



中山大學
SUN YAT-SEN UNIVERSITY

JUNE 18-22, 2023
CVPR VANCOUVER, CANADA

Masked Images Are Counterfactual Samples for Robust Fine-tuning

Yao Xiao Ziyi Tang Pengxu Wei* Cong Liu Liang Lin

*Corresponding author

Sun Yat-sen University

Poster: THU-AM-363

Project/Code: <https://github.com/Coxy7/robust-finetuning/>

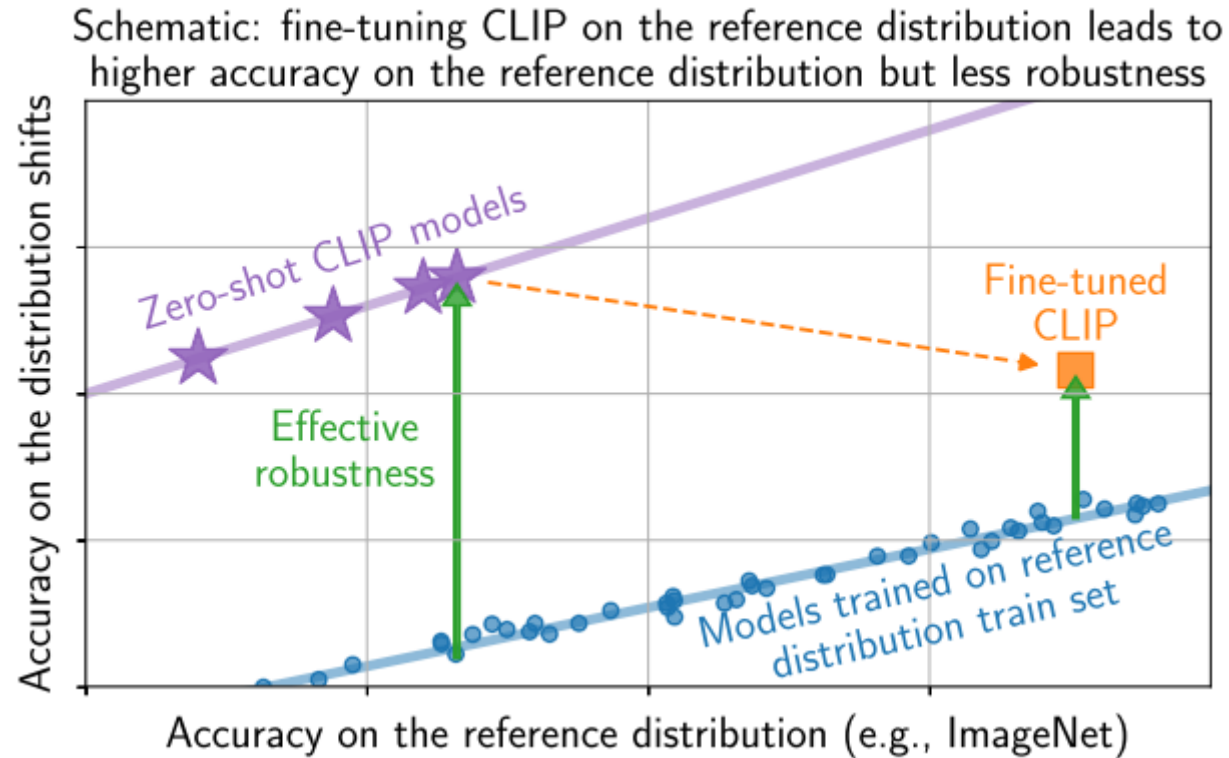
Robustness to distribution shift

		Dataset Examples	ImageNet ResNet101	Zero-Shot CLIP	Δ Score
In-distribution (ID)	ImageNet		76.2	76.2	0%
	ImageNetV2		64.3	70.1	+5.8%
	ImageNet-R		37.7	88.9	+51.2%
Out-of-distribution (OOD)	ObjectNet		32.6	72.3	+39.7%
	ImageNet Sketch		25.2	60.2	+35.0%
	ImageNet-A		2.7	77.1	+74.4%

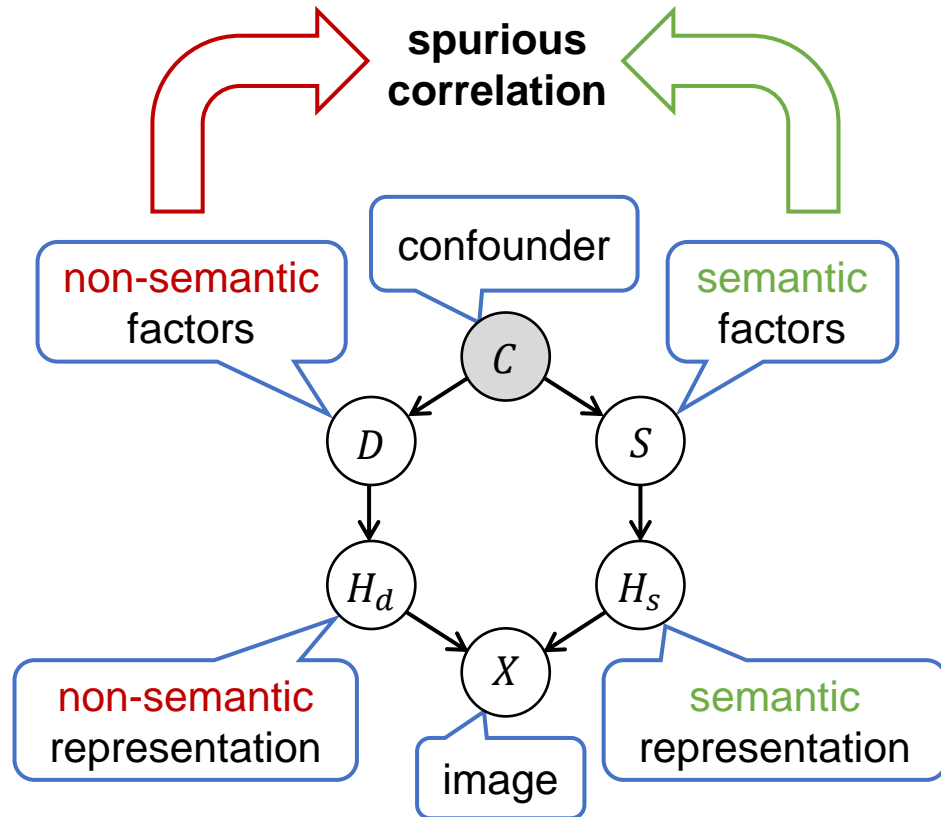
Figure source: Radford, Alec, et al. "Learning transferable visual models from natural language supervision." ICML, 2021.

