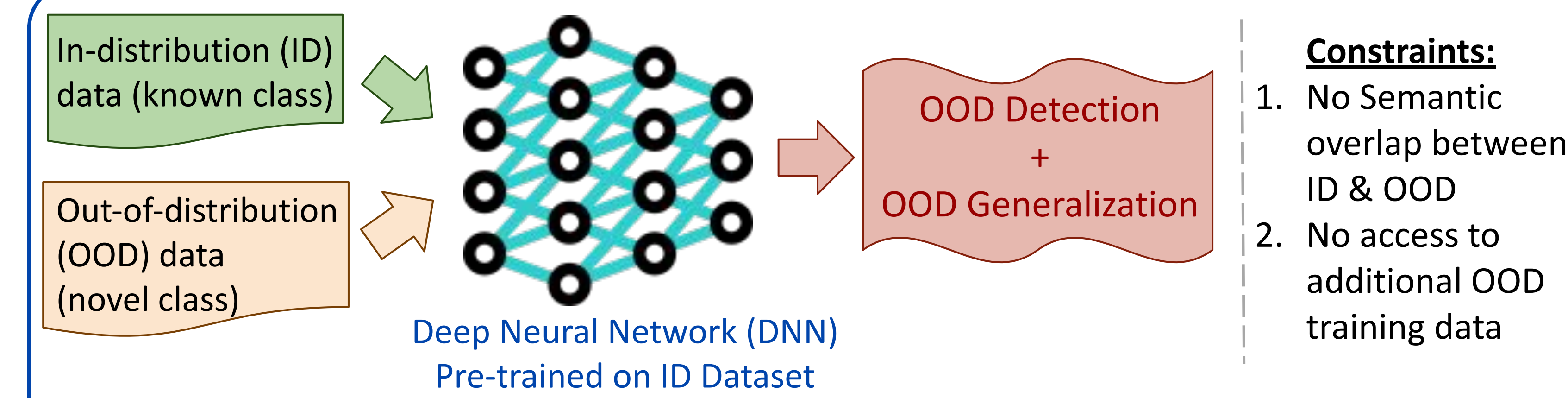


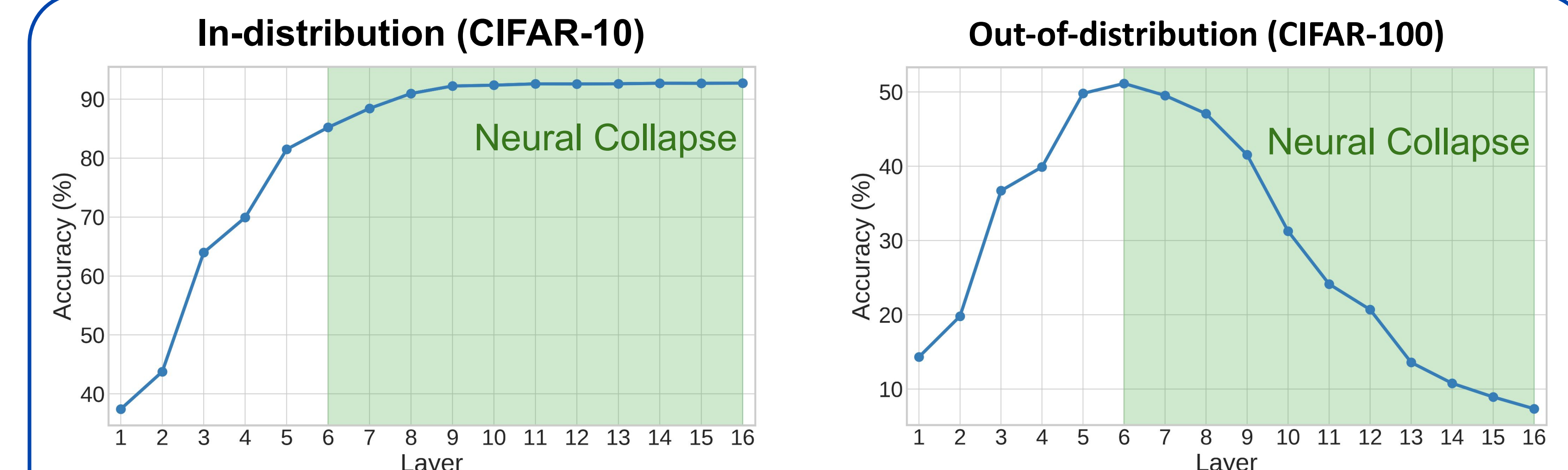


Motivation & Problem Setup



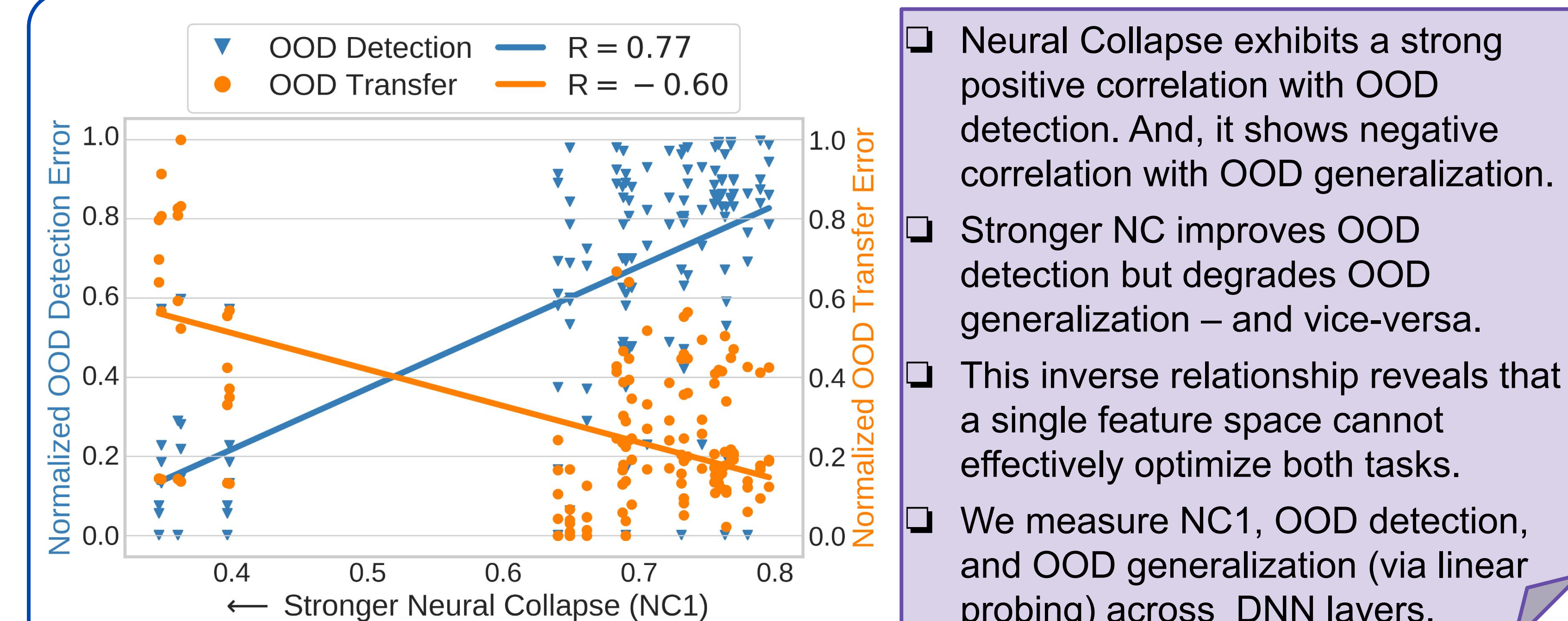
- Goal:** In open-world settings, DNNs must detect novel concepts and maximize forward transfer to facilitate efficient learning.
- Research Question:** How can we build representations in a DNN to simultaneously achieve both OOD detection and generalization?
- Challenge:** Optimizing OOD detection hurts OOD generalization and vice-versa.
- TL;DR:** We developed a method for jointly optimizing the OOD detection and forward transfer (OOD generalization) based on the **Neural Collapse** phenomenon.

Neural Collapse Degrades OOD Generalization



