement Stratification in Computer Vision?

$$P(Y|do(x)) = \sum_{\substack{a \\ \text{Stratify the Confounder}}} P(Y|X, A = a) \underbrace{P(A = a)}_{\text{Weighted by a sample-agnostic prior}}$$

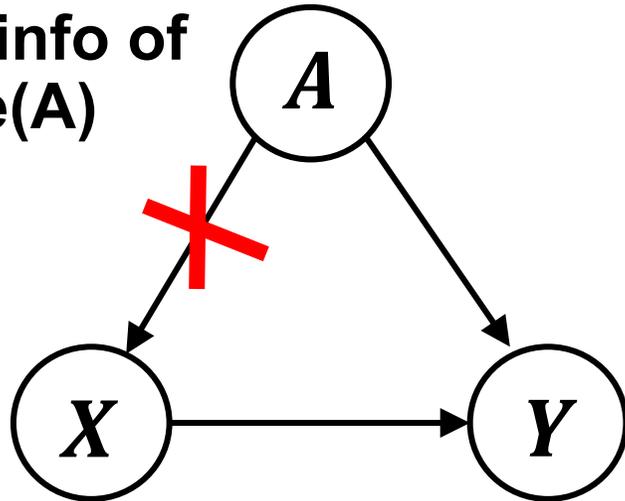
Stratify the Confounder

Weighted by a sample-agnostic prior

Causal Intervention: $P(Y | X) \rightarrow P(Y|do(X))$

Stratification:

Given info of Age(A)



How to implement Stratification in Computer Vision?

-- where we do not have the definition and representation of confounders

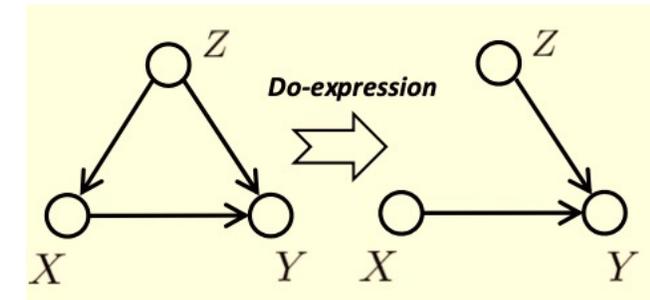
$$P(Y|do(x)) = \sum_{\boxed{a}} P(Y|X, A = a) \underline{P(A = a)}$$

Stratify the Confounder

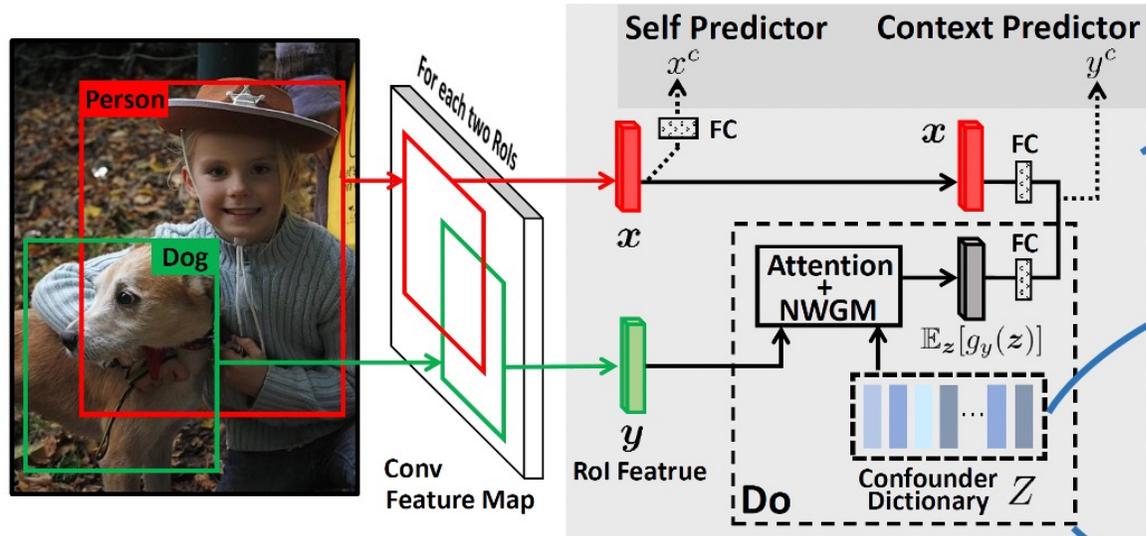
Weighted by a sample-agnostic prior

Causal Intervention: $P(Y | X) \rightarrow P(Y|do(X))$ in Computer Vision

Causal Intervention by debiasing from all Contextual Objects---a Confounder set



X and Y are both image areas, Z is the confounder



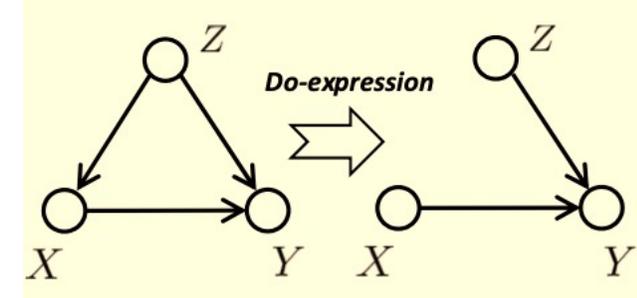
Context Label Prediction Proxy Task

Use an object dictionary to denote confounder Z

$$P(\text{Dog}|do(\text{Person})) = \sum_{z \in Z} P(\text{Dog}|\text{Person}, z)P(z)$$

Causal Intervention: $P(Y | X) \rightarrow P(Y | do(X))$ in Computer Vision

Causal Intervention by debiasing from all Contextual Objects---a Confounder set



Q: Is the girl excited to have a hotdog?



A:Yes



A:Yes



Q: Is his collar buttoned?



A:Yes



A:Yes

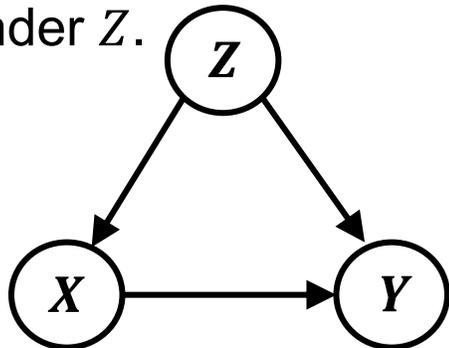


Causal Intervention: $P(Y | X) \rightarrow P(Y|\text{do}(X))$ in Computer Vision?

Causal Intervention: any issues?

Causal Intervention is designed for inference. When used for training:

- Non-positivity: X may never appear with some Z in training set.
e.g., how would you reweight the “black swan” if there is even no “black swan” sample?
- The derivation assumes $P(Y|X)$ to be independent causal mechanisms. That means Z should not influence the $P(Y|X)$ in training. The model should perform **equally** well under different confounder Z .



Causal Intervention: $P(Y | X) \rightarrow P(Y|\text{do}(X))$ in Computer Vision?

Causal Intervention: any issues?

Causal Intervention is designed for inference. When used for training:

- Non-positivity: X may never appear with some Z in training set.
e.g., how would you reweight the “black swan” if there is even no “black swan” sample?
- The derivation assumes $P(Y|X)$ to be independent causal mechanisms. That means Z should not influence the $P(Y|X)$ in training. The model should perform **equally** well under different confounder Z .



