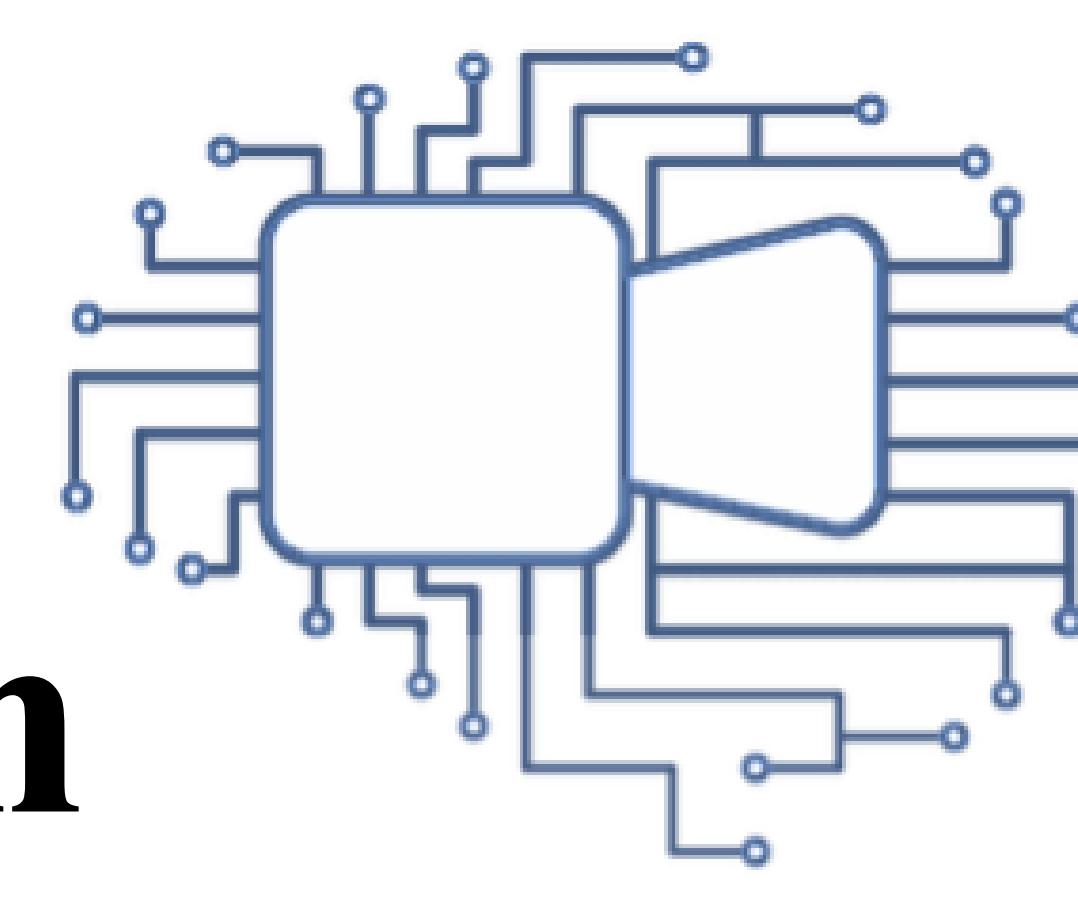


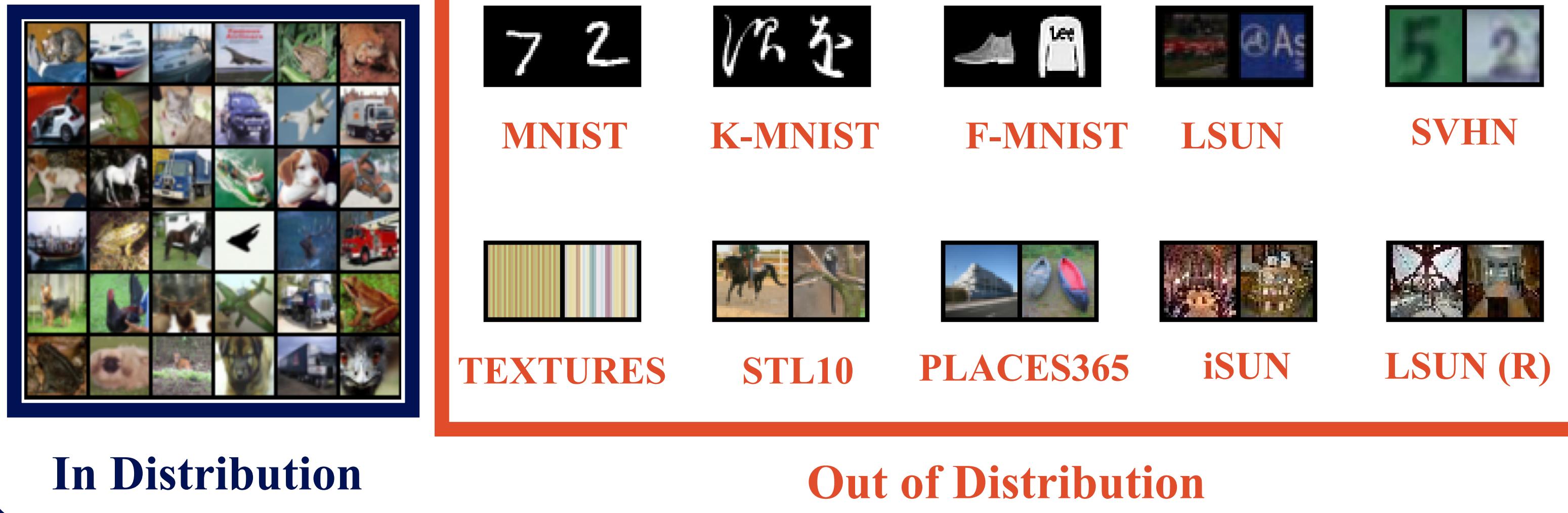
# Learning ODIN

Amir Jevnisek and Shai Avidan



## Problem Statement

**Out-of Distribution Detection Goal** is to identify samples outside of a predetermined distribution, which is defined as In-Distribution. At training time only normal data is available. Area Under ROC Curve (AUROC) and FPR95 are the standard evaluation metrics.