- Neural Collapse (NC):** NC is a phenomenon where DNNs develop compact and structured class representations. While typically seen in the final layer, it can also occur to varying degrees in the last K layers – known as *intermediate* NC.
- In our NeurIPS 2024 paper “*What Variables Affect Out-of-Distribution Generalization in Pretrained Models?*”, we show that intermediate NC degrades OOD generalization.
- As shown above, linear probe ID accuracy monotonically increases as a function of layers, but OOD accuracy only increases until the *Neural Collapse* is reached and then decreases.
- In this follow-up work, we demonstrate that the degree of NC plays a major role in both OOD detection and generalization.

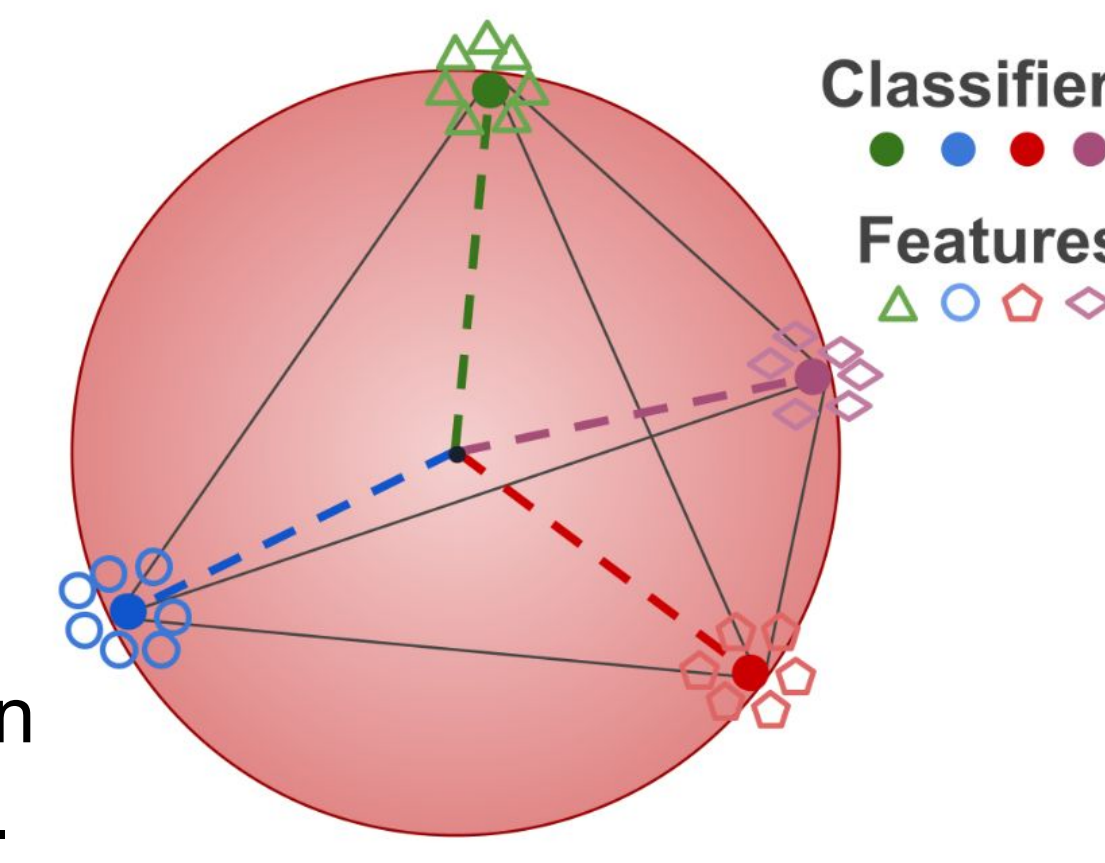
Core Insight: Neural Collapse Relates to OOD Detection & Generalization