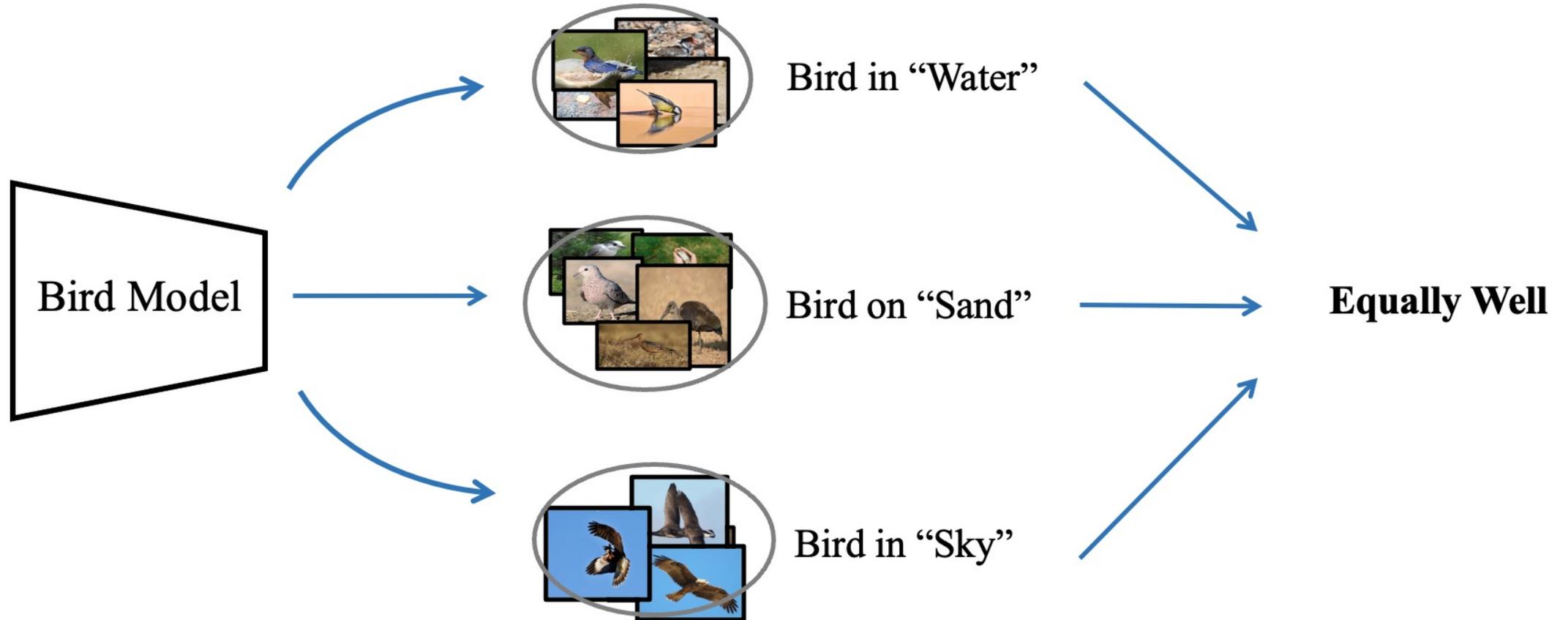
Invariant Learning



Insufficient Data

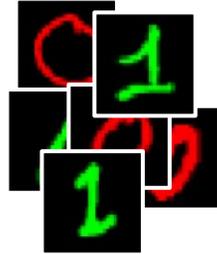
Invariant Learning in Insufficient Data

The model should perform **equally** well under different confounder Z .

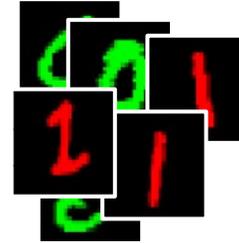


Invariant Learning in Insufficient Data

A toy dataset:



Training Dataset
 $P(Y = 1 | \text{Green}) = 0.8$

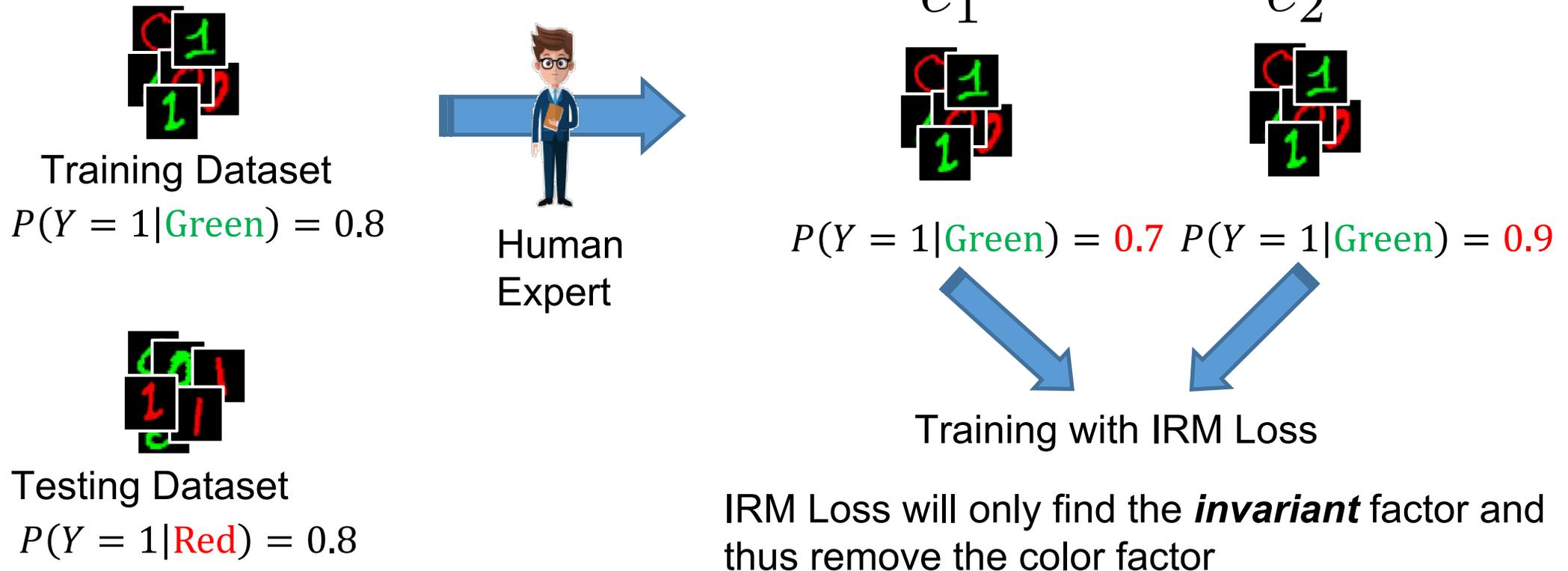


Testing Dataset
 $P(Y = 1 | \text{Red}) = 0.8$

OOD Generalization Problem !

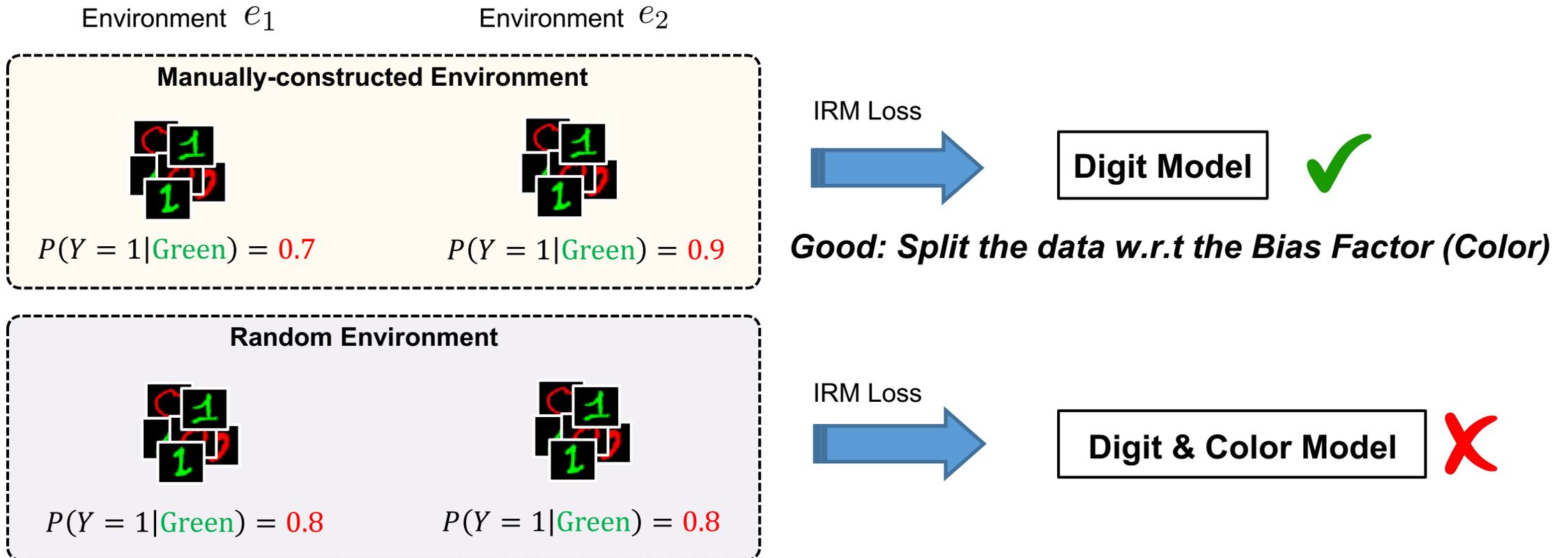
Invariant Learning in Insufficient Data: Invariant Risk Minimization (IRM)

IRM Loss for tackling this problem:



