HOW TO PROBE: SIMPLE YET EFFECTIVE TECHNIQUES FOR IMPROVING POST-HOC EXPLANATIONS





¹ Max Planck Institute for Informatics, Saarland Informatics Campus

²Institute of Science and Technology Austria

²Institute of Science and Technology Austria

Siddhartha Gairola^{1,2}, Moritz Böhle¹, Francesco Locatello^{1,2}, Bernt Schiele¹

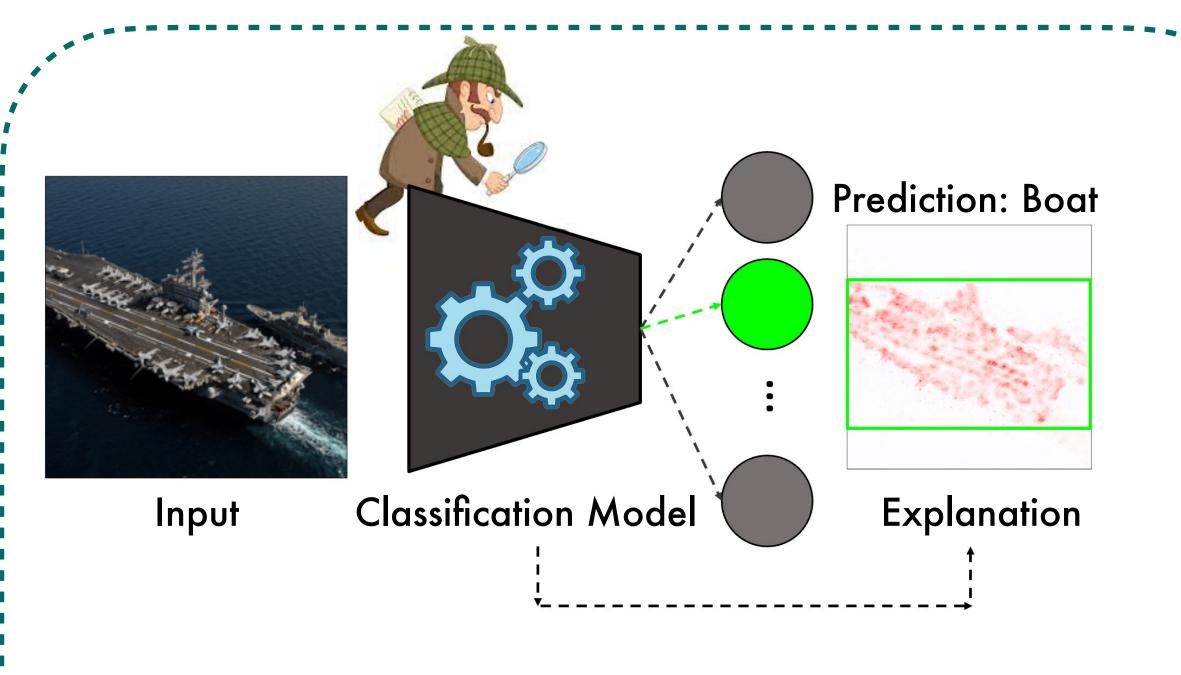












Post-hoc attribution methods are a popular tool for "explaining" DNNs

Key Question

"To what extent does the training objective influence the explanation quality of pre-trained models?"

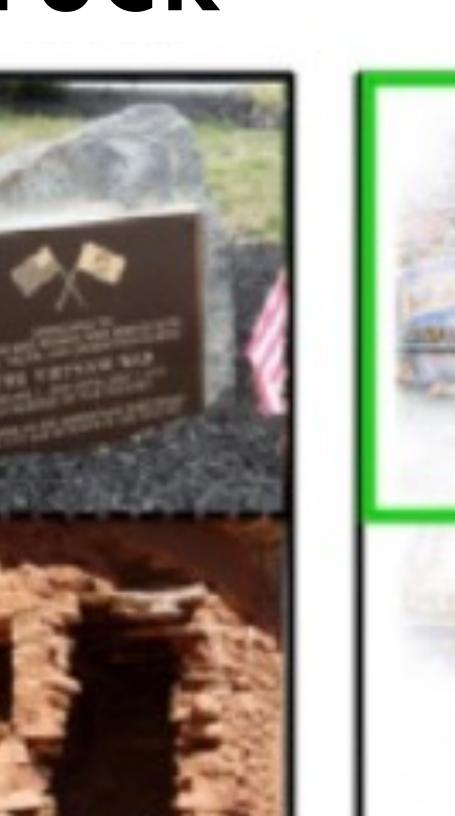
Why This Matters?