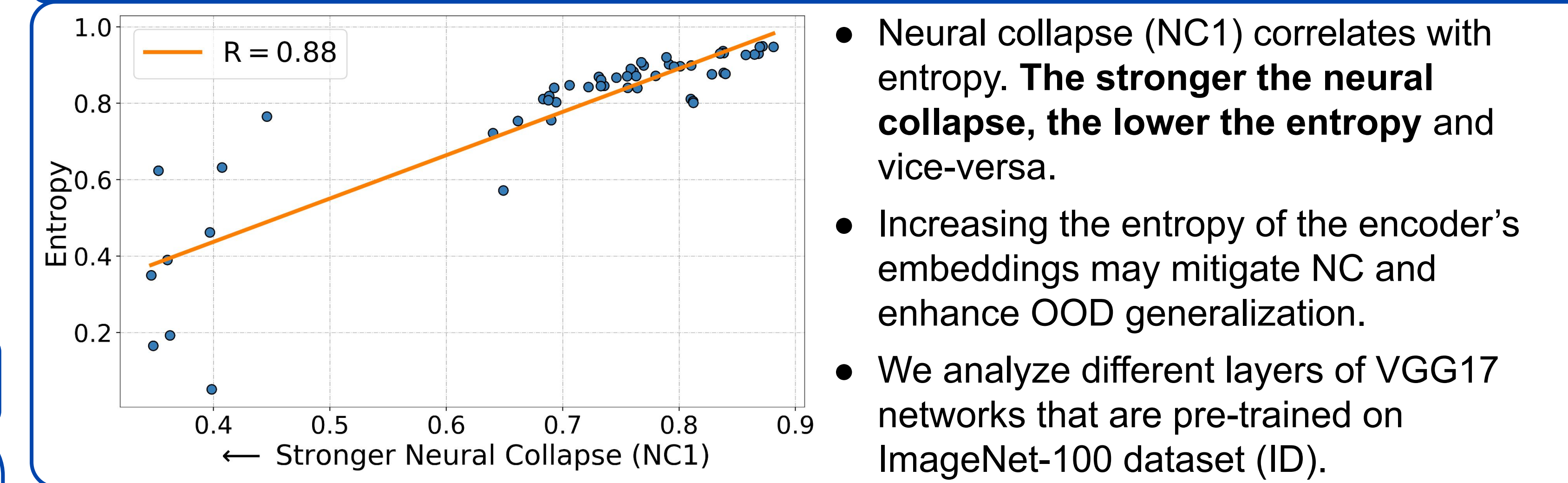
- Neural Collapse exhibits a strong positive correlation with OOD detection. And, it shows negative correlation with OOD generalization.
- Stronger NC improves OOD detection but degrades OOD generalization – and vice-versa.
- This inverse relationship reveals that a single feature space cannot effectively optimize both tasks.
- We measure NC1, OOD detection, and OOD generalization (via linear probing) across DNN layers.

Neural Collapse Criteria

- Feature Collapse (NC1):** Intra-class features collapse to a single mean with low variability.
- Simplex ETF (NC2):** Class means, centered at the global mean, form a maximally spaced simplex on a hypersphere.
- Self-Duality (NC3):** Classifiers align tightly with class means, creating a nearly self-dual configuration.
- Nearest Class Mean (NCM) Decision (NC4):** Classification resembles a NCM scheme, based on class mean proximity.

