

Out-of-Distribution Detection and Neural Collapse

Fares Boudeaa, Tiena Soro

Abstract

This report studies Out-of-Distribution (OOD) detection methods and the Neural Collapse phenomenon in deep networks. We train a ResNet-18 classifier on CIFAR-100 and compare multiple OOD scoring methods, including MSP, Energy, Mahalanobis, and ViM. We then analyze Neural Collapse (NC1–NC5) at the end of training and implement the NECO method for OOD detection.

Introduction

Motivation

Deep neural networks often produce overconfident predictions on out-of-distribution samples. Reliable uncertainty estimation is therefore crucial in safety-critical applications such as autonomous driving, medical imaging, and fraud detection. In these settings, a model that silently assigns high confidence to inputs it has never seen can cause serious harm. OOD detection addresses this by equipping a classifier with the ability to recognize when an input falls outside its training distribution, and to flag it accordingly rather than producing a potentially misleading prediction. In this report, we study several OOD scoring methods and explore how the geometric structure that emerges at the end of training, known as Neural Collapse, can be leveraged for this purpose.

Setup

Residual Networks (ResNets) were originally designed for large-scale image classification on ImageNet, which consists of high-resolution 224×224 images. We adapt ResNet-18 for CIFAR-100 by replacing the first convolutional layer with a smaller kernel and stride, and removing the initial max-pooling layer, so that the spatial resolution of the small 32×32 inputs is preserved through the early layers of the network.

We train a ResNet-18 adapted for CIFAR-100 from scratch and analyze the feature representations of dimension $H = 50000 \times 512$ at the penultimate layer before the classification head.

The CIFAR-100 dataset consists of 60,000 color images of size 32×32 divided into 100 classes, further grouped into 20 superclasses. Every class has 600 examples, 500 for training and 100 for testing. It mainly contains images of

everyday objects, animals, vehicles, and vegetation. We perform data augmentation for a more robust model: random crop with padding, random horizontal flip, and normalization with CIFAR-100 statistics. We then train the ResNet-18 model until the training loss reaches 0, even though the training accuracy already reached 99.9%, to ensure that the model reaches the terminal training phase.

Training Procedure

Optimization We train using Stochastic Gradient Descent (SGD) with momentum combined with hyperparameter tuning, as it outperforms Adam for this setting. The chosen parameters are the following:

- Optimizer: SGD with momentum 0.9
- Weight decay: 5×10^{-4}
- Initial learning rate: 0.1
- LR schedule: cosine decay
- Loss: Cross Entropy Loss

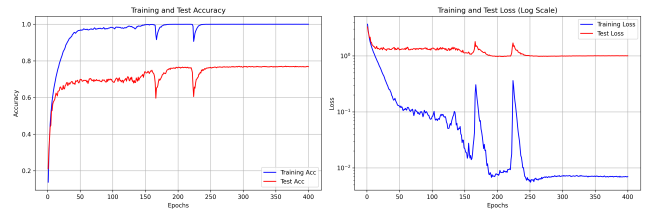


Figure 1: Training and Test Loss and Accuracy

Training Duration In order to evaluate the out-of-distribution metrics as well as evaluate NECO for OOD detection, a different dataset containing classes not present in CIFAR-100 is used. Here we chose the Street View House Numbers (SVHN) dataset, which contains labeled images of house numbers extracted from Google Street View. We compare both datasets' means and standard deviations to confirm that the two datasets are not very similar:

- CIFAR-100 : mean = (0.4914, 0.4822, 0.4465)
- SVHN : mean = (0.4377, 0.4438, 0.4728)
- CIFAR-100 : std = (0.2023, 0.1994, 0.2010)
- SVHN : std = (0.1980, 0.2010, 0.1970)

In the following we use the CIFAR-100 train dataset for the training and the SVHN test dataset as OOD for the evaluation of our different metrics.

We used the same normalization parameters (mean = (0.4914, 0.4822, 0.4465), std = (0.2023, 0.1994, 0.2010)) as in the paper [(2)]. The training curves and test accuracy can be observed in Fig.[1]

OOD Scoring Methods

Problem Setup

Let $(x, y) \sim \mathcal{D}_{\text{ID}}$ denote in-distribution data. At test time, inputs may instead come from \mathcal{D}_{OOD} .