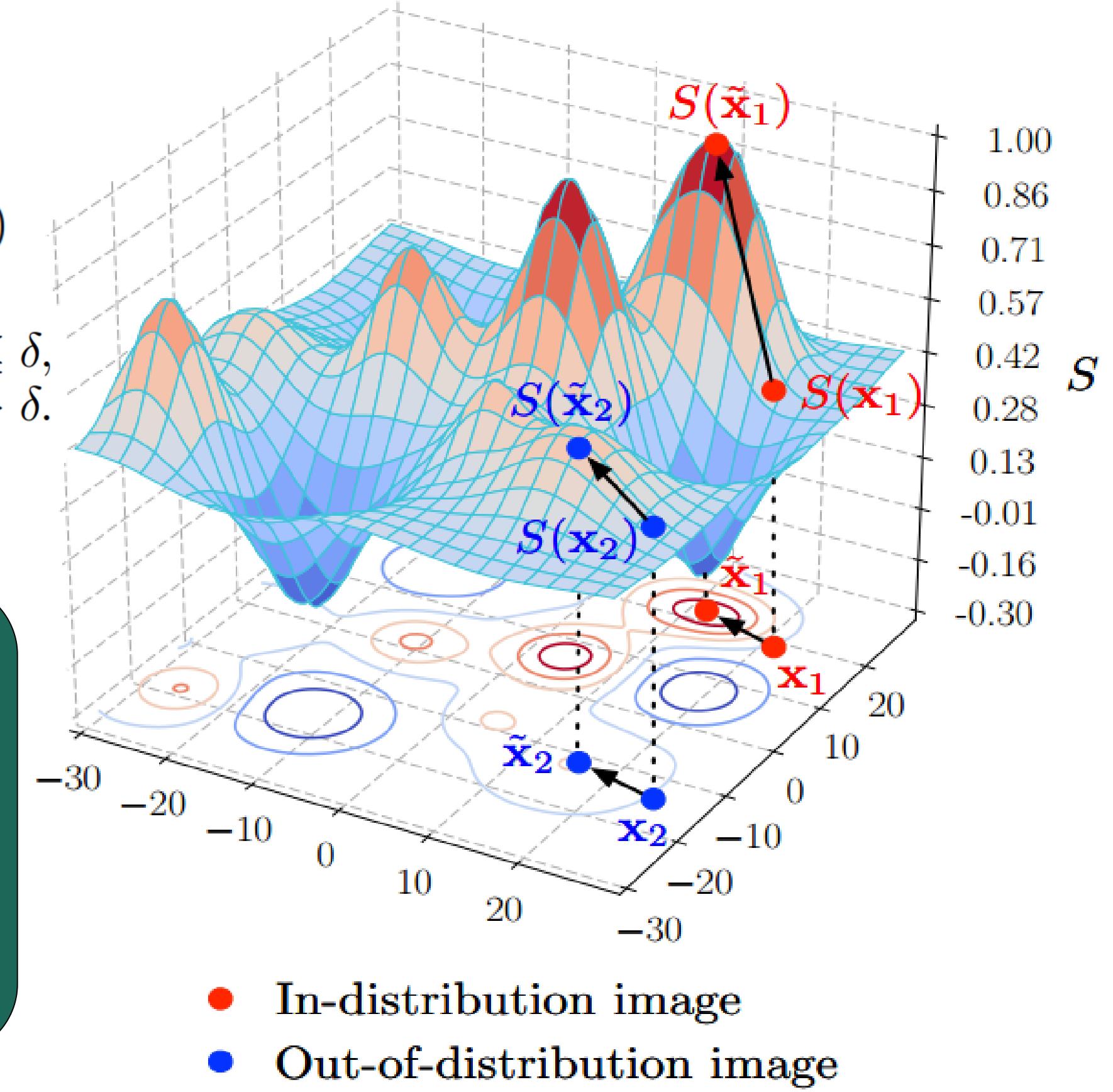
## Introduction

### ODIN OOD Detector:

$$(1) S_i(\mathbf{x}; T) = \frac{\exp(f_i(\mathbf{x})/T)}{\sum_{j=1}^N \exp(f_j(\mathbf{x})/T)},$$

$$(2) \tilde{\mathbf{x}} = \mathbf{x} - \varepsilon \text{sign}(-\nabla_{\mathbf{x}} \log S_{\hat{y}}(\mathbf{x}; T))$$

$$(3) g(\mathbf{x}; \delta, T, \varepsilon) = \begin{cases} 1 & \text{if } \max_i p(\tilde{\mathbf{x}}; T) \leq \delta, \\ 0 & \text{if } \max_i p(\tilde{\mathbf{x}}; T) > \delta. \end{cases}$$



ODIN is a prevailing OOD framework.

ODIN assumes that the logit gap is maximal for In-Distribution samples.

We encode this assumption into the loss function.  
So we learn what ODIN is expecting.

## Implementation

Code is available!



## Our NEW Regularization Term

We Reformalize ODIN's Decision Rule:

$$|S_{\hat{y}}(\tilde{\mathbf{x}}) - S_{\hat{y}}(\mathbf{x})| \stackrel{ID}{\gtrless} Th.$$

We take the Taylor expansion of the softmax with the perturbed image  $S_{\hat{y}}(\tilde{\mathbf{x}})$  around  $\mathbf{x}$ :