Fine-tuning can reduce the robustness

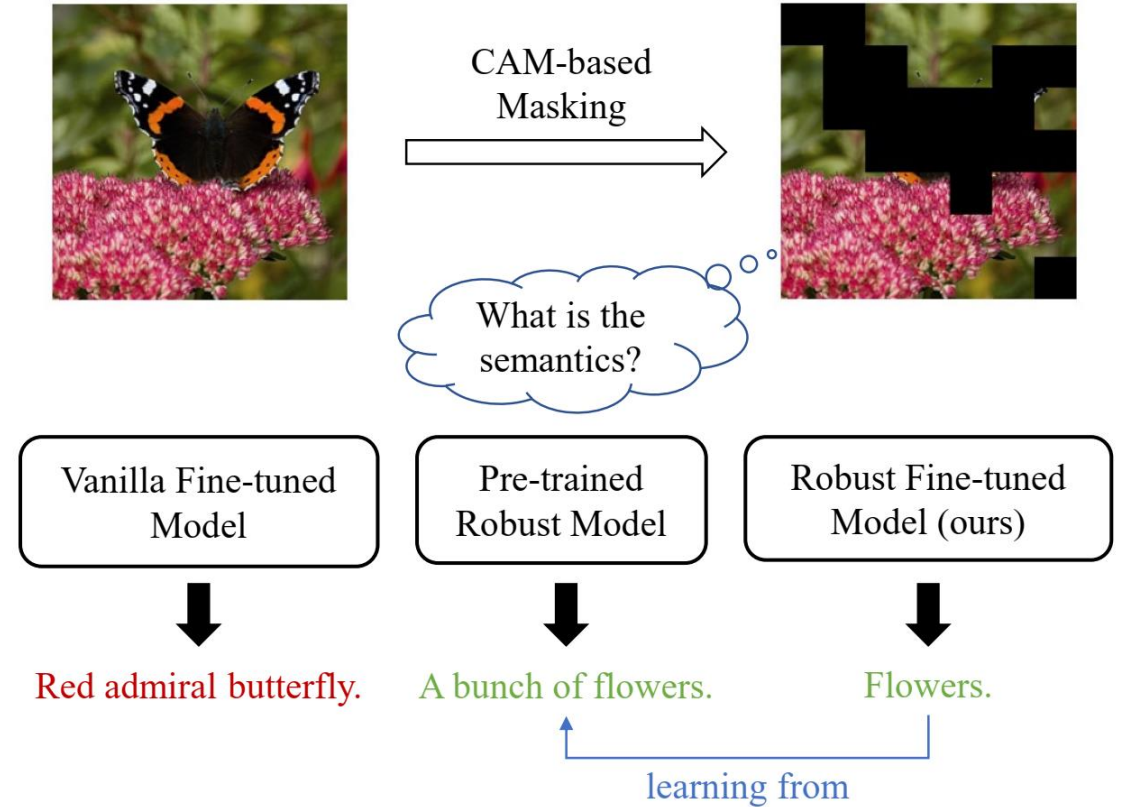
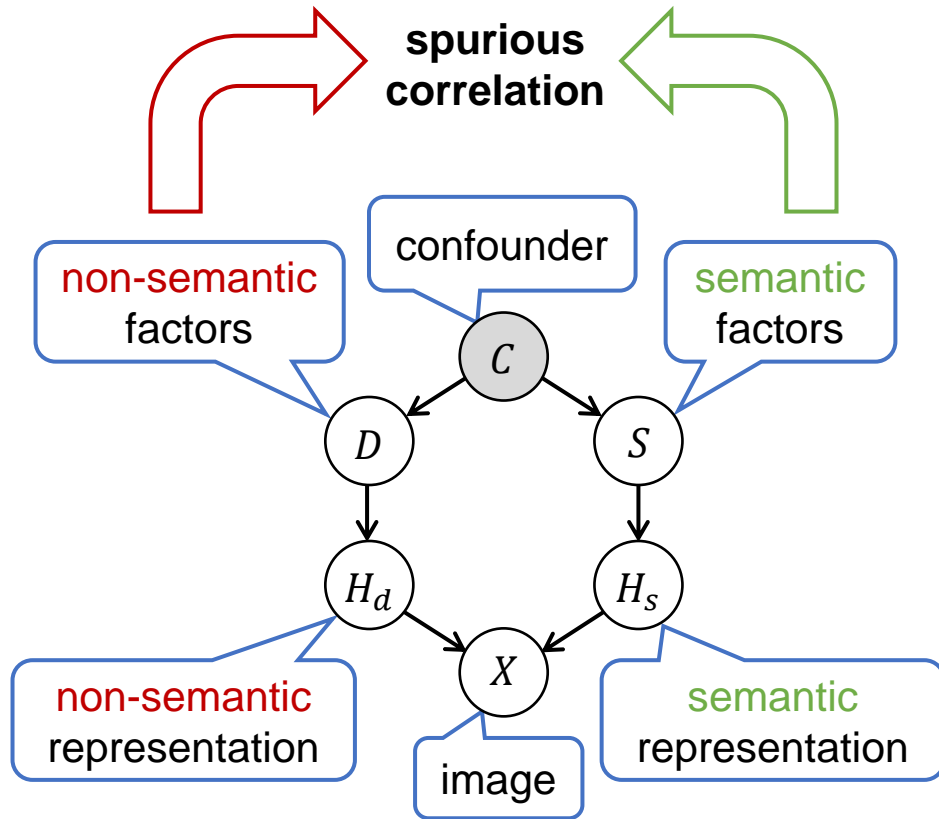
- Trade-off between ID and OOD performance