The goal of OOD detection is to define a score function $S(x)$ such that:

$$S(x_{\text{ID}}) > S(x_{\text{OOD}})$$

Evaluation metrics: *AUROC*, *AUPR*, *FPR95*

Let $f(x) \in \mathbb{R}^K$ be the logits of the network.

Max Softmax Probability (MSP)

$$S_{\text{MSP}}(x) = \max_k \text{softmax}(f(x))_k$$

Maximum Logit Score

$$S_{\text{Logit}}(x) = \max_k f_k(x)$$

Energy Score

$$S_{\text{Energy}}(x) = -T \log \sum_{k=1}^K e^{f_k(x)/T}$$

Mahalanobis Score

Let μ_c be class means and Σ the shared covariance matrix.

$$S_{\text{Mahalanobis}}(x) = -\min_c (h(x) - \mu_c)^T \Sigma^{-1} (h(x) - \mu_c)$$

ViM

ViM decomposes features into principal and residual subspaces:

$$h(x) = h_{\parallel}(x) + h_{\perp}(x)$$

The score combines logit magnitude and residual norm.

NECO: Neural Collapse Inspired OOD Detection

Motivation

OOD samples violate the geometric structure induced by Neural Collapse. In-distribution samples, at the end of training, have features that cluster tightly around their respective class means, which are arranged in a Simplex ETF configuration. OOD samples do not belong to any of the trained classes and therefore do not align with this structure. NECO exploits this by measuring how much of a feature vector lies in the subspace spanned by the class means: an ID sample should project strongly onto this subspace, while an OOD sample should not.

Method

Let P denote the orthogonal projector onto the subspace spanned by the centered class mean vectors. The NECO score is defined as the relative projection of the feature onto this subspace:

$$\text{NECO}(x) = \frac{\|Ph_{\omega}(x)\|}{\|h_{\omega}(x)\|} = \frac{\sqrt{h_{\omega}(x)^{\top} P P^{\top} h_{\omega}(x)}}{\sqrt{h_{\omega}(x)^{\top} h_{\omega}(x)}}$$

This score is large when the feature vector aligns with the class-mean subspace, which happens for in-distribution inputs. For OOD inputs, the feature vector has a significant component orthogonal to the class-mean subspace, making the score smaller. The score is scale-invariant by construction, which makes it robust to the norm growth described in NC1.

We evaluate the scoring methods with AUROC (Area Under the ROC Curve), FPR 95 (False Positive Rate at 95% TPR), and AUPR (Area Under the Precision-Recall Curve). In order to compute and evaluate these scores, we used the framework of OOD detection(OOD is the positive class). In the benchmark Fig.[9], we also compute AUPR-Out which is the other case(where ID is positive class). So our AUPR is in fact AUPR-In in this benchmark.

Among all tested methods, NECO achieves the best performance for every metric used. Looking at the distributions in Fig. [9], we see that NECO separates the ID and OOD data better than other methods. This confirms the fact that the NECO scoring method is suitable for OOD detection. The performance of each scoring method is displayed in Table.[1]. We also have a complete visualization of these metrics in Fig. [9]

Comparison of OOD Methods

| Method | AUROC \uparrow | AUPR \uparrow | FPR95 \downarrow |
|-------------|------------------|-----------------|--------------------|
| MSP | 83.55 | 86.05 | 48.71 |
| MLS | 84.5 | 86.92 | 46.47 |
| Energy | 85.00 | 87.24 | 46.18 |
| Mahalanobis | 87.99 | 89.04 | 45.02 |
| ViM | 84.39 | 87.87 | 41.66 |
| NECO | 99.65 | 99.92 | 1.778 |

Table 1: OOD detection results on CIFAR-100 vs. SVHN.

Neural Collapse

Neural Collapse emerges at the terminal phase of training when cross-entropy loss approaches zero.

Let $h_{i,c}$ denote features of sample i from class c .