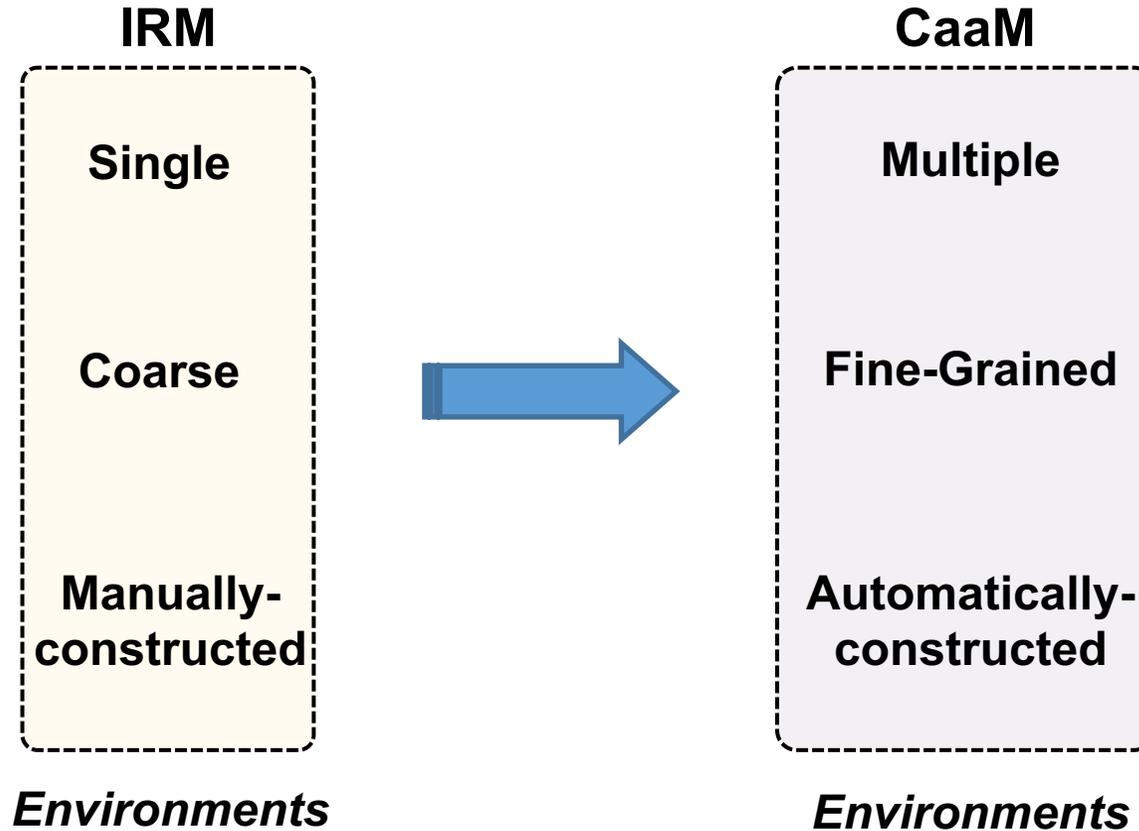
Invariant Learning in Insufficient Data: Invariant Risk Minimization (IRM)

IRM Loss for tackling this problem: by constructing good env



Invariant Learning in Insufficient Data: Causal Attention Module (CaaM)

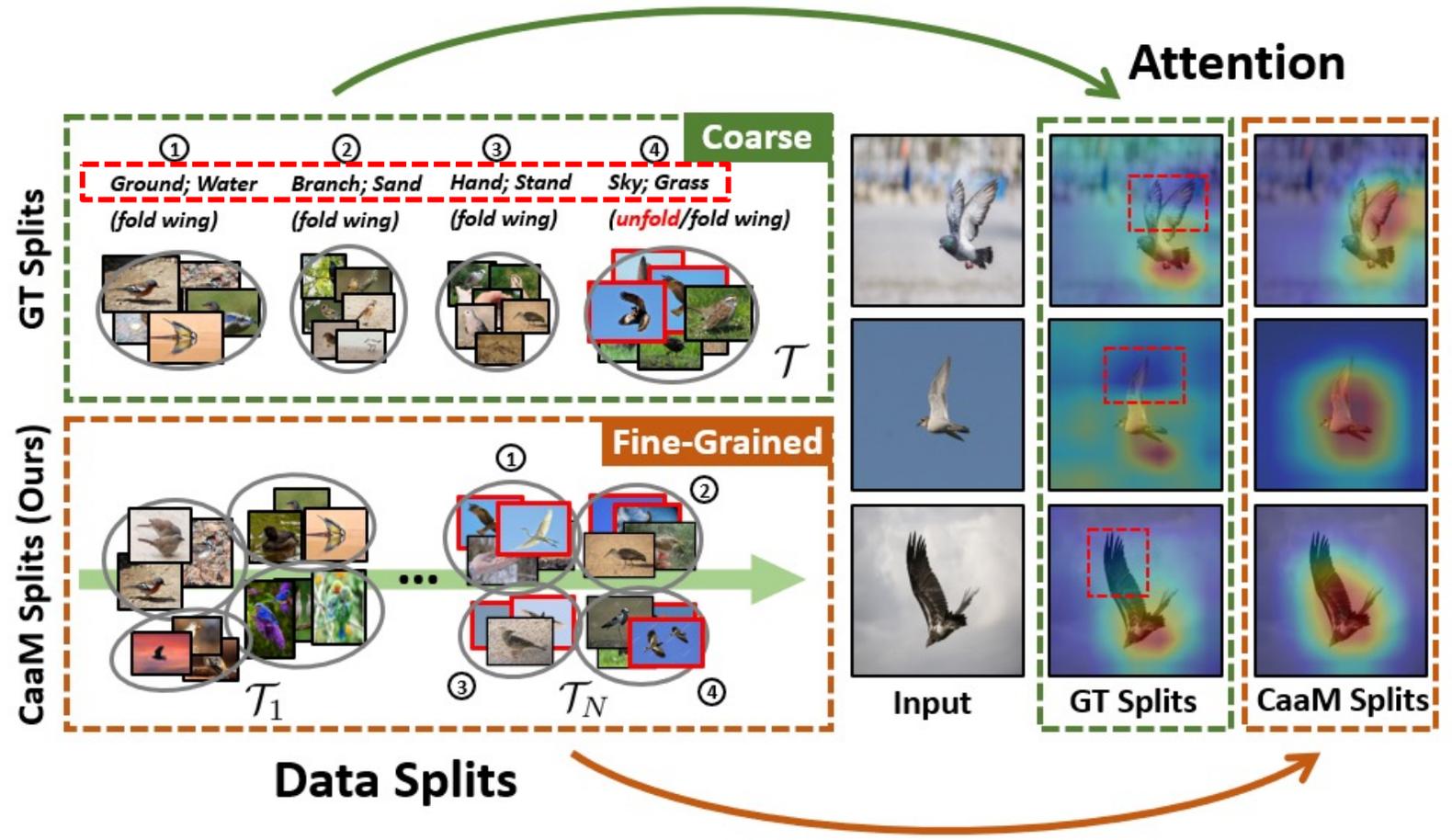
IRM vs. CaaM:



To learn invariances across environments, find a data representation such that the optimal classifier on top of that representation matches for all environments.

Invariant Learning in Insufficient Data: Causal Attention Module (CaaM)