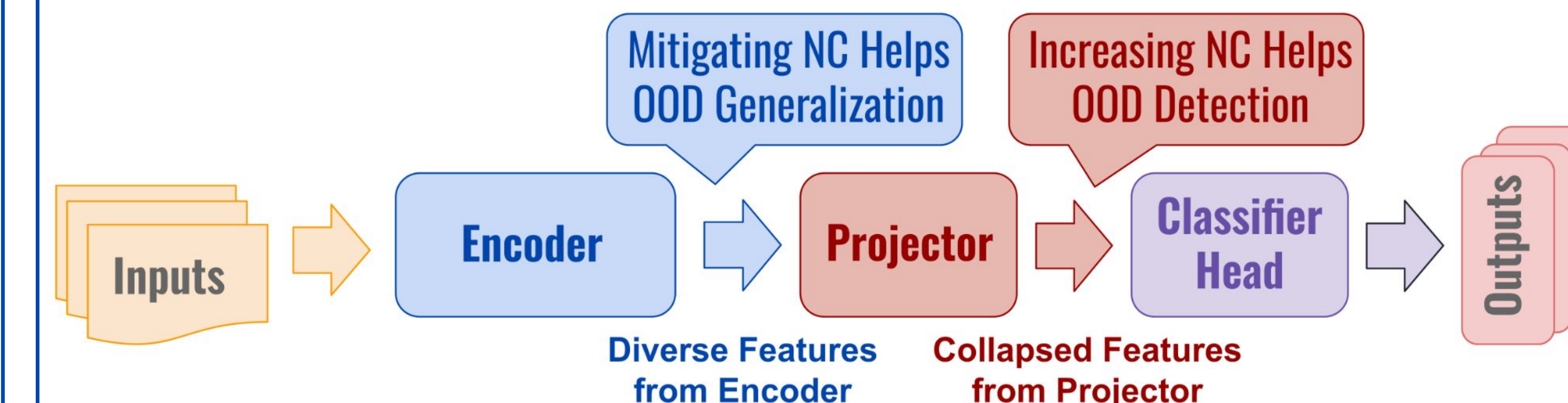


Correlation between Neural Collapse and Entropy

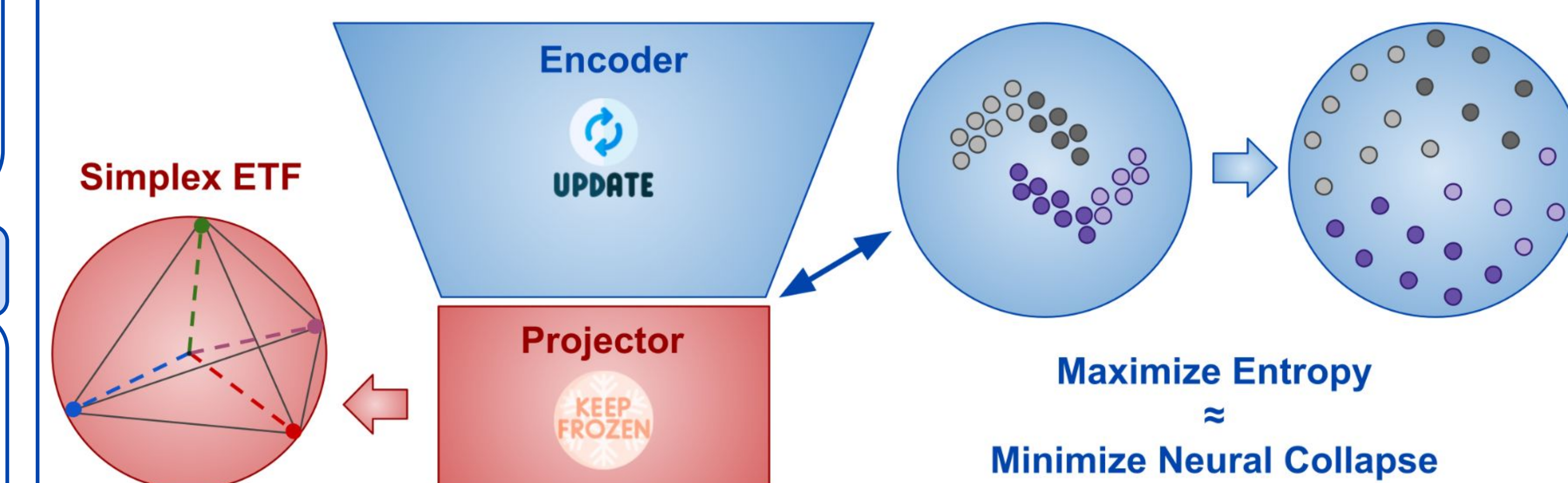


- Neural collapse (NC1) correlates with entropy. **The stronger the neural collapse, the lower the entropy** and vice-versa.
- Increasing the entropy of the encoder's embeddings may mitigate NC and enhance OOD generalization.
- We analyze different layers of VGG17 networks that are pre-trained on ImageNet-100 dataset (ID).

Method Overview: Controlling Neural Collapse



- A single feature space cannot effectively achieve both OOD detection and generalization.
- To address this, we control NC at different DNN layers, using an encoder optimized for generalization and a projector tailored for detection.



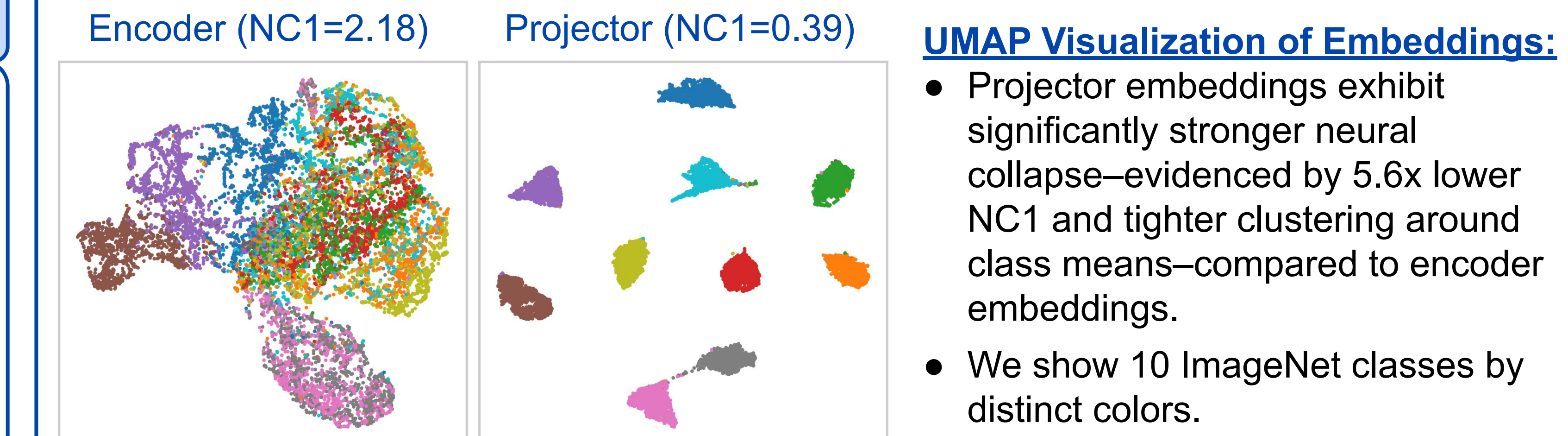
- Layer for OOD generalization:** we introduce entropy regularization to mitigate NC in the encoder, improving feature diversity for OOD generalization.
 - We develop a theoretical framework that explains how entropy regularization mitigates NC. In particular, we show that collapsing implies entropy diverges to negative infinity.
 - For the entropy regularization, we leverage nearest-neighbor-based density estimation.
- Layer for OOD detection:** we leverage a fixed simplex Equiangular Tight Frame (ETF) projector to induce NC in the final layer, improving feature compactness for detection.
 - For the ETF projector, we configure a two-layer MLP as a simplex ETF (equinorm and maximum equiangularity) and keep it frozen during training.