Overview

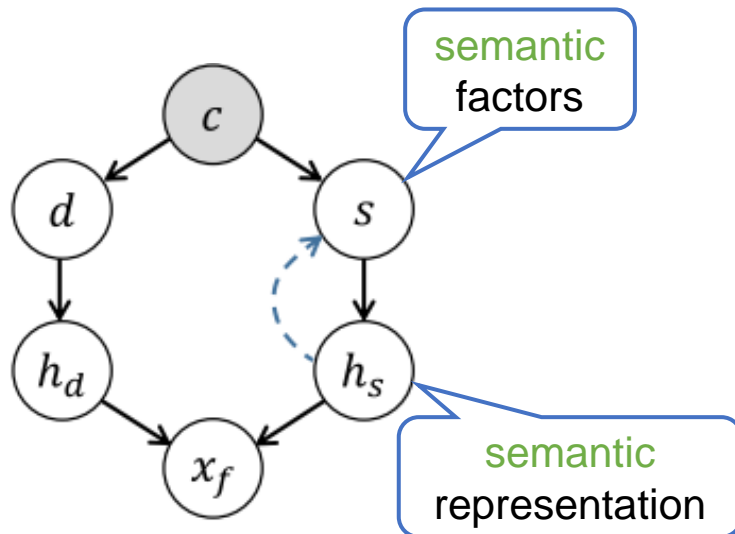


Overview

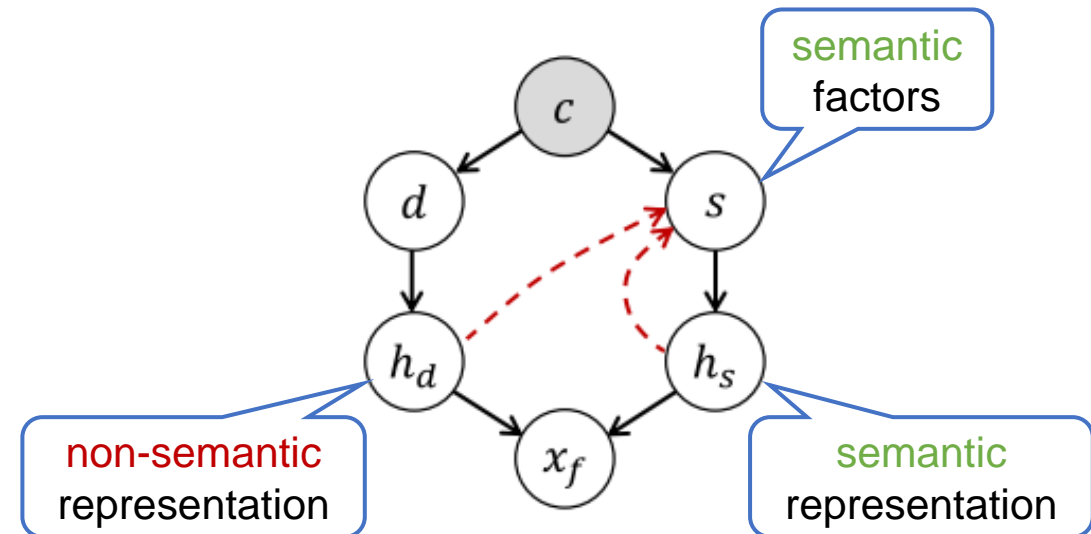


Fine-tuning can reduce the robustness: a causal perspective

- Fine-tuned models tend to rely on both semantic & non-semantic representations (h_s , h_d) for the prediction of image semantics s



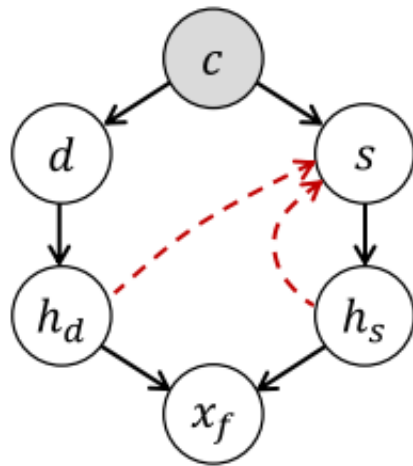
Pre-trained models:
Predict s from h_s



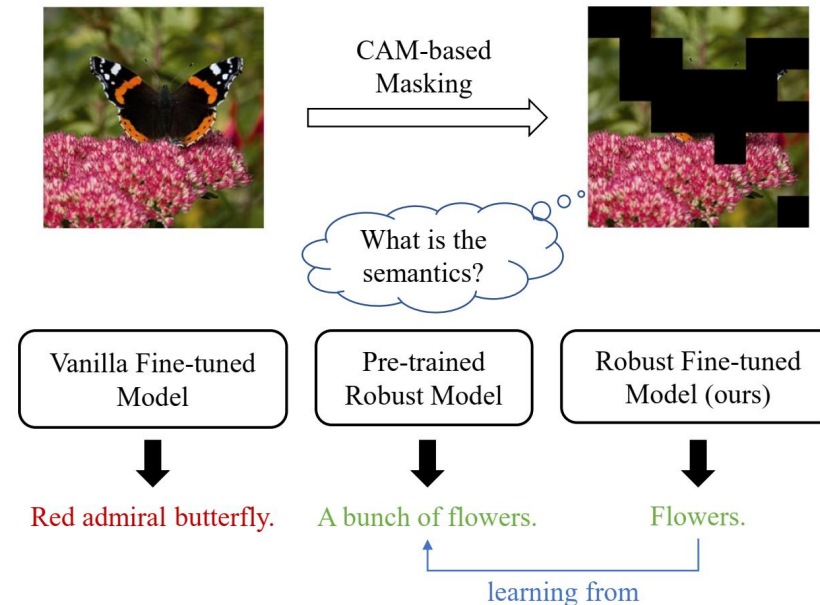
Fine-tuned models:
Predict s from h_s & h_d due to
the *spurious correlation*