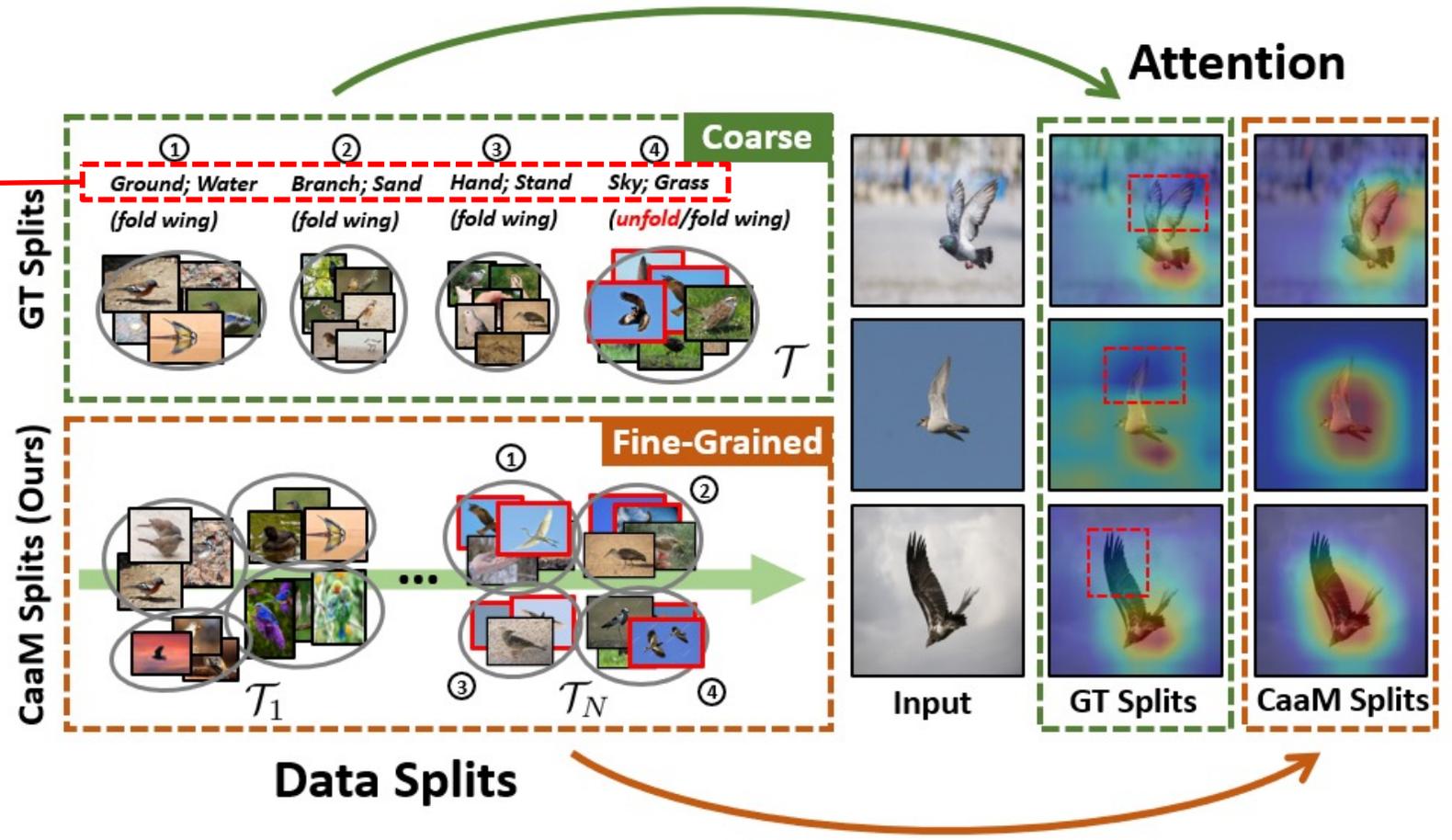
How to understand CaaM:



Invariant Learning in Insufficient Data: Causal Attention Module (CaaM)

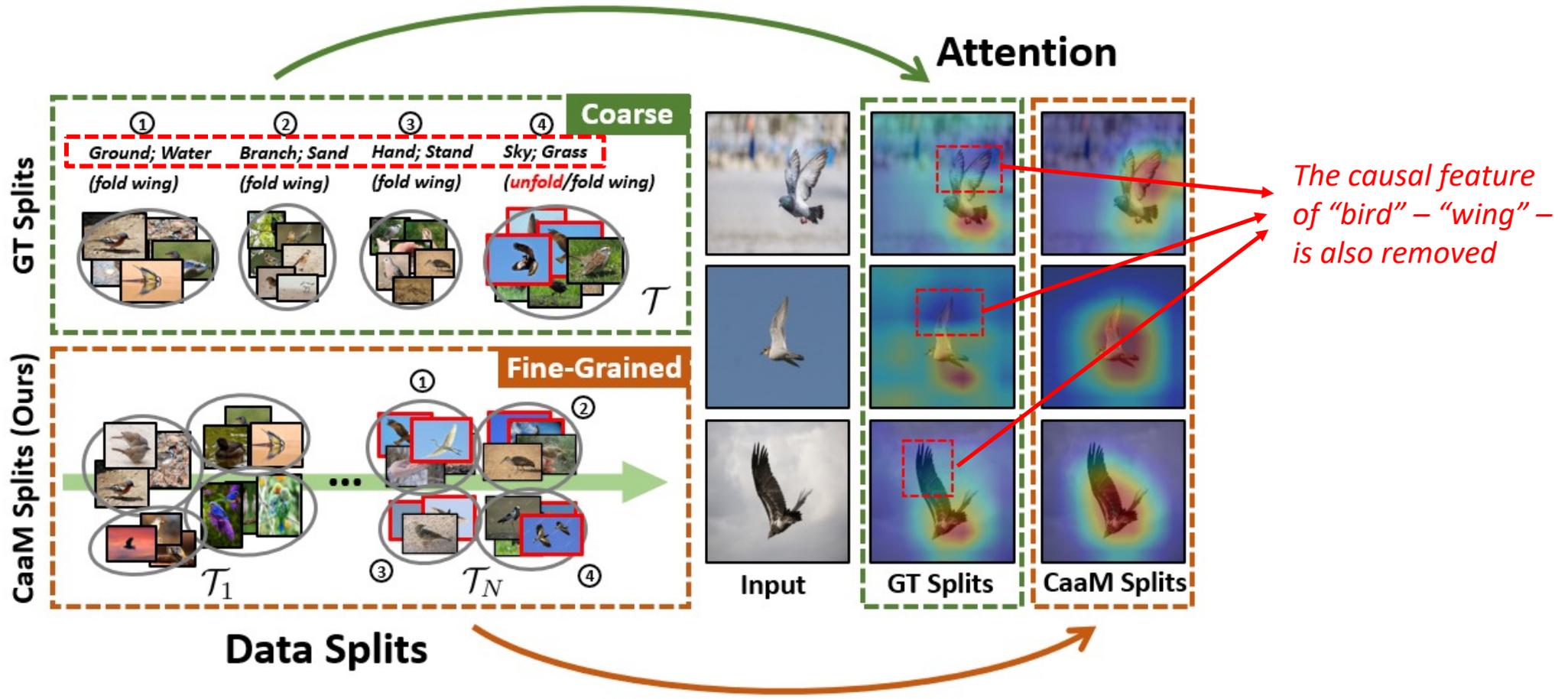
How to understand CaaM:

IRM: Group the images according to the ground truth context (single split)



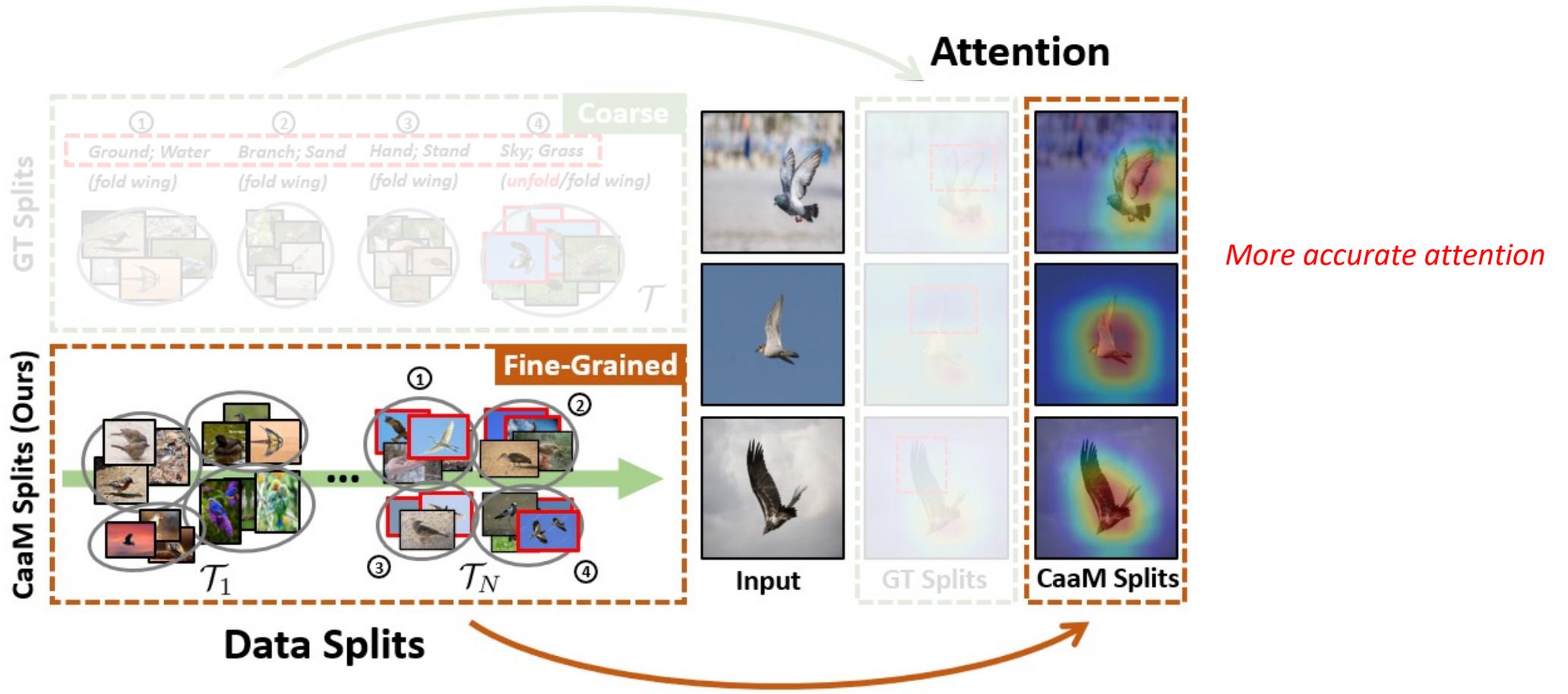
Invariant Learning in Insufficient Data: Causal Attention Module (CaaM)