Experimental Setup

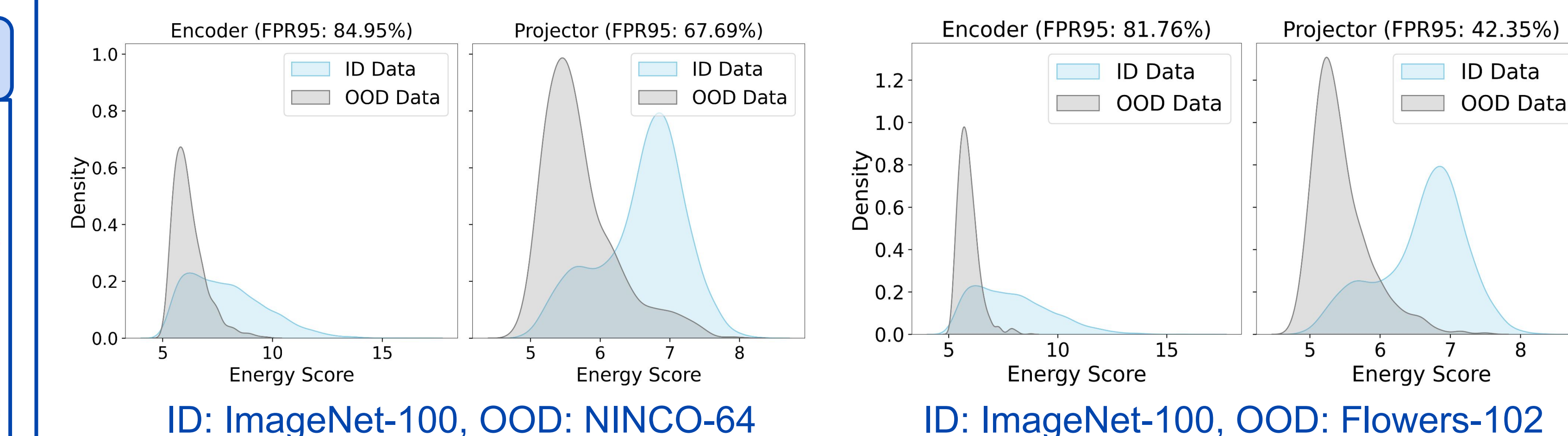
- Dataset:** Used ImageNet-100 as ID dataset and eight OOD datasets: NINCO, CUB-200, CIFAR-100, ImageNet-R, Oxford 102 Flowers, Aircrafts, Oxford Pets, and STL-10
- Architecture:** VGG17, ResNet18, ResNet34, ViT-Tiny, and ViT-Small
- NC Evaluation:** Four NC metrics NC1, NC2, NC3, and NC4 characterized by NC criteria. A lower NC indicates stronger Neural Collapse and vice-versa.
- Metrics:** ID generalization error, OOD generalization error, OOD detection error.