NC1: Variability Collapse

Within-class covariance:

$$\Sigma_W = \frac{1}{N} \sum_{i,c} (h_{i,c} - \mu_c)(h_{i,c} - \mu_c)^{\top} \rightarrow 0$$

In practice we use :

$$NC1 = \text{Tr} \left[\frac{\Sigma_W \Sigma_B^\dagger}{C} \right]$$

We use this metric instead of computing Σ_W directly because as training progresses the norms of $h_{i,c}$ increase. This is due to gradient descent: $w = w - \eta \nabla L$. As L decreases, the gradient is negative, so $-\eta \nabla L$ is positive and the weights keep growing in magnitude. This is seen in the Fig.8. In classification, Σ_W captures within-class noise while Σ_B represents the between-class covariance:

$$\begin{aligned} h_{i,c} &= \mu_c + \epsilon_{i,c} \\ \Sigma_B &= \frac{1}{C} \sum_c (\mu_c - \mu_G)(\mu_c - \mu_G)^T \\ \Sigma_W &= \frac{1}{N} \sum_{i,c} \epsilon_{i,c} \epsilon_{i,c}^T \end{aligned}$$

By using $\text{Tr}[\Sigma_W \Sigma_B^\dagger / C]$, we obtain a scale-invariant measure of within-class variability relative to between-class spread. As this ratio converges to zero, the within-class noise becomes negligible compared to the separation between classes, which is precisely the condition of variability collapse.

During our experiments with the Resnet-18 model, we got a $NC1$ value of 0.04. We can observe the evolution of $NC1$ through epochs on Fig.[2]

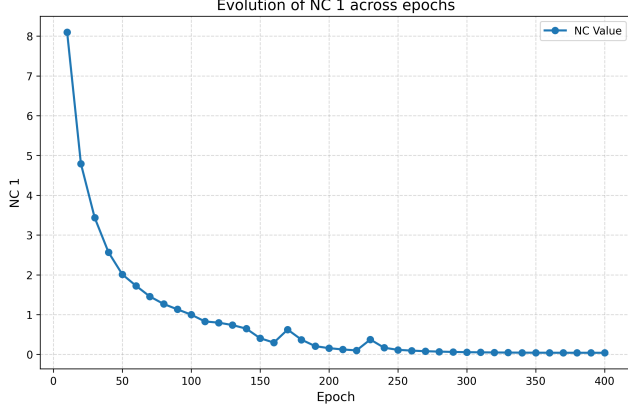


Figure 2: Evolution of $NC1$ through epochs.

NC2: Simplex ETF Structure

Here, the class means become equiangular and equal norm:

$$\begin{aligned} \|\mu_c - \mu_G\|_2 - \|\mu_{c'} - \mu_G\|_2 &\rightarrow 0 \quad \forall c, c' \\ \langle \tilde{\mu}_c, \tilde{\mu}_{c'} \rangle &\rightarrow \frac{C}{C-1} \delta_{c,c'} - \frac{1}{C-1} \quad \forall c, c' \end{aligned}$$

In practice we use:

$$\text{EN}_{\text{class-means}} = \frac{\text{std}_c \{\|\mu_c - \mu_G\|_2\}}{\text{avg}_c \{\|\mu_c - \mu_G\|_2\}}$$

and equi-angularity

$$\text{Equi_ang}_{c-m} = \text{Avg}_{c,c'} \left| \frac{\langle \mu_c - \mu_G, \mu_{c'} - \mu_G \rangle + \frac{1}{C-1}}{\|\mu_c - \mu_G\|_2 \|\mu_{c'} - \mu_G\|_2} \right|$$

The practical metrics are self-evident. For the first, we compute the standard deviation of the norms of the class mean vectors centered around the global mean, divided by the average norm to account for the increasing magnitude mentioned earlier. The geometric interpretation of this metric approaching 0 is that all class mean vectors are at equal distance from the global mean in the high-dimensional feature space. The second metric computes the average pairwise deviation from the ideal cosine similarity $-\frac{1}{C-1}$. Geometrically, as this metric approaches 0, the angle between any two centered class means converges to the maximum possible value, which corresponds to the Simplex ETF structure.