How to understand CaaM:



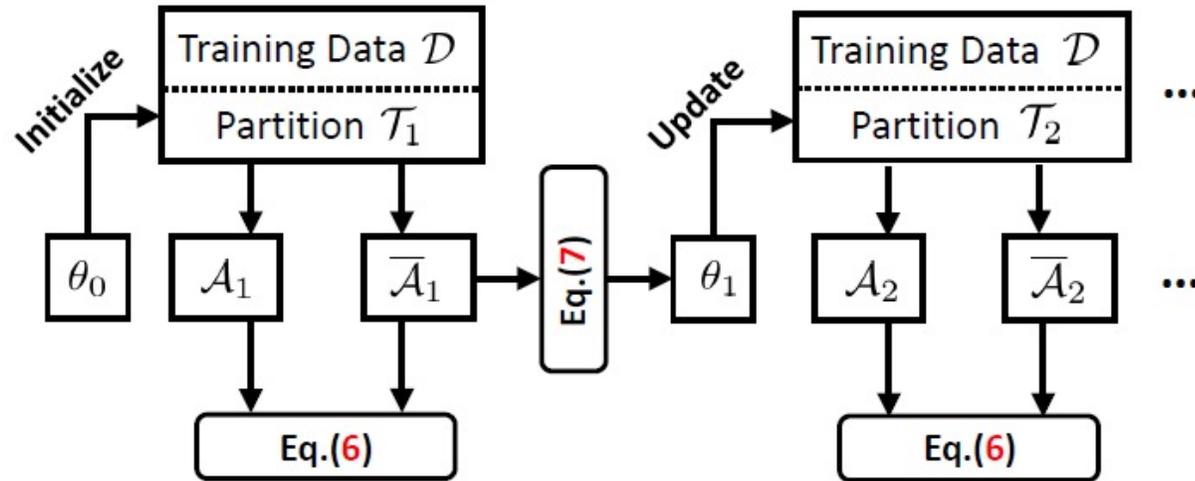
Invariant Learning in Insufficient Data: Causal Attention Module (CaaM)

How to understand CaaM:



Invariant Learning in Insufficient Data: Causal Attention Module (CaaM)

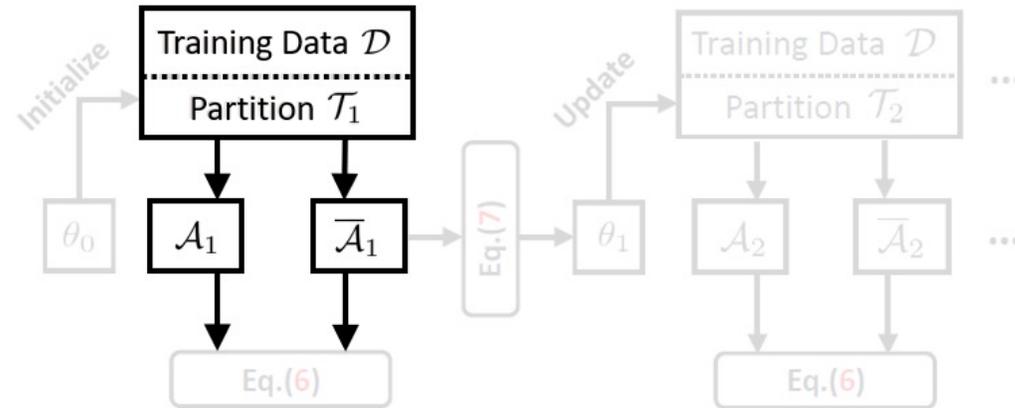
How to implement CaaM:



A : Causal Effect
 \bar{A} : Confounding Effect

Invariant Learning in Insufficient Data: Causal Attention Module (CaaM)

How to implement CaaM:

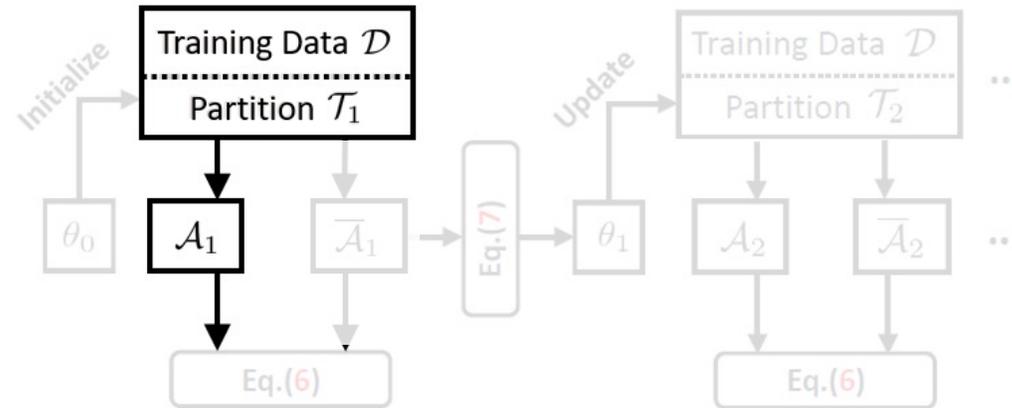


Cross-Entropy Loss (XE): This loss is to ensure that A and \bar{A} combined will capture the total effects

$$XE(f, \tilde{x}, \mathcal{D}) = \mathbb{E}_{(x,y) \in \mathcal{D}} \ell(f(\tilde{x}), y) \quad \tilde{x} = \mathcal{A}(x) \circ \bar{\mathcal{A}}(x)$$

Invariant Learning in Insufficient Data: Causal Attention Module (CaaM)

How to implement CaaM:

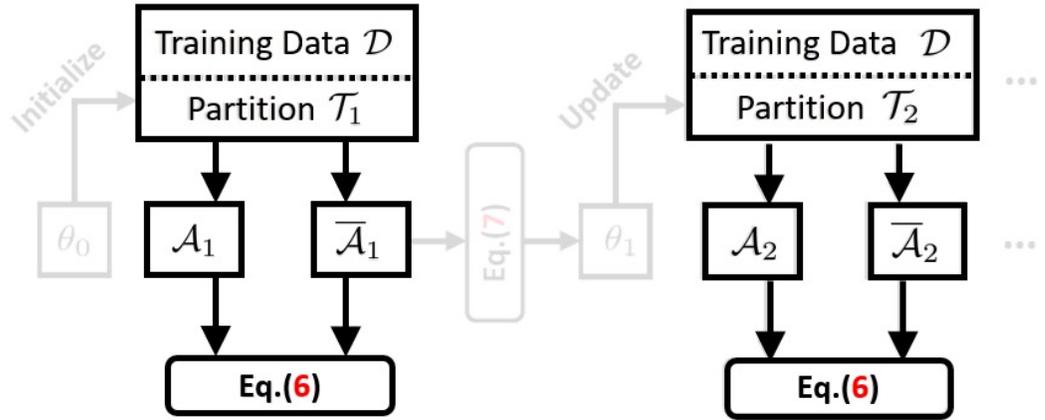


Invariant Loss (IL): This loss is for learning A that is split invariant made by the causal intervention with the data partition T_i .

$$IL(g, \mathcal{A}(x), \mathcal{T}_i) = \sum_{t \in \mathcal{T}_i} XE(g, \mathcal{A}(x), t) + \lambda \underbrace{\|\nabla_{\mathbf{w}=1.0} XE(\mathbf{w}, \mathcal{A}(x), t)\|_2^2}_{\text{Follow the original IRM}}$$