Post-hoc attribution methods are popular but training-agnostic.

Diversity in pre-training requires revisiting that assumption.

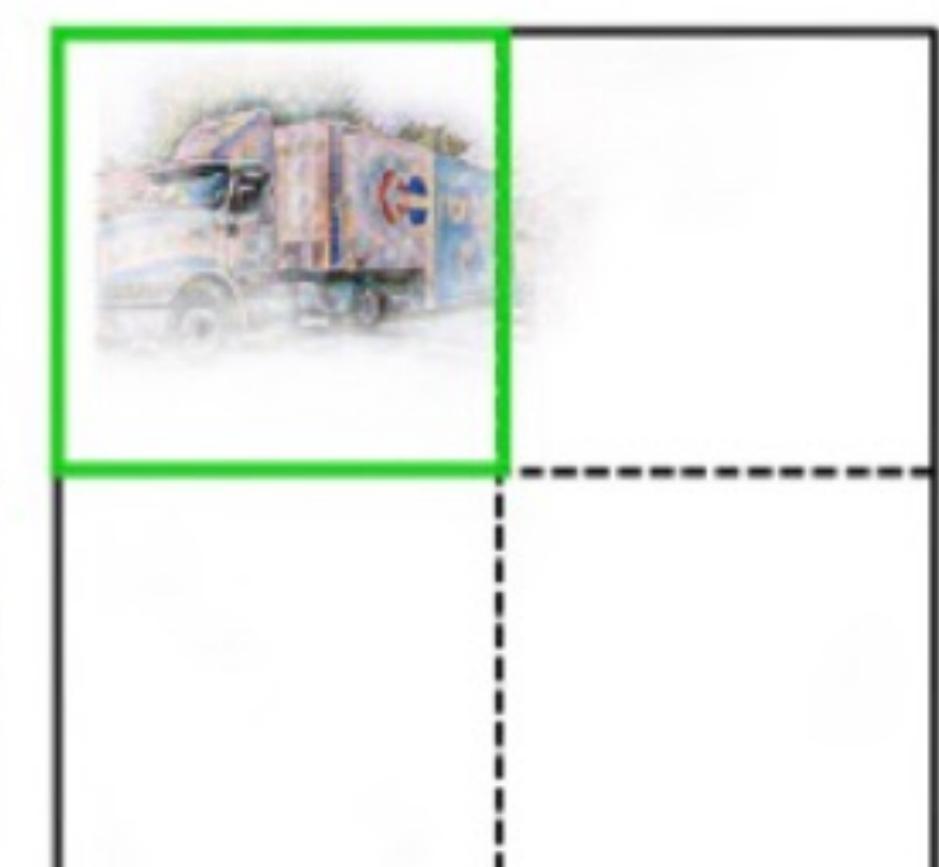
Interpretability Isn't Plug-and-Play

Trailer Truck





CE Loss

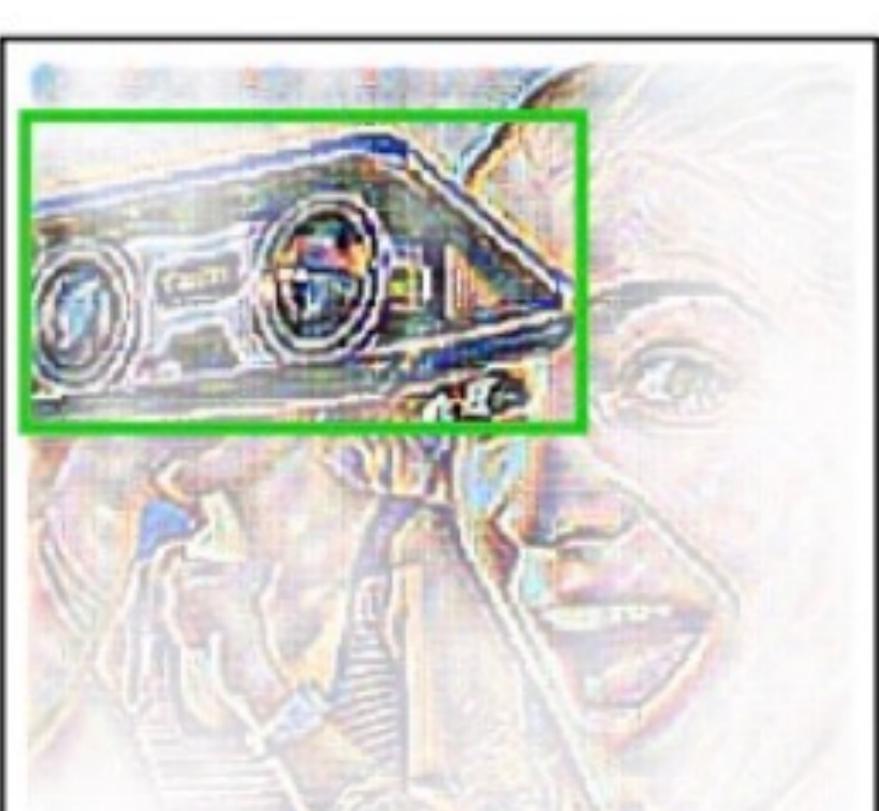


BCE Loss

BCE probes localize better than CE probes.