For our model, we got a $\text{EN}_{\text{class-means}} = 0.029$ and $\text{Equi_ang}_{c-m} = 0.052$. As expected, the values are very small. This validates the $NC2$ property. We can now observe the evolution of these values through epochs Fig.[3]. We also plot the cosine similarity heatmap and the distributions of pairwise angles in Fig.[8]

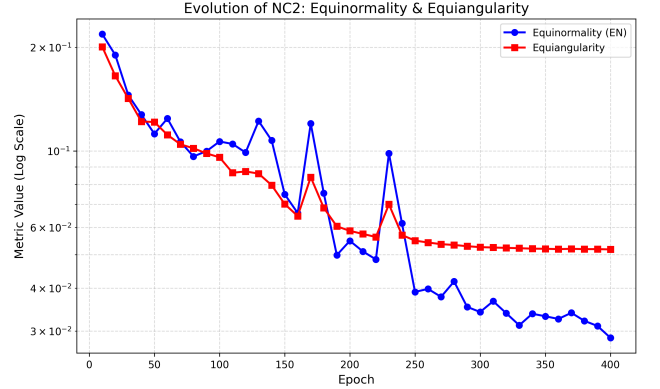


Figure 3: Evolution of Equinormality and Equiangularity through epochs

NC3 : Convergence to Self-Duality

Here, the classifier weights w_c align with class means:

$$NC3 = \left\| \frac{W^T}{\|W\|_F} - \frac{\dot{M}}{\|\dot{M}\|_F} \right\|_F \rightarrow 0$$

NC3 states that the classifier weight matrix and the class mean matrix become proportional at the end of the training. The normalized weight matrix $\frac{W^T}{\|W\|_F}$ converges to the normalized class mean matrix $\frac{\dot{M}}{\|\dot{M}\|_F}$, meaning that the classifier weights are the dual basis of the class means. In other words, the linear classification head no longer learns a separate geometric structure from the features; it exactly mirrors the arrangement of the class mean vectors. This is a direct consequence of the simplex ETF structure, since both the

class means and the classifier weights converge to the same equiangular tight frame.

In our experiments, we got $NC3 = 0.178$. This value is small, but not that small. Looking at the evolution of $NC3$ values through epochs in Fig.[4], we see that these values seem to stabilize around 0.16 from 250 epochs. So maybe we needed to train for more epochs to get smaller values.

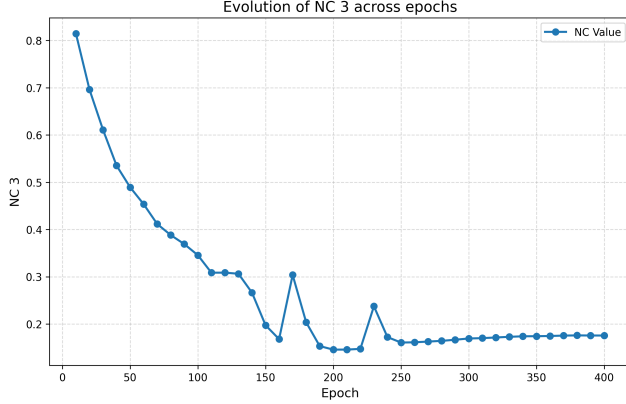


Figure 4: Evolution of $NC3$ through epochs

NC4: Nearest Class Center Rule

Classification becomes equivalent to:

$$\arg \max_{c'} \langle w_{c'}, h \rangle + b_{c'} \rightarrow \arg \min_{c'} \|h - \mu_{c'}\|_2$$

In practice, we measure the "disagreement" between the actual model predictions and the NCM (Nearest Class Mean) predictions. As collapse occurs, we have:

$$NC4 = \frac{1}{N} \sum_{i=1}^N 1 \left(\text{model}(x_i) \neq \arg \min_c \|h(x_i) - \mu_c\|_2 \right) \rightarrow 0$$

NC4 is a direct consequence of NC2 and NC3. Since the weight vectors align with the class means and all class means have equal norm, the dot product between a feature vector and a class weight vector becomes a monotone function of the Euclidean distance to the corresponding class mean. The linear classifier therefore degenerates to a simple nearest-class-center decision rule: classifying a sample amounts to finding the class whose mean is closest to the feature vector in Euclidean distance. This behavioral simplification is what makes the classifier's geometry interpretable and is also exploited by OOD methods such as Mahalanobis and NECO.