Invariant Learning in Insufficient Data: Causal Attention Module (CaaM)

How to implement CaaM:



Mini-Game: Optimize the Model

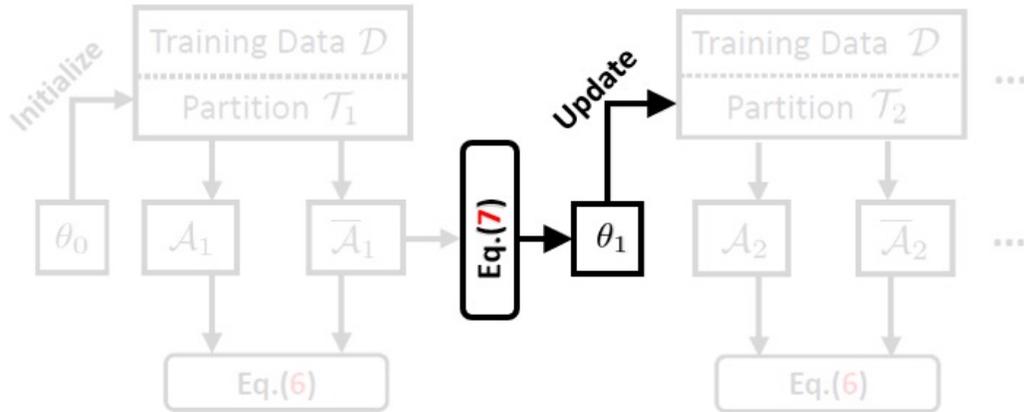
$$\min_{\mathcal{A}, \bar{\mathcal{A}}, f, g, h} \underbrace{\text{XE}(f, \tilde{x}, \mathcal{D}) + \text{IL}(g, \mathcal{A}(x), \mathcal{T}_i)}_{\text{Learn the causal feature } A(x)} + \underbrace{\text{XE}(h, \bar{\mathcal{A}}(x), \mathcal{D})}_{\text{Learn the bias classifier } h \text{ by fitting XE loss on } \bar{A}(x)},$$

Learn the causal feature $A(x)$.

Learn the bias classifier h by fitting XE loss on $\bar{A}(x)$.

Invariant Learning in Insufficient Data: Causal Attention Module (CaaM)

How to implement CaaM:



$$\max_{\theta} \text{IL}(h, \bar{\mathcal{A}}(x), \mathcal{T}_i(\theta))$$

A **good** partition update should capture the bias factor that is currently **NOT** split invariant.

Maxi-Game: Optimize a New Partition

Invariant Learning in Insufficient Data: Causal Attention Module (CaaM)

How to implement CaaM:

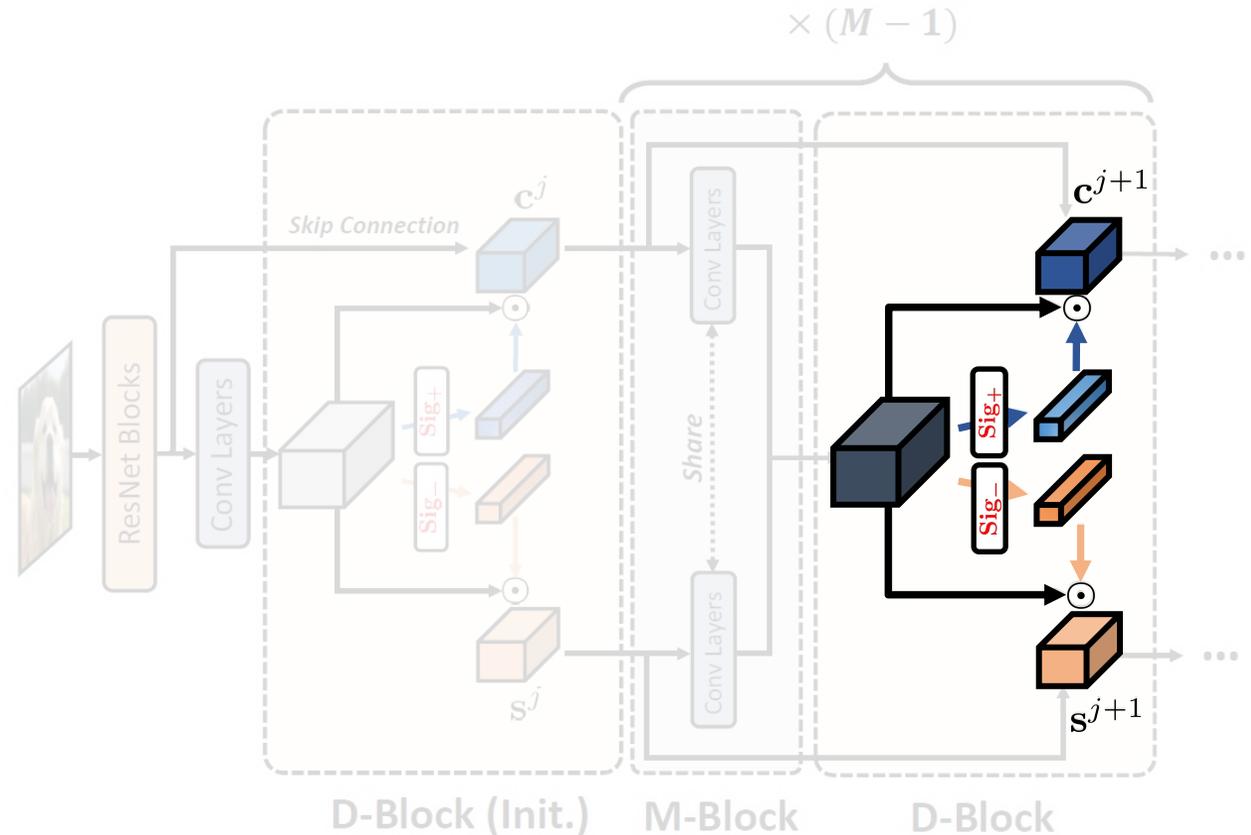
CaaM Attention Calculus for CNN

$$\text{CaaM} : \begin{cases} \mathbf{z} = \text{CBAM}(\mathbf{x}), \\ \mathbf{c} = \text{Sigmoid}(\mathbf{z}) \odot \mathbf{x}, \\ \mathbf{s} = \text{Sigmoid}(-\mathbf{z}) \odot \mathbf{x} \end{cases}$$

$$\text{Sigmoid}(z) = 1 - \text{Sigmoid}(-z)$$



Complementary Attention

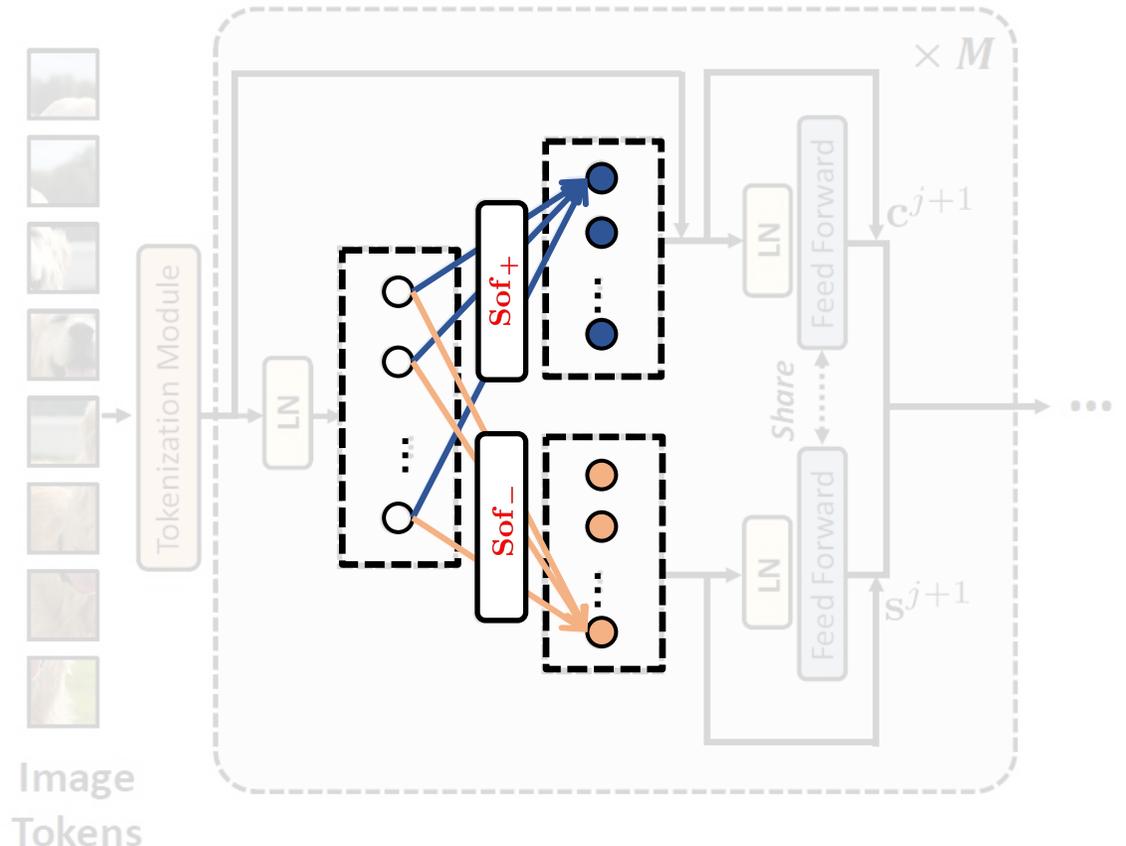


Invariant Learning in Insufficient Data: Causal Attention Module (CaaM)