$$S_{\hat{y}}(\tilde{\mathbf{x}}) = S_{\hat{y}}(\mathbf{x}) + \varepsilon \cdot \nabla_{\mathbf{x}} S_{\hat{y}}(\mathbf{x}) + \text{higher orders of } \varepsilon$$

Plugging this into the ODIN decision rule, we get:

$$\begin{aligned} |S_{\hat{y}}(\tilde{\mathbf{x}}) - S_{\hat{y}}(\mathbf{x})| &= |S_{\hat{y}}(\mathbf{x}) + \varepsilon \cdot \nabla_{\mathbf{x}} S_{\hat{y}}(\mathbf{x}) - S_{\hat{y}}(\mathbf{x})| = \\ &= |\varepsilon \cdot \nabla_{\mathbf{x}} S_{\hat{y}}(\mathbf{x})| \gtrless Th. \end{aligned}$$

The term 'large' is ill-defined ('large' compared to what?). Instead, we introduce a Gradient Quotient (GQ) loss term:

$$\mathcal{L}_{GQ} = \frac{\|\nabla_{\mathbf{x}} S_{\hat{y}}\|}{\sum_{y_j \neq \hat{y}} \|\nabla_{\mathbf{x}} S_{y_j}\|}$$

This way, the term 'large' refers to the softmax scores of the rest of the class labels.

We suggest to train networks with this as a regularization term:

$$L = L_{CE}(y, \hat{y}) + \lambda \cdot \frac{\|\nabla_{\mathbf{x}} S_y\|_1}{\sum_{y_j \neq y} \|\nabla_{\mathbf{x}} S_{y_j}\|_1}$$