We got a $NC4 = 8e - 05$ which is very small and thus validates the fourth property of NC

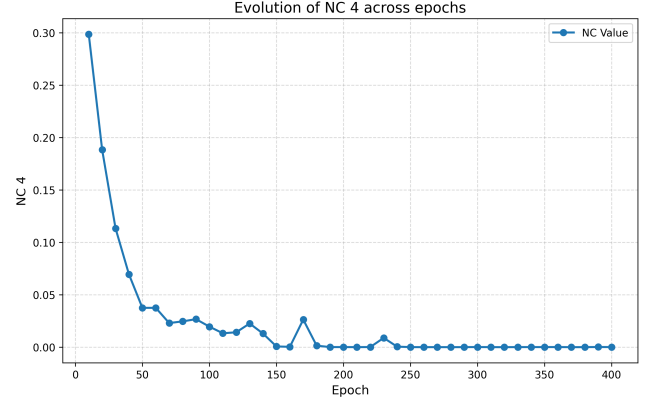


Figure 5: Evolution of $NC4$ through epochs

NC5 : ID/OOD Orthogonality

NC5 characterizes the relationship between the in-distribution class mean vectors and the global mean of OOD features.

$$\forall c, \frac{\langle \mu_c, \mu_G^{\text{OOD}} \rangle}{\|\mu_c\|_2 \|\mu_G^{\text{OOD}}\|_2} \rightarrow 0$$

In practice, we use the following metric:

$$NC5 = \text{OrthoDev}_{\text{classes-OOD}} = \text{Avg}_c \left| \frac{\langle \mu_c, \mu_G^{\text{OOD}} \rangle}{\|\mu_c\|_2 \|\mu_G^{\text{OOD}}\|_2} \right|$$

NC5 states that as training converges, the in-distribution class mean vectors become approximately orthogonal to the global mean of OOD features. Intuitively, the training objective pushes the class means to spread out in a symmetric ETF configuration centered near the origin, so their average direction in feature space becomes diffuse. An OOD sample, whose feature vector does not align with any of the trained class means, will tend to have a non-negligible projection onto the OOD global mean. This orthogonality property provides a geometric basis for separating ID and OOD samples and is directly exploited by the NECO score described in Section.[NECO].

We got $NC5 : 0.089$ which is very small and validates the fifth property of NC. Looking at the evolution of NC through epochs in Fig.[6], we find that these values seem to stabilize around 0.088 when the number of epochs increases. This can be explained by the fact that the OOD dataset (here SVHN) might share low-level features (textures, edges) with ID data, creating a "floor" for how low the $NC5$ can go.

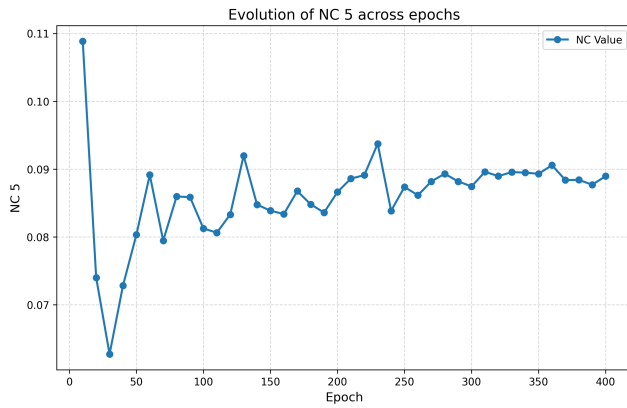


Figure 6: Evolution of NC5 values through epochs

Conclusion

Our experiments on CIFAR-100 reproduce the Neural Collapse properties reported in (1) with high fidelity. The scale-invariant NC1 metric converges to zero, the equinorm and equiangularity metrics for both class means and classifier weights collapse as predicted, the nearest-class-center disagreement vanishes at the terminal phase of training and the classifier weights align with class means. The last property of NC, the ID/OOD orthogonality, introduced in (2) is also verified. However, the result here differs slightly from the idealized theory. Indeed, the values of $NC5 = \text{OrthoDev}_{\text{classes-OOD}}$ decreases over training but stabilize near 0.088. This residual non-zero value can be attributed to the shared low-level visual primitives such as textures and edges between the ID data and the SVHN OOD dataset, which establishes a geometric "floor" for the orthogonality between the two subspaces.