Qualitative Results: Encoder Vs. Projector

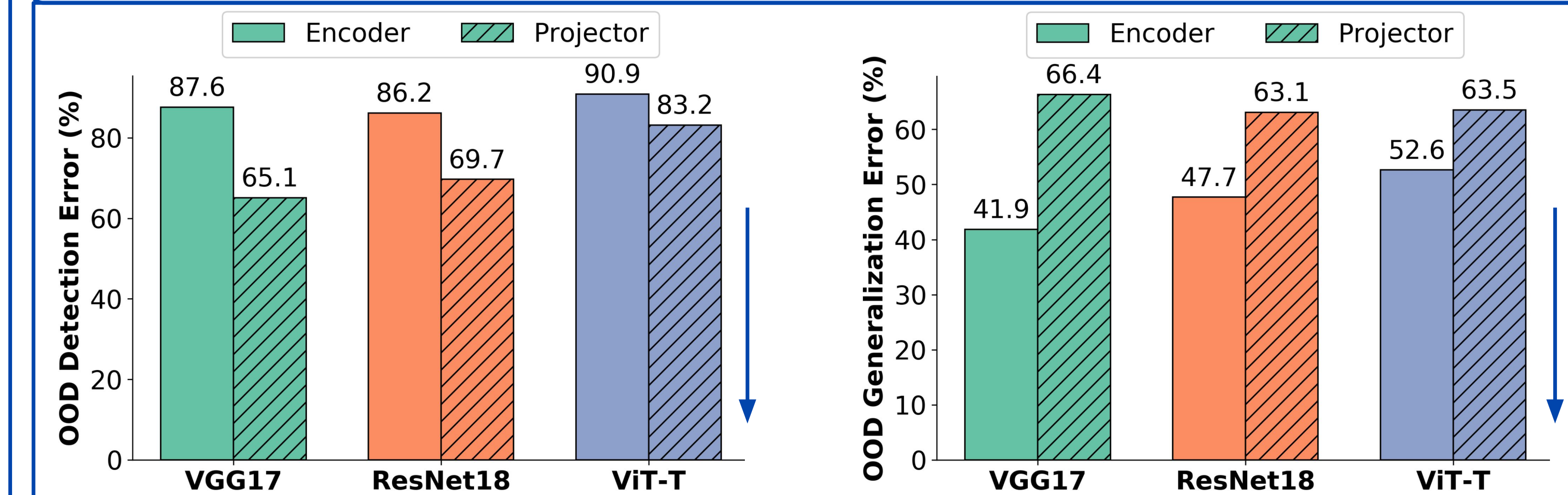
The following results are based on VGG17 models pre-trained on ImageNet-100 dataset (ID).



OOD Detection: As shown in energy score distribution, the projector creates a greater separation between ID & OOD data and achieves a lower FPR95 than the encoder.



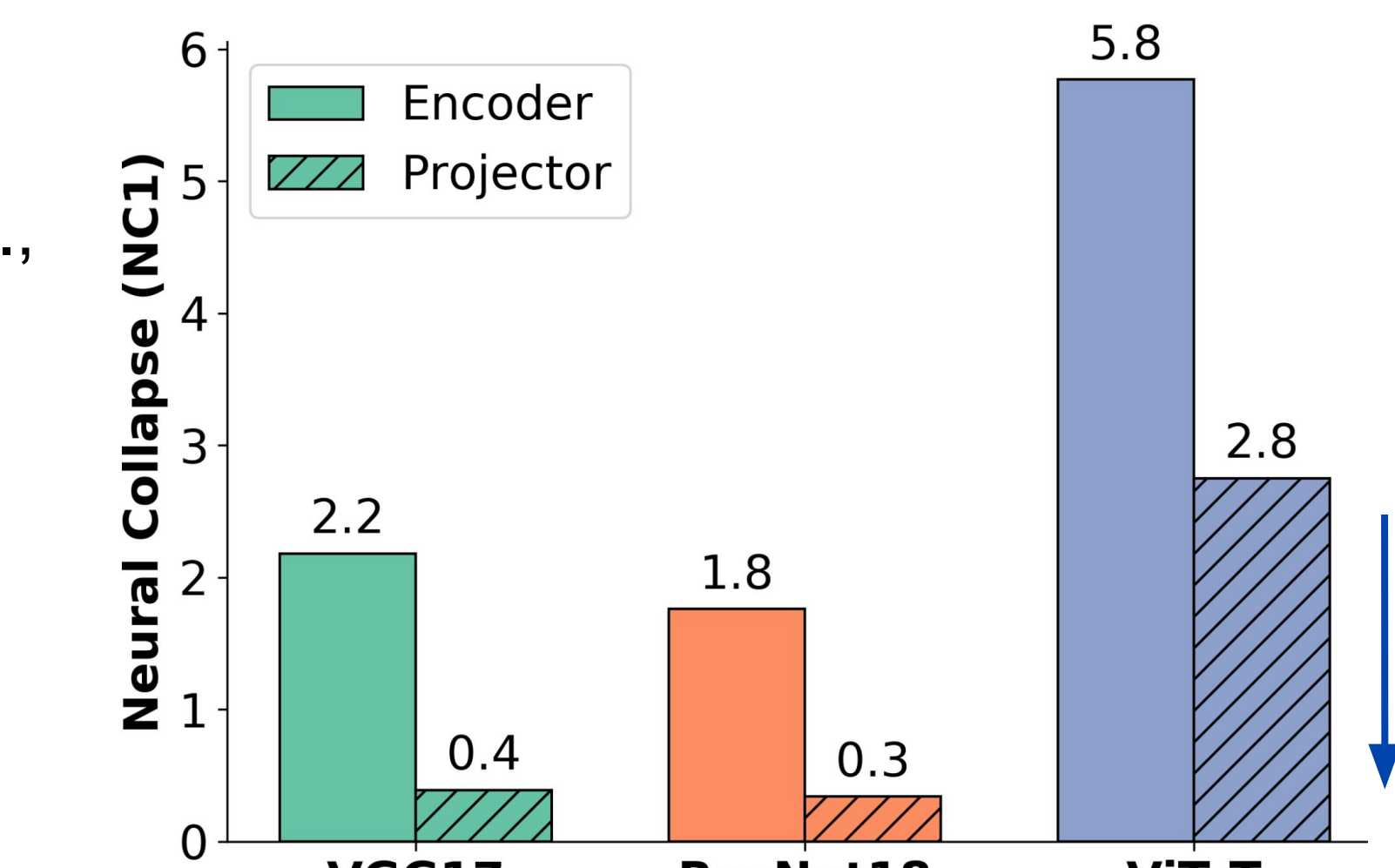
Quantitative Results: Encoder Vs. Projector