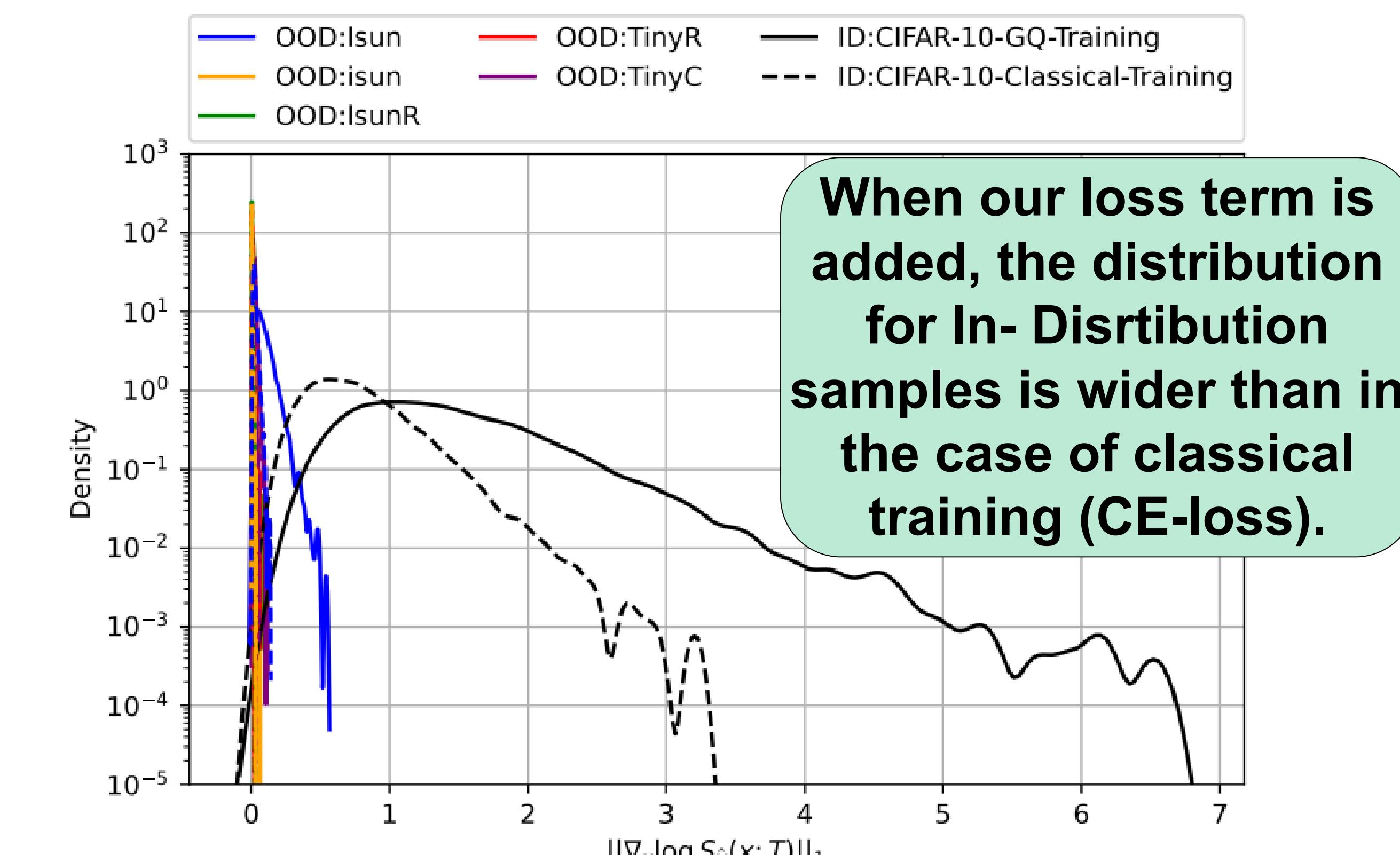
## OOD Algorithm

### Algorithm 1 Learning ODIN

**Inputs** (1) In-Distribution dataset  $\mathcal{D}_{ID}^{train}$ , (2) classifier  $f_{\theta}(\cdot)$  with parameters  $\theta$   
**Outputs:** (1) An OOD detector, (2) An In-Distribution classifier.

- 1: Train  $f_{\theta}(\cdot)$  with the following loss:  $L = L_{CE} + \lambda \cdot \mathcal{L}_{GQ}$
- 2: Create  $g(x)$  an ODIN OOD-detector on  $f_{\theta}$ 's logits.  
**return**  $g(x)$  is the OOD-detector,  $f_{\theta}$  is the In-Distribution classifier

## Gradient Distribution



## OOD Performance

Architecture	Method	AUROC ↑ (%)	FPR95 ↓ (%)	ID-Accuracy ↑ (%)
WideResNet-40-4	MSP	88.99	56.81	94.93
	ODIN	90.11	35.31	94.93
	Mahalanobis	89.33	35.48	94.93
	Energy	90.04	35.26	94.93
MSDNet Exit@last	MSP	89.72	49.87	94.09
	ODIN	90.33	29.30	94.09
	Mahalanobis	82.84	75.19	94.09
	Energy	90.48	33.62	94.09
MSDNet (dynamic exit)	MOOD	91.26	28.05	94.13
	Madry's Resnet /+ GQ	90.51 / 91.94	36.57 / 30.19	92.19 / 93.29
	ResNet18 /+ GQ	89.09 / 93.55	27.75 / 25.14	95.27 / 94.26
ResNet34 /+ GQ	ODIN	89.03 / 89.08	35.33 / 27.63	95.21 / 95.56

**CIFAR-10 as ID:** Adding GQ to any of the architectures improves results. Our approach achieves best results for all metrics compared to baseline models.

## High Res. Images

GQ also works for **High Resolution** images (ImageNet).

It is possible to use GQ for **fine tuning**, which lets us work with **pre-trained models**.