Overall, these results confirm that ResNet-18 trained on CIFAR-100 undergoes Neural Collapse and reaches the terminal training phase.

Also, from our experiments, NECO is a better OOD scoring method when comparing CIFAR-100 and SVHN. The NECO achieved best performance compared for each evaluation metric used. The empirical success of NECO confirms its mathematical soundness and suggests a promising trajectory for future research focused on OOD detection .

References

- [1] V. Pappayan, H. Xue, and D. L. Donoho, *Prevalence of neural collapse during the terminal phase of deep learning training*, Proceedings of the National Academy of Sciences (PNAS), vol. 117, no. 40, pp. 24652–24663, 2020. 5
- [2] M Ben Ammar, N Belkhir, S Popescu, A Manzanera, G Franchi, *NECO: Neural Collapse Based Out-of-Distribution Detection*, International Conference on Learning Representations (ICLR), 2024. 2, 5

Additional Figures

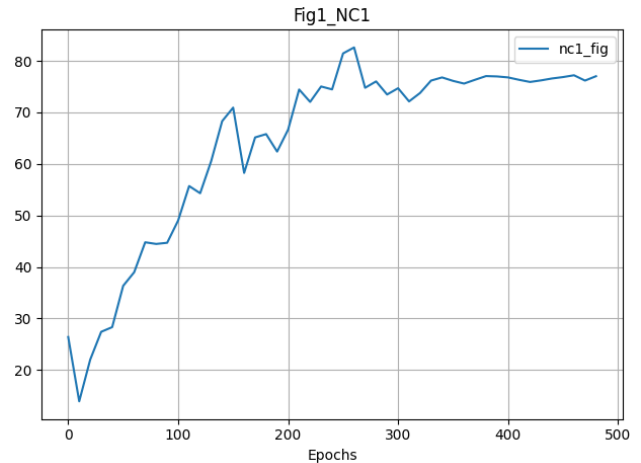


Figure 7: Computation of $\|\Sigma_W\|$ through the epochs. The norm keeps growing, which shows its unreliability as a direct measure of collapse.

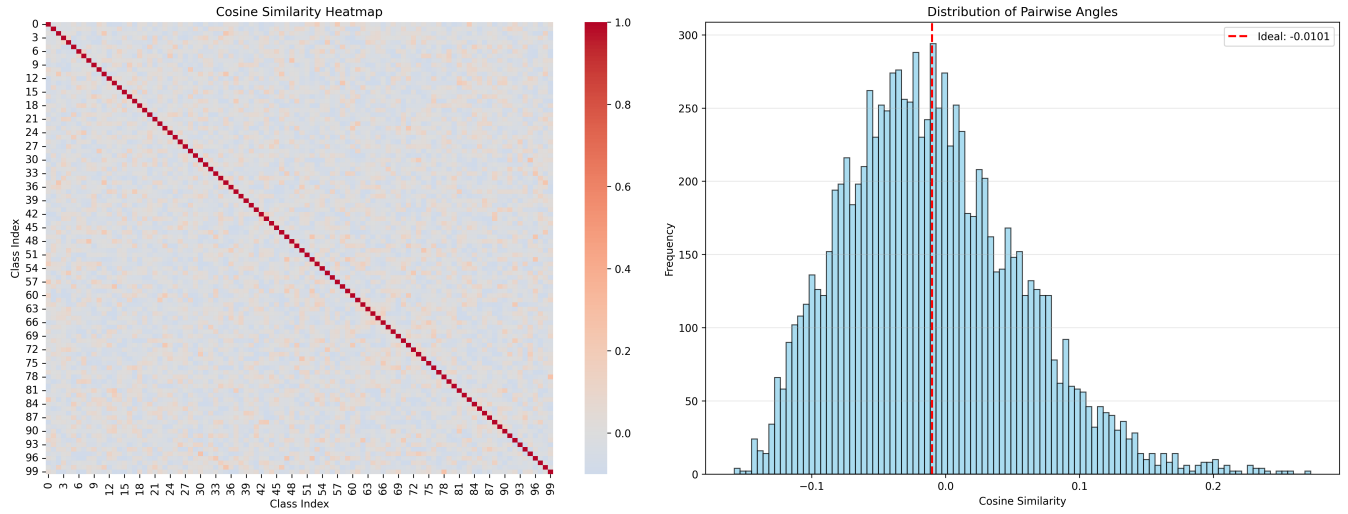


Figure 8: Heatmap and distribution of pairwise angles in the study of NC 2.