How to implement CaaM:

CaaM Attention Calculus for ViT

$$\text{CaaM} : \begin{cases} \mathbf{q}, \mathbf{k}, \mathbf{v} = \mathbf{W}_q \mathbf{x}, \mathbf{W}_k \mathbf{x}, \mathbf{W}_v \mathbf{x}, \\ \mathbf{c} = \text{Softmax}\left(\frac{\mathbf{q}\mathbf{k}^T}{\sqrt{d_K}}\right)\mathbf{v}, \\ \mathbf{s} = \text{Softmax}\left(-\frac{\mathbf{q}\mathbf{k}^T}{\sqrt{d_K}}\right)\mathbf{v} \end{cases}$$



Invariant Learning in Insufficient Data: Causal Attention Module (CaaM)

Evaluate CaaM:

Model		CNN-Based					ViT-Based				
		NICO		ImageNet-9 [31]		ImageNet-A [23]	NICO		ImageNet-9 [31]		ImageNet-A [23]
		Val	Test	Biased	Unbiased [5]	Test	Val	Test	Biased	Unbiased [5]	Test
Conv.	ResNet18 [20]	43.77	42.61	95.00	94.40	33.67	–	–	–	–	–
	ResNet18+CBAM [57]	42.15	42.46	94.81	94.09	34.31	–	–	–	–	–
	T2T-ViT7 [63]	–	–	–	–	–	36.23	35.62	88.76	88.35	31.28
	RUBi [8]	43.86	44.37	94.81	94.27	34.13	35.27	34.15	87.95	87.48	29.90
	ReBias [5]	44.92	45.23	95.20	94.89	34.26	35.28	35.74	88.99	88.32	29.33
	Cutout [15]	43.69	43.77	95.24	94.81	34.68	35.31	33.69	87.52	86.47	27.97
	Mixup [67]	44.85	41.46	95.43	94.79	37.71	37.85	34.31	89.72	88.66	30.73
w/H.A. \mathcal{T}	IRM [4]	40.62	41.46	94.13	94.41	33.52	36.46	34.38	89.43	88.87	30.17
	REx [34]	41.00	41.15	94.15	94.28	33.18	36.23	33.46	88.52	87.26	29.18
	Unshuffle [49]	43.15	43.00	94.71	94.33	34.41	37.38	36.00	87.38	86.86	28.61
	CaaM (Ours)	45.46	45.77	95.52	94.96	35.60	38.08	37.54	90.05	89.35	32.01
w/o H.A. \mathcal{T}	IRM [4]	40.54	41.23	94.09	94.32	33.39	33.76	33.77	89.62	88.98	29.25
	REx [34]	40.85	41.52	93.26	93.79	32.84	35.62	34.00	88.68	87.01	28.72
	Unshuffle [49]	41.69	41.61	94.81	94.30	34.04	33.62	32.92	88.38	87.39	28.52
	CaaM (Ours)	46.38	46.62	96.19	95.83	38.55	38.00	37.61	90.33	90.01	32.38

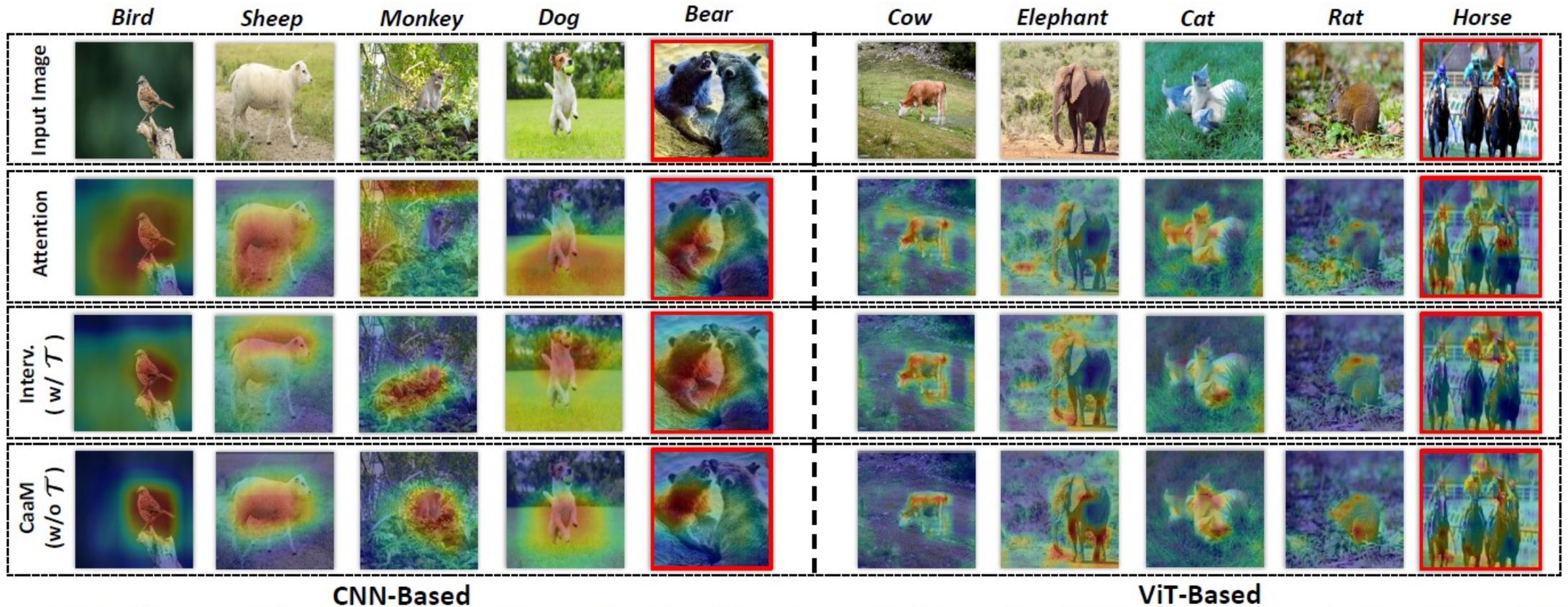
Invariant Learning in Insufficient Data: Causal Attention Module (CaaM)

Evaluate CaaM:

Model		CNN-Based					ViT-Based				
		NICO		ImageNet-9 [31]		ImageNet-A [23]	NICO		ImageNet-9 [31]		ImageNet-A [23]
		Val	Test	Biased	Unbiased [5]	Test	Val	Test	Biased	Unbiased [5]	Test
Conv.	ResNet18 [20]	43.77	42.61	95.00	94.40	33.67	-	-	-	-	-
	ResNet18+CBAM [57]	42.15	42.46	94.81	94.09	34.31	-	-	-	-	-
	T2T-ViT7 [63]	-	-	-	-	-	36.23	35.62	88.76	88.35	31.28
	RUBi [8]	43.86	44.37	94.81	94.27	34.13	35.27	34.15	87.95	87.48	29.90
	ReBias [5]	44.92	45.23	95.20	94.89	34.26	35.28	35.74	88.99	88.32	29.33
	Cutout [15]	43.69	43.77	95.24	94.81	34.68	35.31	33.69	87.52	86.47	27.97
	Mixup [67]	44.85	41.46	95.43	94.79	37.71	37.85	34.31	89.72	88.66	30.73
w/H.A. \mathcal{T}	IRM [4]	40.62	41.46	94.13	94.41	33.52	36.46	34.38	89.43	88.87	30.17
	REx [34]	41.00	41.15	94.15	94.28	33.18	36.23	33.46	88.52	87.26	29.18
	Unshuffle [49]	43.15	43.00	94.71	94.33	34.41	37.38	36.00	87.38	86.86	28.61
	CaaM (Ours)	45.46	45.77	95.52	94.96	35.60	38.08	37.54	90.05	89.35	32.01
w/o H.A. \mathcal{T}	IRM [4]	40.54	41.23	94.09	94.32	33.39	33.76	33.77	89.62	88.98	29.25
	REx [34]	40.85	41.52	93.26	93.79	32.84	35.62	34.00	88.68	87.01	28.72
	Unshuffle [49]	41.69	41.61	94.81	94.30	34.04	33.62	32.92	88.38	87.39	28.52
	CaaM (Ours)	46.38	46.62	96.19	95.83	38.55	38.00	37.61	90.33	90.01	32.38