Fine-tuning can reduce the robustness: a causal perspective

- Fine-tuned models tend to rely on both semantic & non-semantic representations (h_s, h_d) for the prediction of image semantics s
- The correlation between h_d and s is unstable under distribution shift

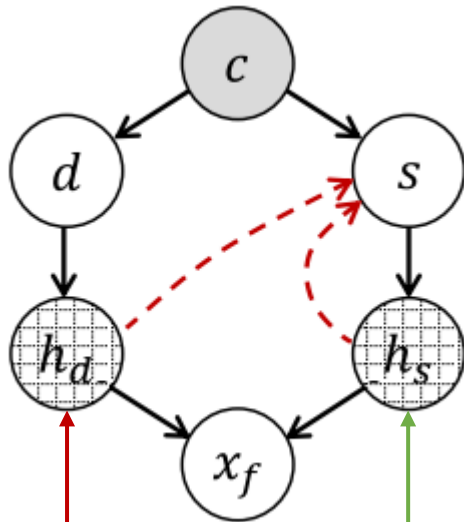


Fine-tuned models:
Predict s from h_s & h_d due to
the spurious correlation

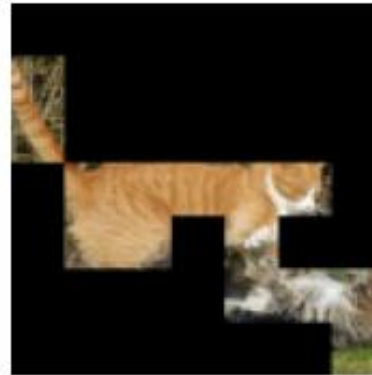


Masked images as counterfactual samples

- Masking based on Class Activation Map (CAM):
 - Context-mask: h_d
 - Object-mask: h_s



context-mask



object-mask



Robust fine-tuning with counterfactual samples

- Distillation with the pre-trained model based on masked images x_{cf}

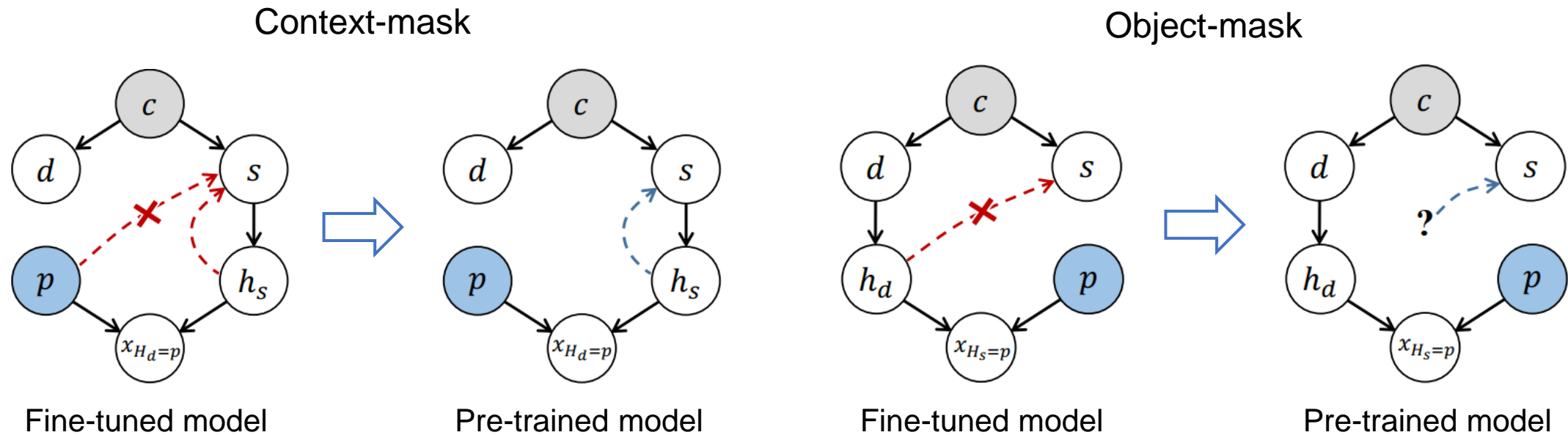
$$\mathcal{L} = \mathcal{L}_{\text{CE}}(g(f(x)), y) + \beta \mathcal{L}_{\text{MSE}}(\hat{f}(x_{cf}), f(x_{cf}))$$

*x : factual sample; x_{cf} : counterfactual sample; y : label
 f, \hat{f} : backbone of fine-tuning/pre-trained model; g : classifier*