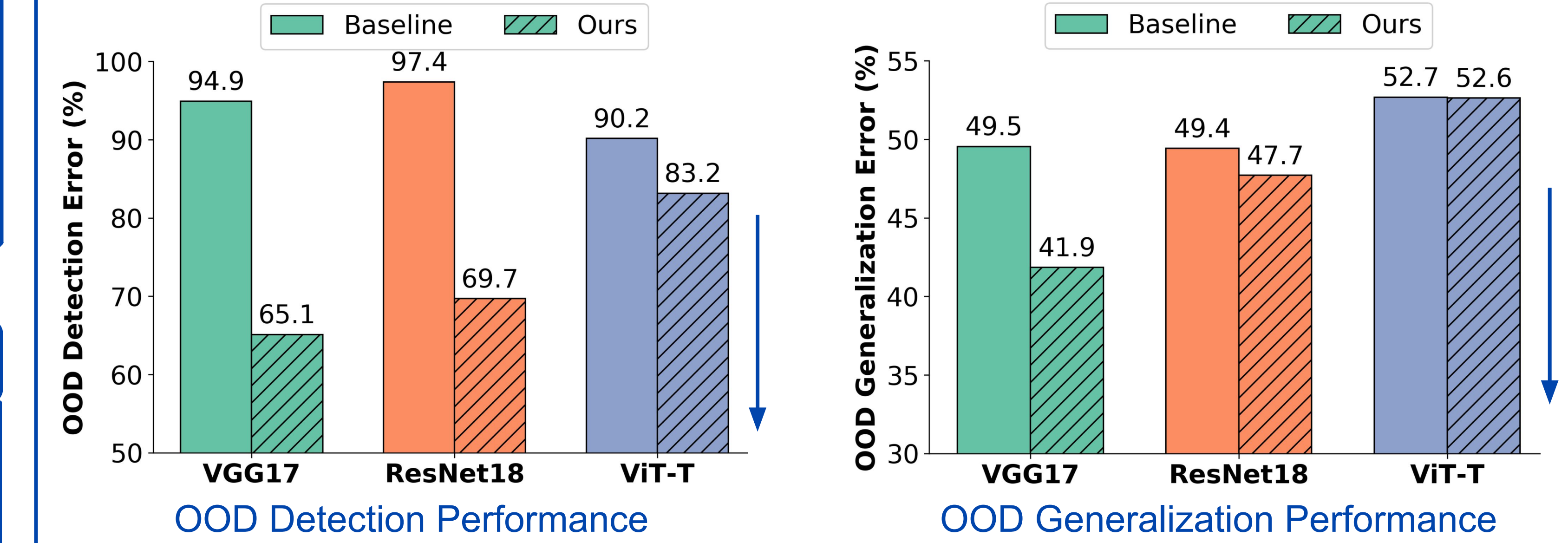
- The encoder mitigates NC and becomes a better OOD generalizer than the projector.
- The projector intensifies NC and becomes a better OOD detector than the encoder.

Neural Collapse Evaluation:

- The projector exhibits lower NC1 values (i.e., stronger neural collapse) than the encoder across DNN architectures.
- We report NC1 (feature collapse), the most dominant indicator of neural collapse.
- All above results are averaged across eight OOD datasets