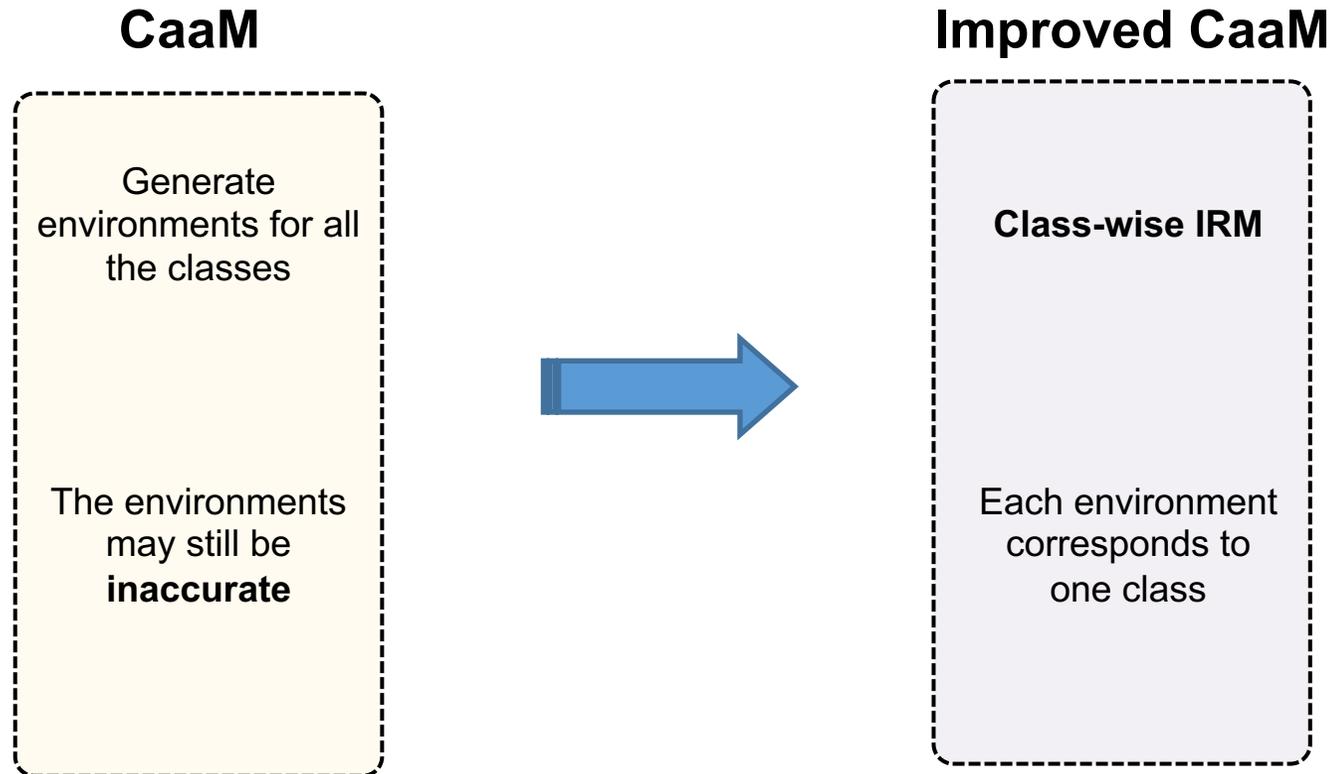
Invariant Learning in Insufficient Data: Causal Attention Module (CaaM)

Evaluate CaaM:



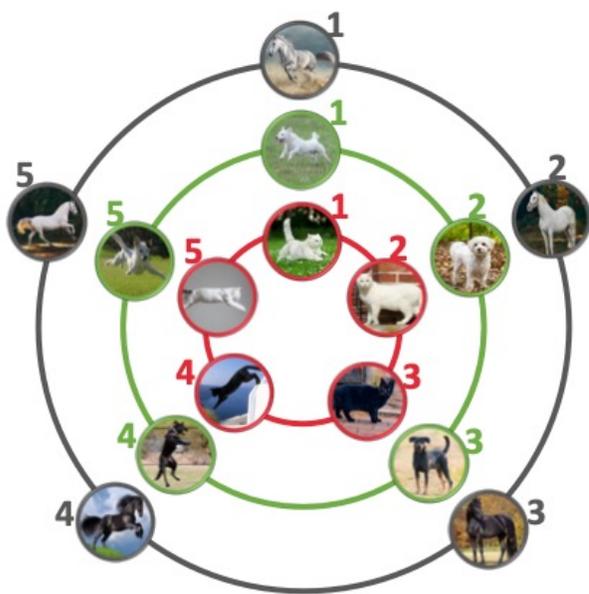
Invariant Learning in Insufficient Data: Causal Attention Module (CaaM)

Improved CaaM (our ongoing work):

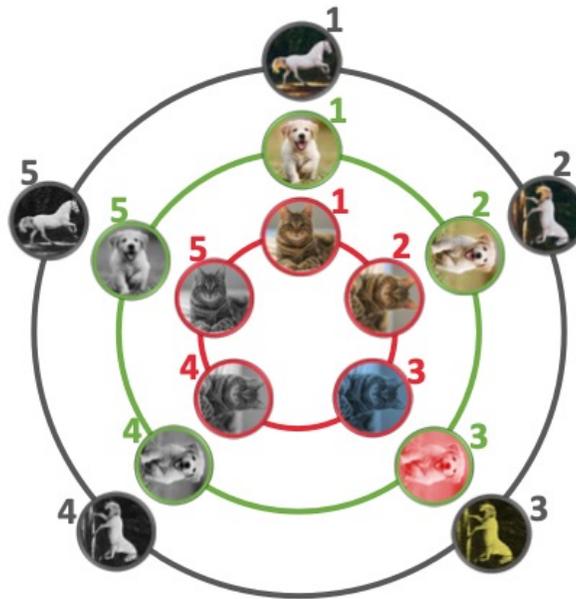


Invariant Learning in Insufficient Data: Representation Learning

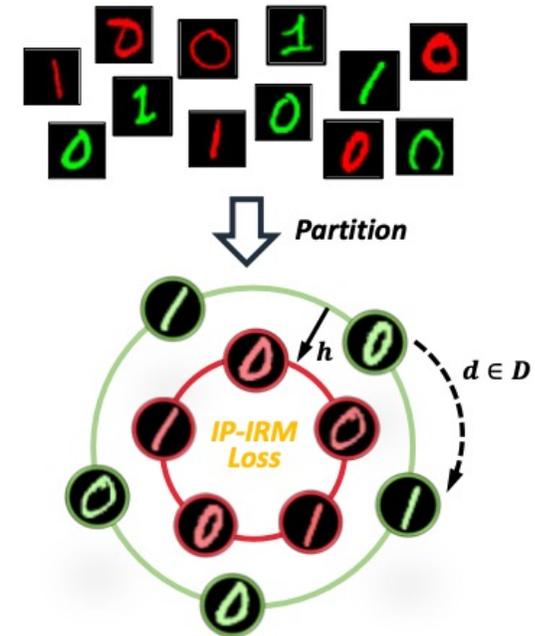
Representation Learning (our published work):



(a) Supervised Learning



(b) Self-Supervised Learning



(c) Our IP-IRM

Invariant Learning in Insufficient Data

**An Interesting Challenge: NICOCHALLENGE
(also in ECCV'22 workshop)**

<https://nicochallenge.com/>

What is NICO CHALLENGE?

The goal of NICO Challenge is to facilitate the OOD (Out-of-Distribution) generalization in visual recognition through promoting the research on the intrinsic learning mechanisms with native invariance and generalization ability. The training data is a mixture of several observed contexts while the test data is composed of unseen contexts. Participants are tasked with developing reliable algorithms across different contexts (domains) to improve the generalization ability of models.