Linear Probe

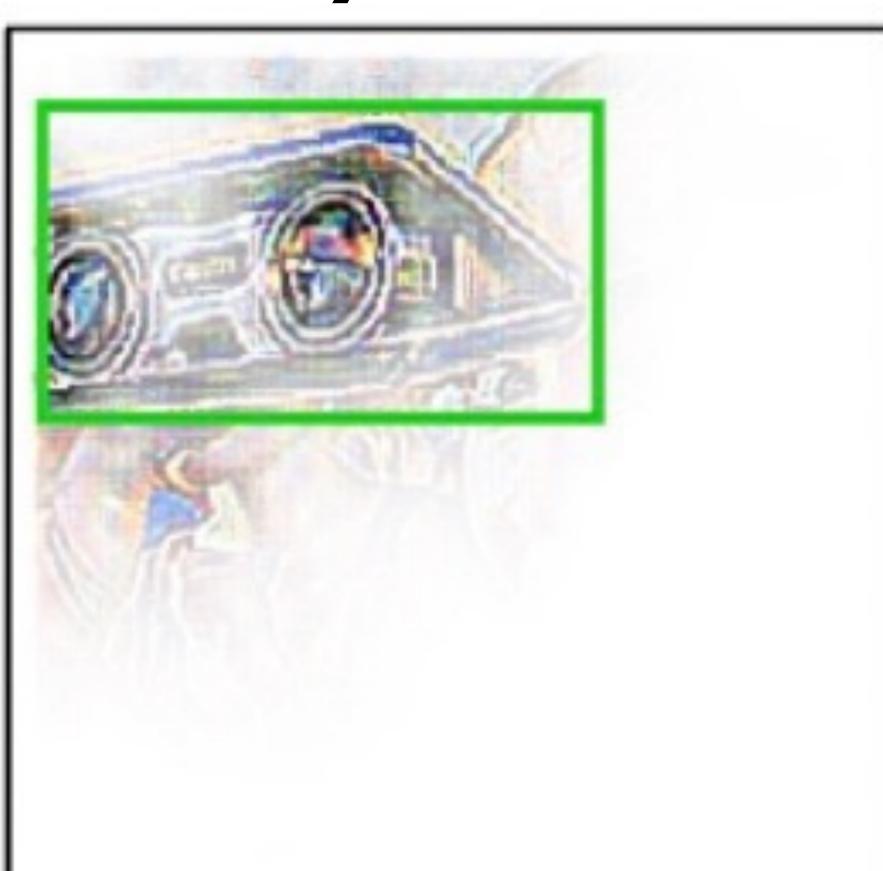
Binoculars



3-layer MLP

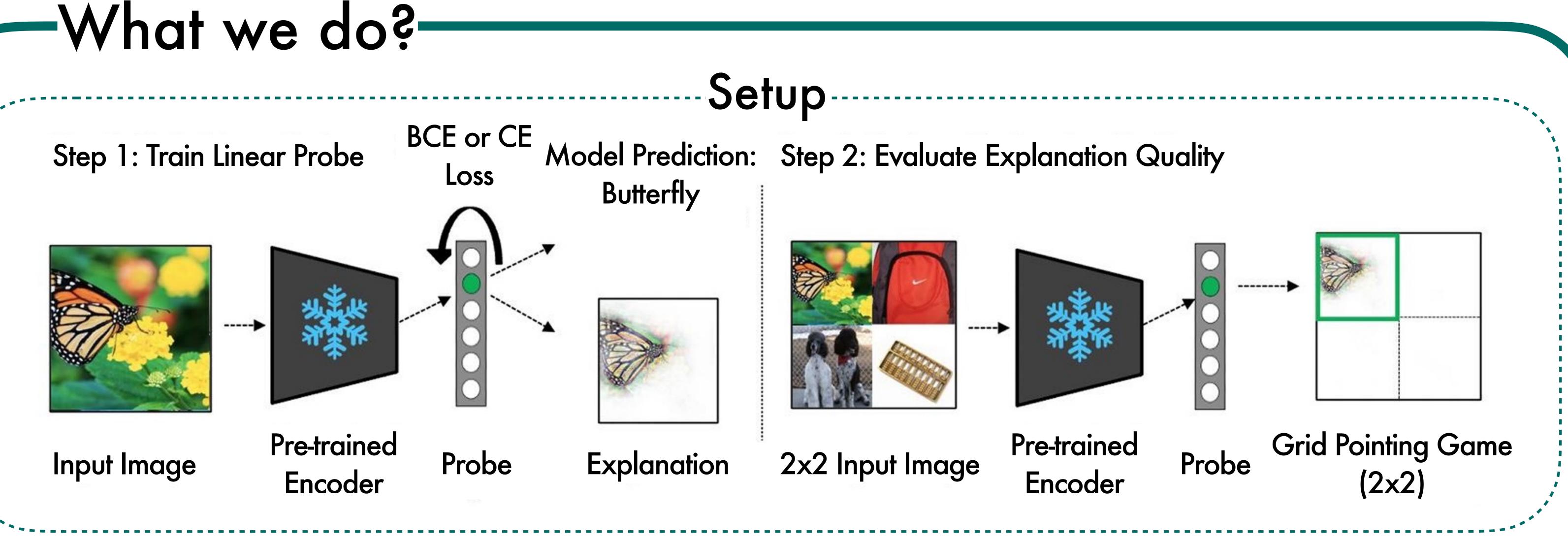




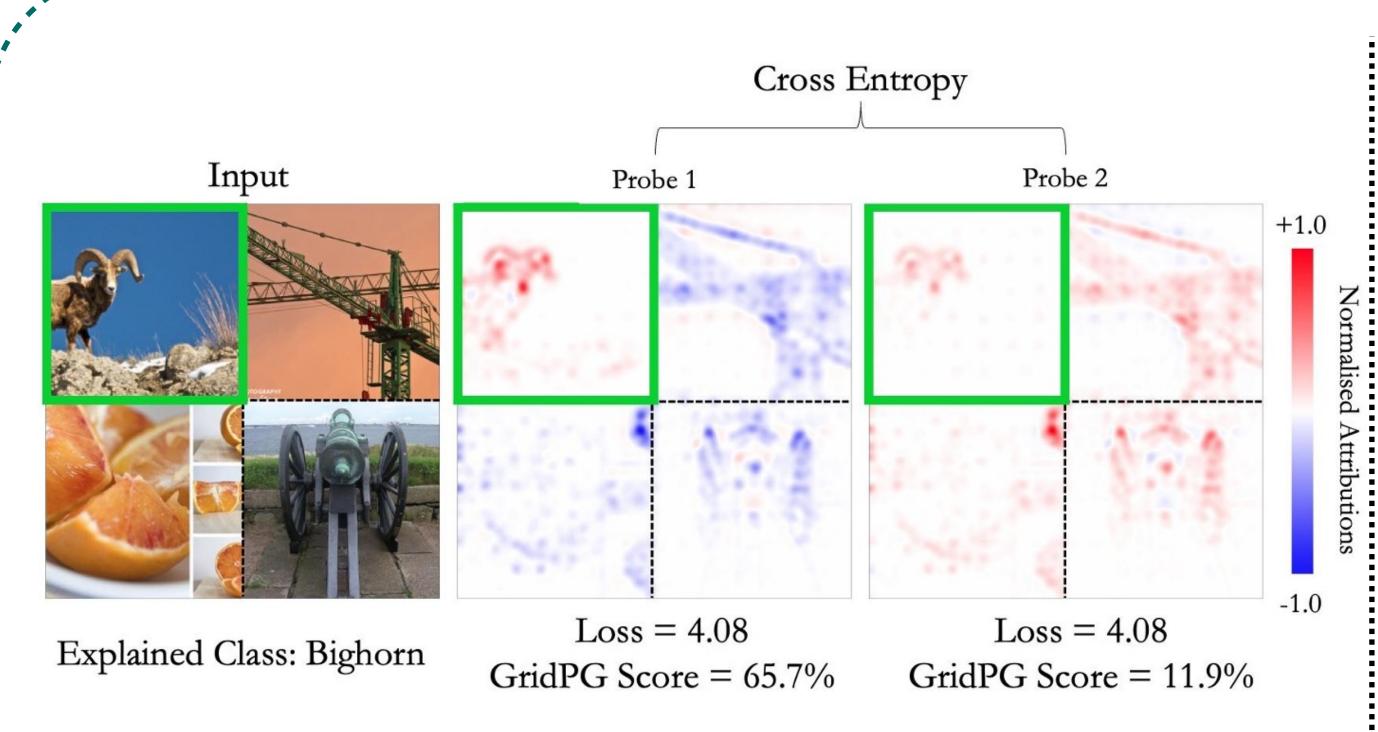


B-cos MLPs lead to improved explanations.

Training of the classification head matters more than pre-training. Just < 10% of parameters can define explanation quality.