Robust fine-tuning with counterfactual samples

- Distillation with the pre-trained model based on masked images x_{cf}

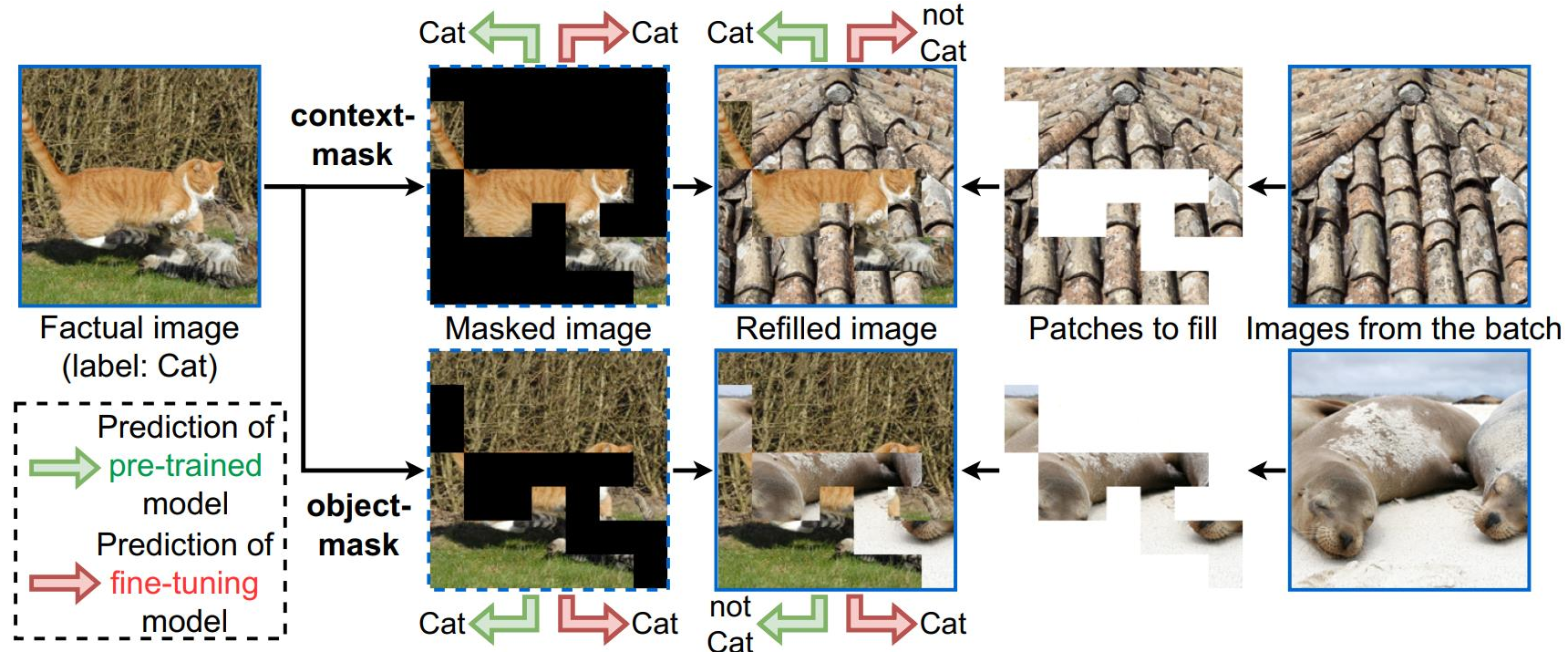
$$\mathcal{L} = \mathcal{L}_{\text{CE}}(g(f(x)), y) + \beta \mathcal{L}_{\text{MSE}}(\hat{f}(x_{cf}), f(x_{cf}))$$



Robust fine-tuning with counterfactual samples

- Refilling

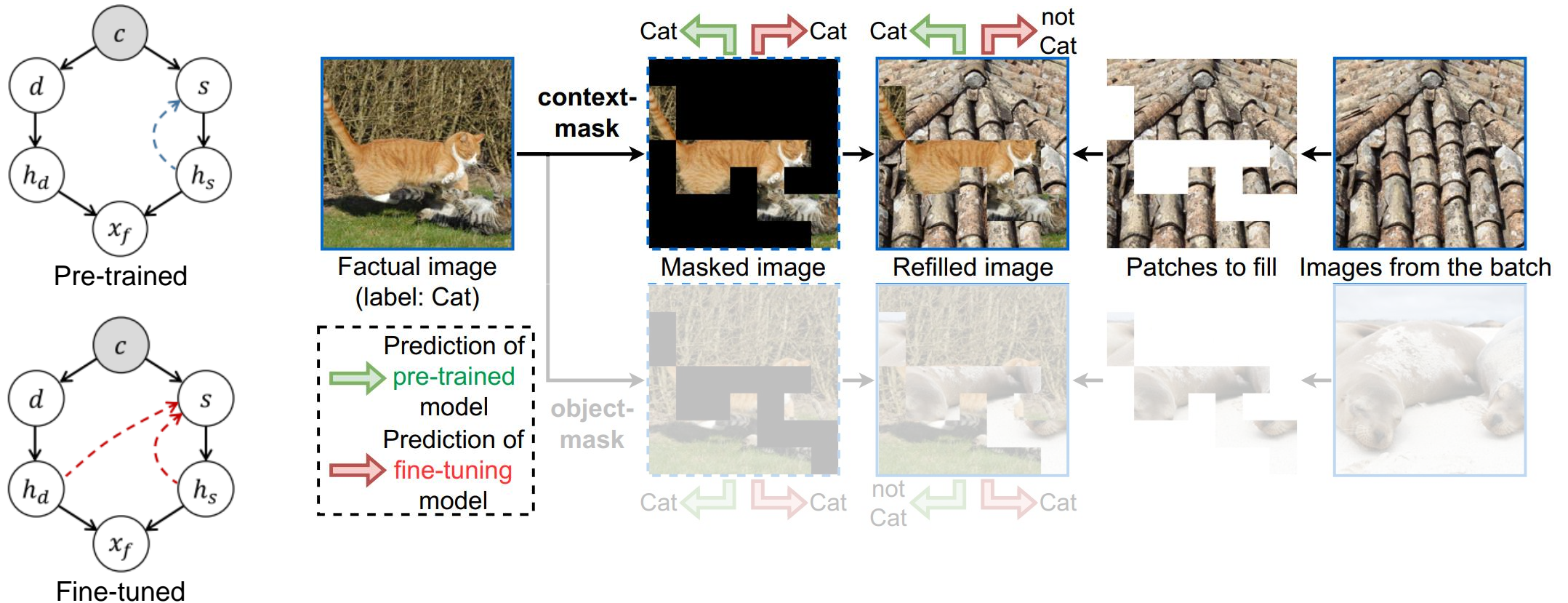
- To enlarge the **disagreement** between fine-tuning and pre-trained model on x_{cf}