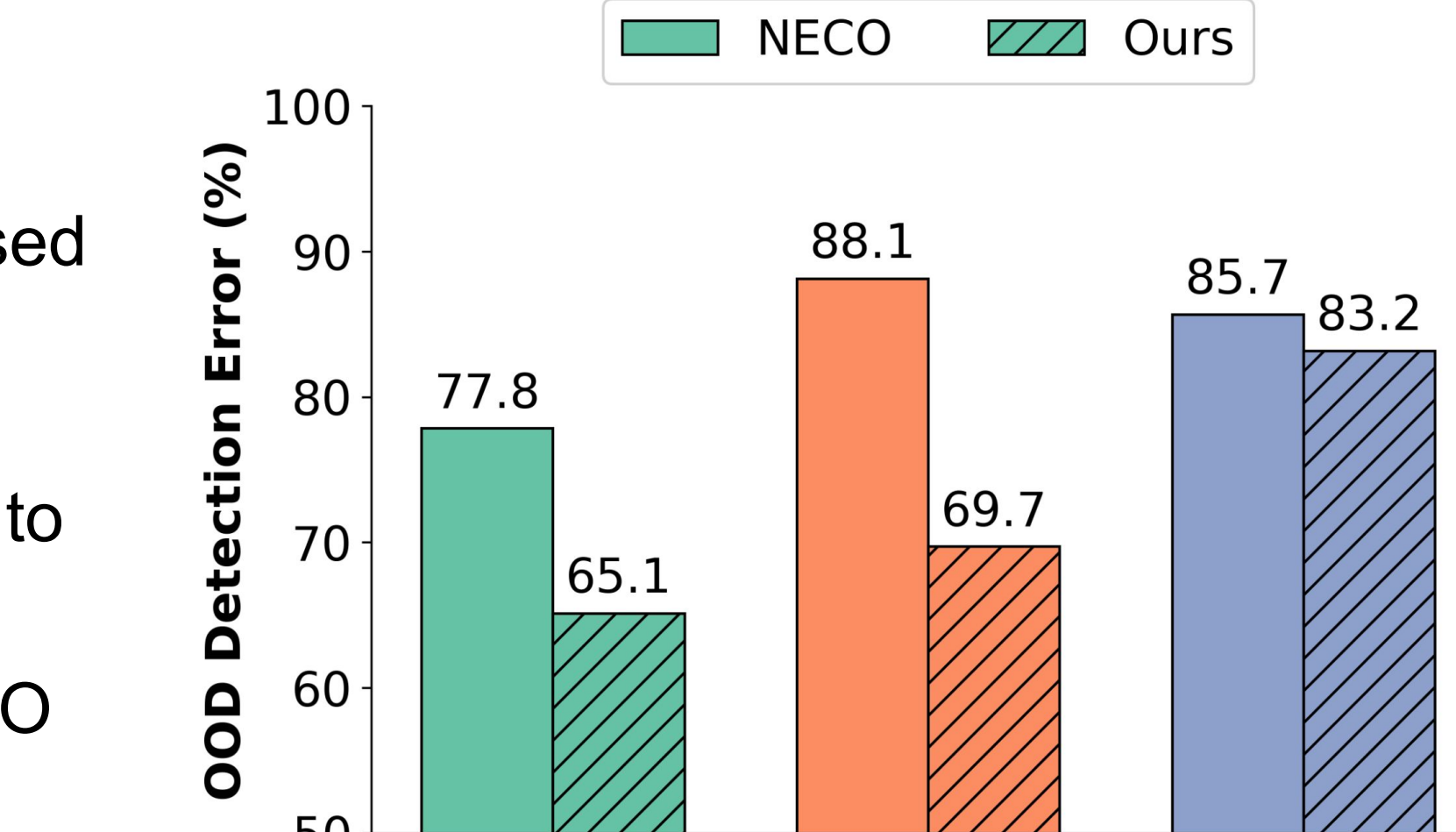
Comparison with Baseline



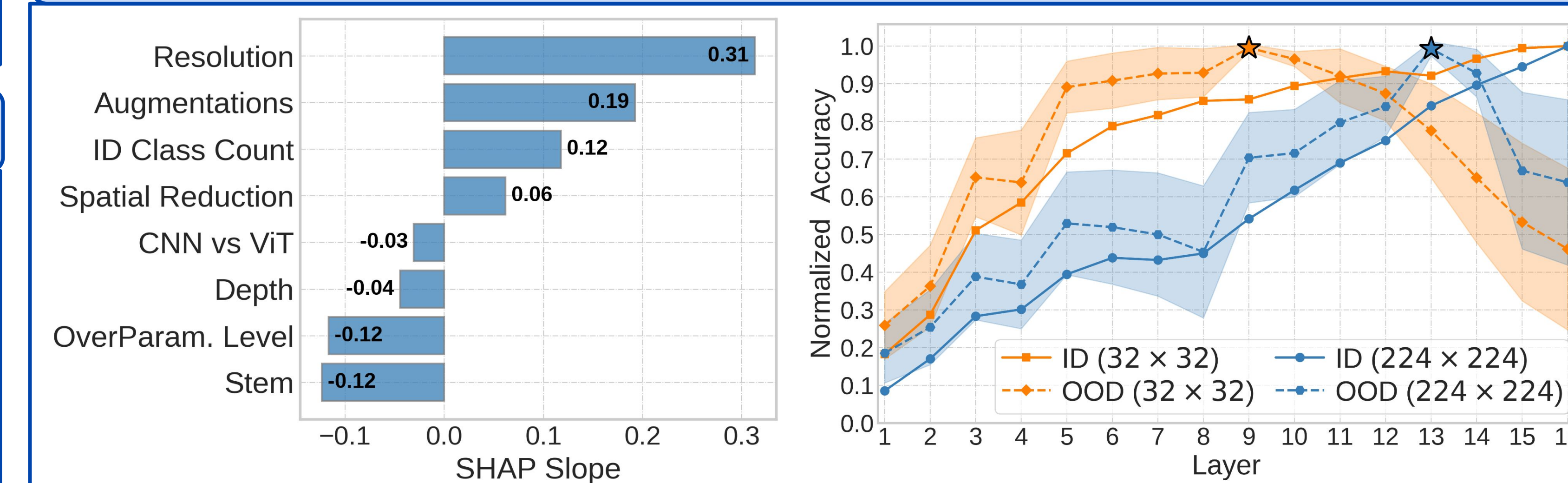
- Baseline DNNs lack mechanisms to control NC, resulting in poor performance.
- Our method controls NC and achieves significant improvements over these baselines.
- Results are based on DNNs pre-trained on ImageNet-100 dataset (ID), with performance averaged across eight OOD datasets.

Comparison with SOTA OOD Detector:

- We compare our method with NECO, a state-of-the-art OOD detection method based on NC properties.
- Since NECO does not address OOD generalization, we restrict this comparison to OOD detection only.
- Our method consistently outperforms NECO across all settings.



How Do Variables Impact Neural Collapse?



In our NeurIPS-2024 paper, we study how variables impact NC and find that:

- Our SHAP analysis reveals that image resolution is the most dominant variable followed by augmentations and ID class count in terms of reducing NC and enhancing transfer.
- The characteristics of toy datasets e.g., CIFAR lead to *sub-optimal* representations that hinder OOD generalization, explaining why methods successful on such datasets frequently fail on real-world datasets e.g., ImageNet.
- Increasing ID class count (*between-class diversity*), using augmentations (*within-class diversity*), and using higher image resolution (*hierarchical features*) greatly reduce NC.
- Increasing dataset diversity significantly reduces NC, and with sufficient diversity, it can be entirely prevented.
- In this follow-up work, we show that entropy regularization offers an alternative means to mitigate NC and enhance OOD generalization.

Acknowledgments

We thank US National Science Foundation for supporting our research.