-Softmax Invariance-Issue-



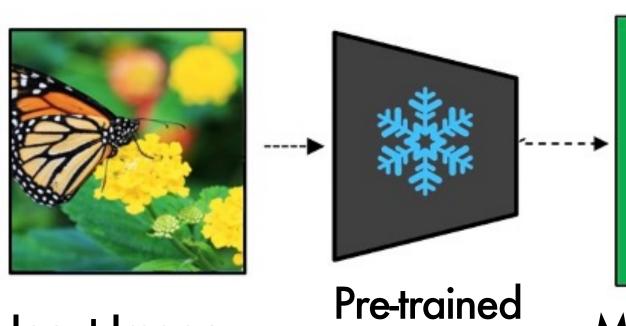
CE Loss is invariant to adding a shift δ to all output logits.

Image-specific shifts result from $w_k \rightarrow w'_k = w_k + \Delta w$.

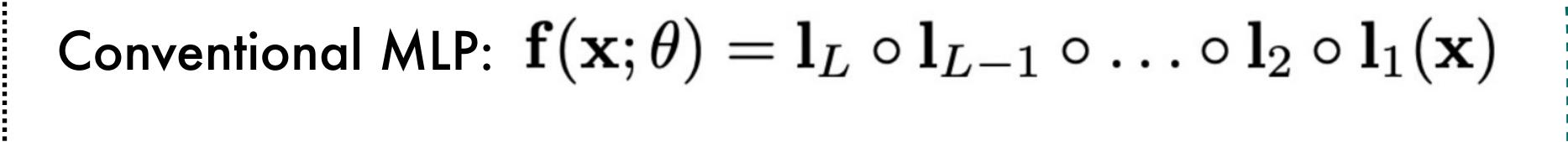
 $\mathcal{L}_{\text{BCE},i} = -\sum_{c,i} \log \left(\sigma\left(\hat{y}_{c,i}\right)\right) + \left(1 - t_{c,i}\right) \log \left(1 - \sigma\left(\hat{y}_{c,i}\right)\right)$

itive and Negative attributions NOT calibrated. BCE penalizes non-target shifts, improving class localization.

Complex Probes







B-cos Layer:
$$\mathbf{l}_l^*(\mathbf{a}_l;\mathbf{W}_l) = |c(\mathbf{a}_l;\mathbf{W}_l)|^{\mathrm{B}-1} \odot \mathbf{W}_l \, \mathbf{a}_l$$

What we find?