Robust fine-tuning with counterfactual samples

- Refilling

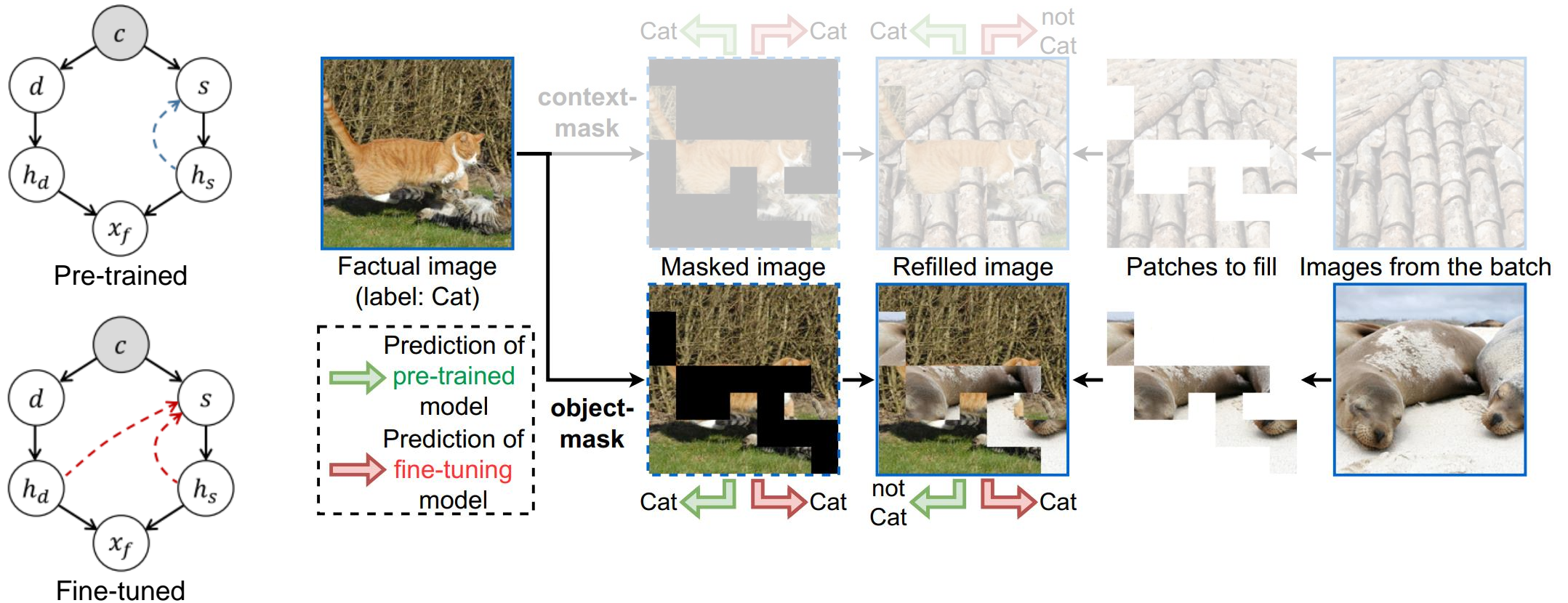
- To enlarge the **disagreement** between fine-tuning and pre-trained model on x_{cf}



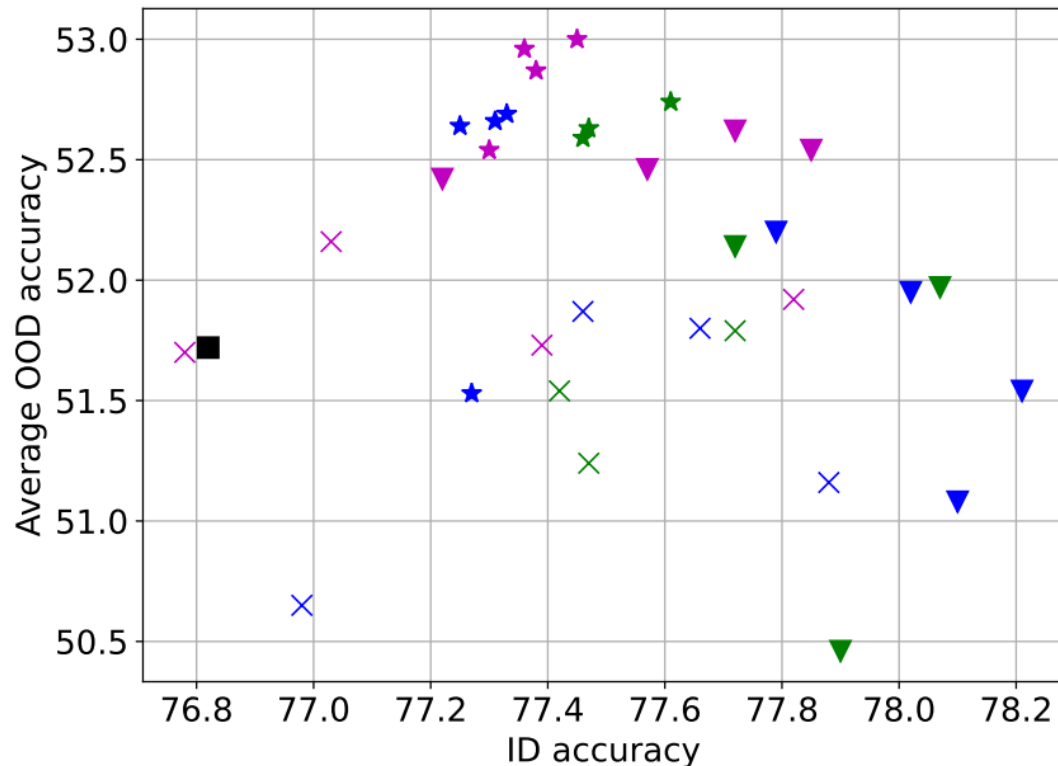
Robust fine-tuning with counterfactual samples

