



Spurious-Aware Prototype Refinement for Reliable Out-of-Distribution Detection

Reihaneh Zohrabi*, Hosein Hasani*, Mahdieh Soleymani, Anna
Rohrbach, Marcus Rohrbach, Mohammad Hossein Rohban

TU Darmstadt, Sharif University of Technology

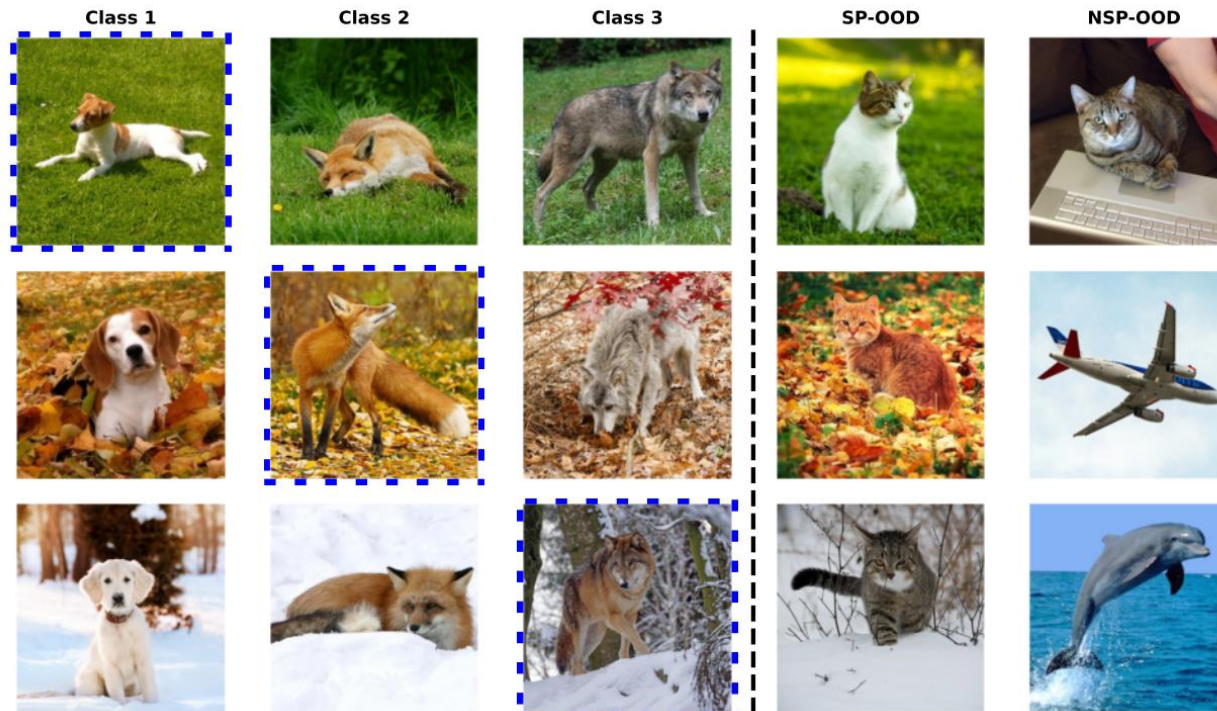
NeurIPS 2025



TECHNISCHE
UNIVERSITÄT
DARMSTADT

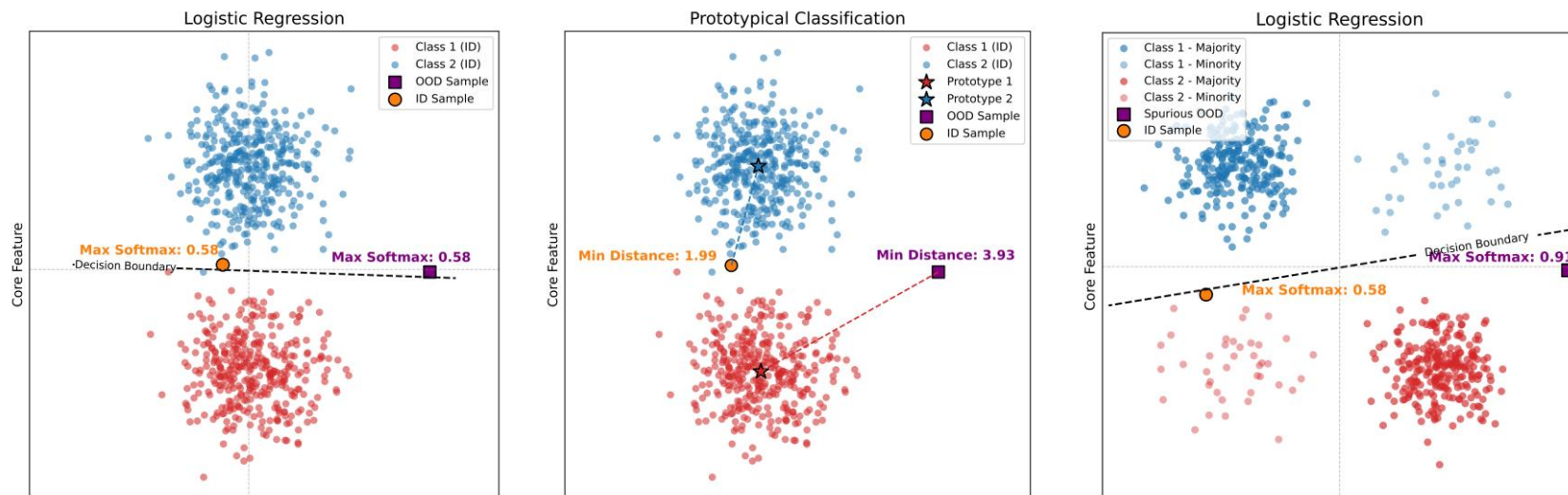
Problem Definition & Background

- **OOD Detection:** Identify inputs outside the training distribution.
- **Spurious Features:** Irrelevant cues learned by models.
- **SP-OOD:** OOD samples sharing spurious cues with in-distribution data.