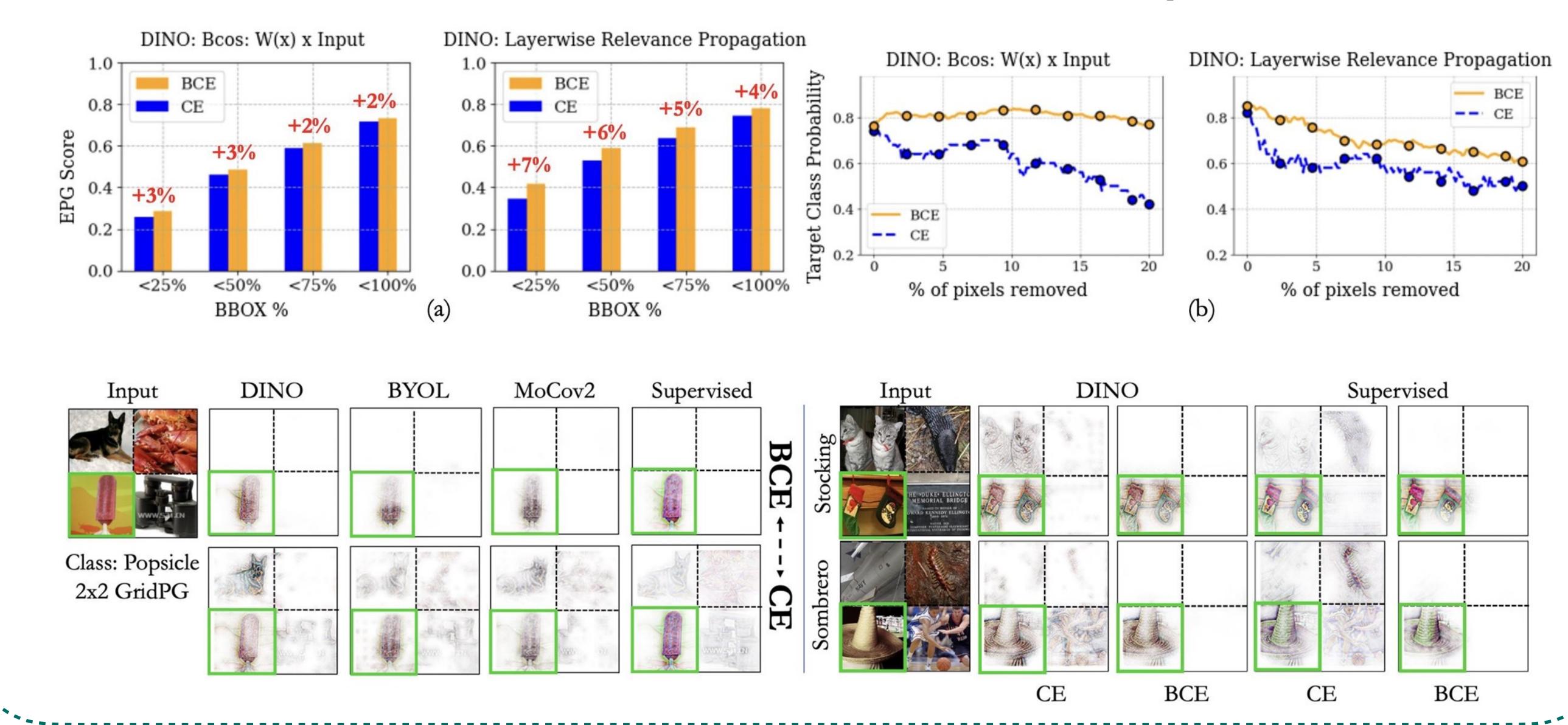
Pre-training Frameworks

- 1. Fully supervised Learning
- 2. Vision-Language Learning (CLIP [1])
 - 3. Self-supervised Learning
 - (MoCov2, BYOL, DINO)

Explanation Methods

- 1. Gradient Based: LRP, IxG, IntGrad
- 2. Activation Based: GradCAM
- 3. Perturbation Based: LIME, ScoreCAM
- 4. Inherently Interpretable: B-cos
- 5. Vision Transformer Based: CheferLRP, etc.

BCE Probes localize better than CE probes.



-B-cos MLPs improve class-specific localization.