- Refilling

- To enlarge the **disagreement** between fine-tuning and pre-trained model on x_{cf}



Validation of masking and refilling strategies



Hyper-parameters:

Random-mask:

- Masking rate in {0.25, 0.5, 0.75}

Context/Object-mask:

- CAM score threshold in {0.3, 0.4, 0.5, 0.6}

Conclusions:

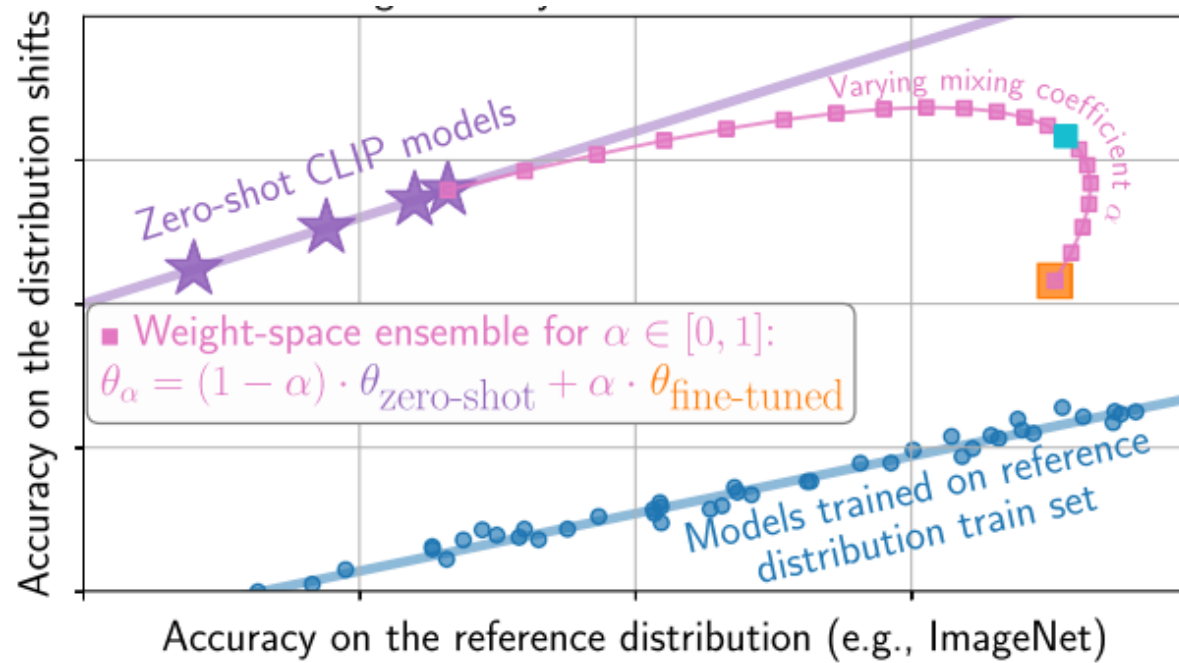
1. Masking > no masking
2. Refilling > no refilling
3. Object-mask > random/context mask

Comparison with existing methods

Model	Method	In-distribution (ID)		Out-of-distribution (OOD)				OOD avg.
		IN	IN-V2	IN-R	IN-Sketch	ObjectNet	IN-A	
CLIP ViT-B/32	Zero-shot [32]	63.4	55.9	69.3	42.3	44.5	31.4	48.7
	Vanilla fine-tuning	75.9	64.7	57.0	39.8	39.5	20.0	44.2
	WiSE-FT [†] [43]	76.6	66.6	70.2	<u>47.1</u>	46.3	31.9	52.4
	Uniform soup [‡] [42]	80.0	68.6	66.6	47.7	46.1	29.2	51.6
	Ours (multi-fill)	<u>77.9</u>	<u>67.7</u>	68.1	46.6	<u>47.5</u>	<u>33.0</u>	<u>52.6</u>
	Ours (single-fill)	<u>77.5</u>	<u>67.1</u>	<u>69.7</u>	46.9	48.0	33.8	53.1
CLIP ViT-B/16	Zero-shot [32]	68.3	61.9	77.6	48.3	54.0	50.1	58.4
	Vanilla fine-tuning	80.7	70.4	64.0	45.1	49.1	35.2	52.8
	LP-FT [23]	81.7	71.6	72.9	48.4	/	49.1	/
	WiSE-FT [43]	81.7	72.8	78.7	53.9	<u>57.3</u>	<u>52.2</u>	<u>63.0</u>
	Ours (multi-fill)	82.5	<u>73.4</u>	76.4	52.7	56.8	52.0	62.3
	Ours (single-fill)	<u>82.4</u>	73.4	<u>78.1</u>	<u>53.4</u>	57.9	53.5	63.3