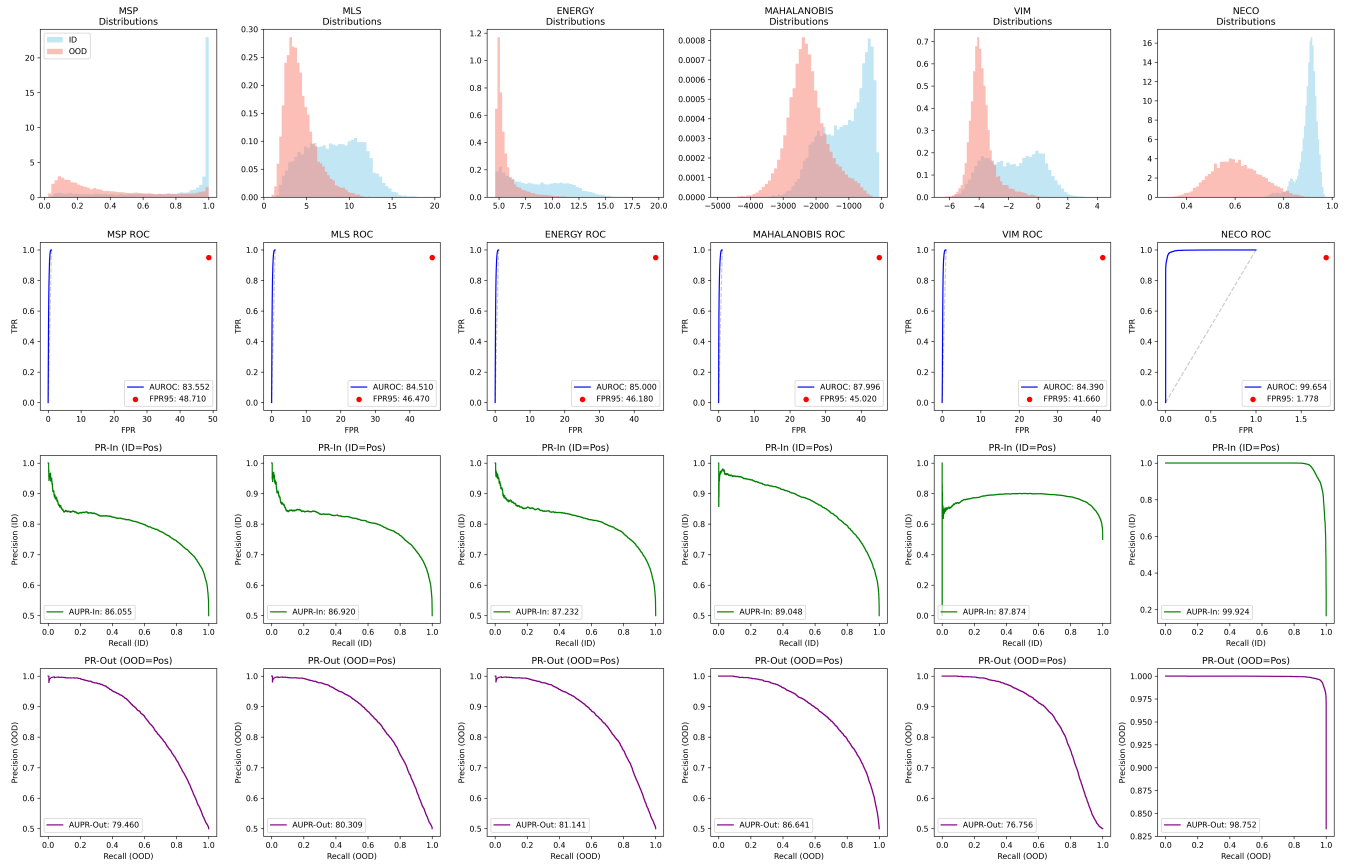


Figure 9: Comprehensive OOD detection benchmark results. The results demonstrate the superior performance of NECO across various In-Distribution and Out-of-Distribution dataset pairs compared to established baselines.