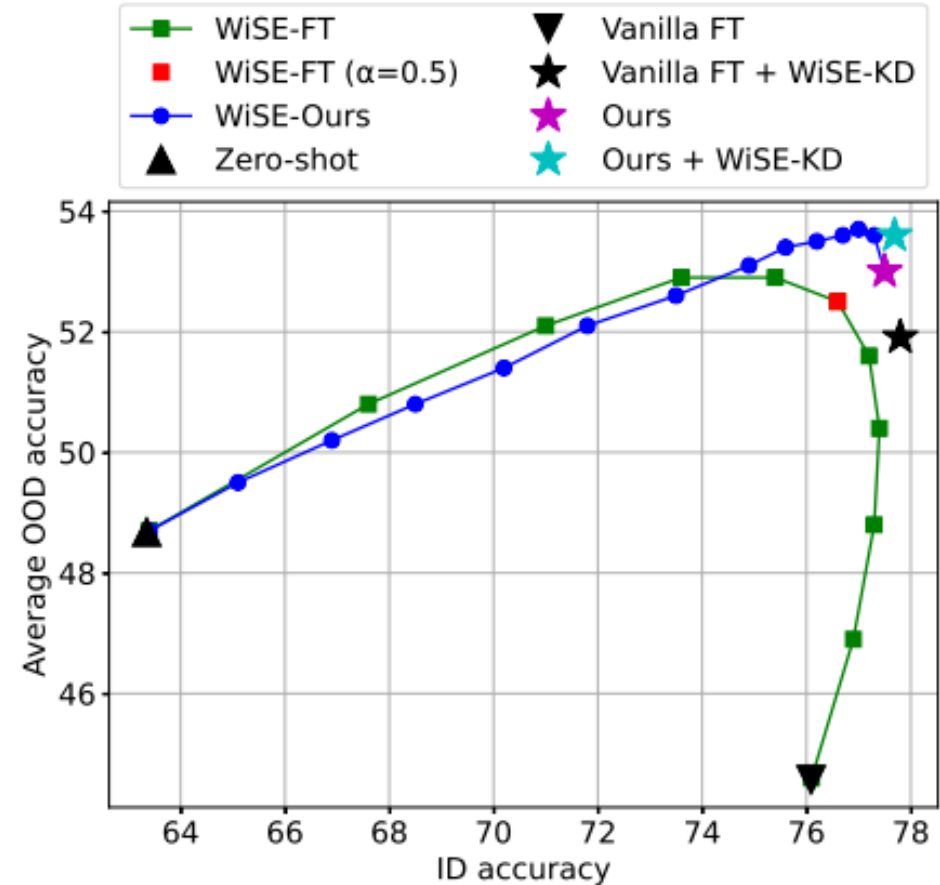
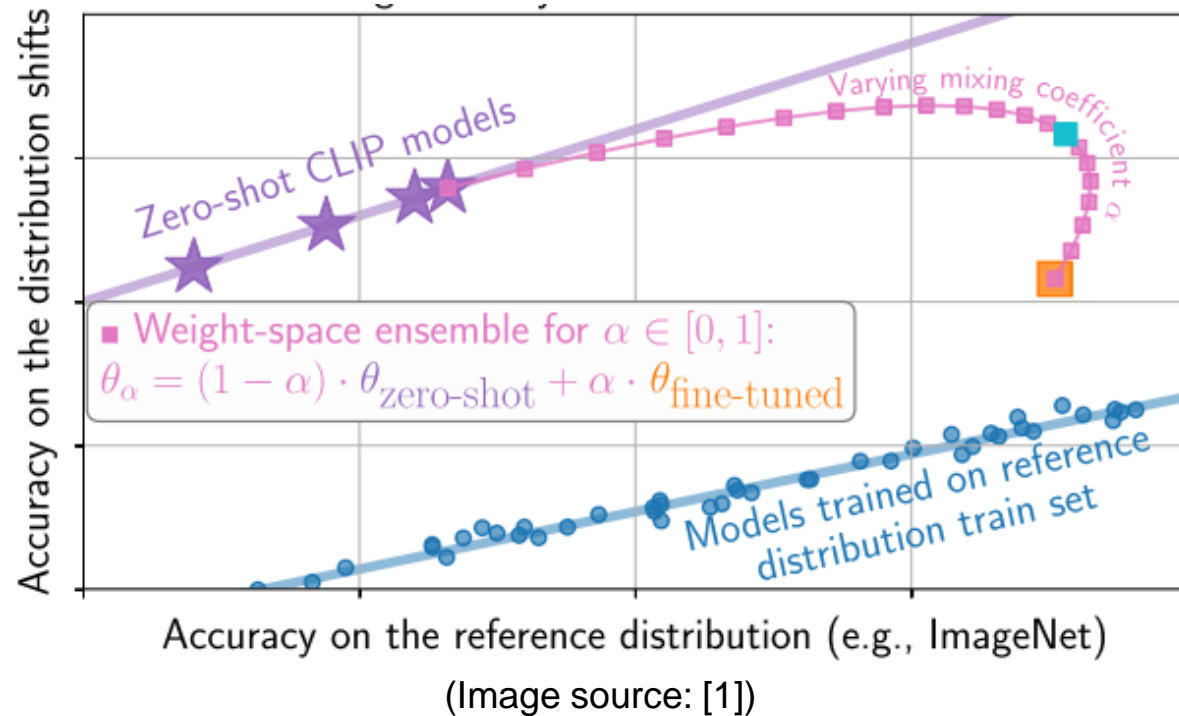
- The proposed method is particularly effective when objects are shown in **unusual contexts** (ObjectNet, ImageNet-A)

Discussion: WiSE-FT [1]



(Image source: [1])

Discussion: WiSE-FT [1]



- Weight-space ensemble of the zero-shot model & our model is less meaningful
- Adding an distillation loss with WiSE-FT teacher slightly improves our model

Conclusion

- The **spurious correlation** between semantic and non-semantic factors in downstream data may account for the **robustness degradation** in fine-tuning.
- **Masked images** can be effective **counterfactual samples** for robust fine-tuning, breaking the spurious correlation.
- **Weight-space constraints** may be sufficient but **not necessary** for maintaining the robustness of the pre-trained model.

Thanks!

Paper: <https://arxiv.org/abs/2303.03052>

Project/Code: <https://github.com/Coxy7/robust-finetuning/>