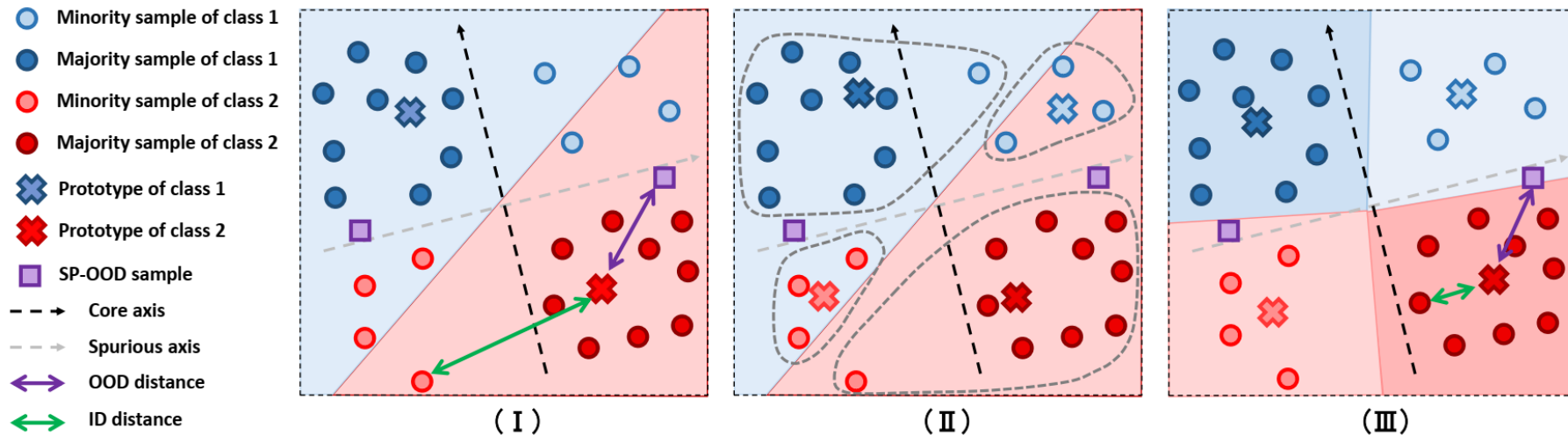


Motivation: Score Calculation

- **Score Calculation:**
a scalar score $S(x)$ that separates ID and OOD samples.
- Two Perspectives:
 - **Discriminative:** defining score based on $p(y|x)$; e.g. softmax probabilities.
 - **Generative:** defining score based on $p(x|y)$; e.g. distance from prototypes.
- Discriminative approach has certain **failure cases**.



SPROD Overview



- **Stage 1 – Initial Prototypes:**
Compute class prototypes from training features.
- **Stage 2 – Classification-Aware Prototypes:**
Separate correctly and incorrectly classified samples to capture feature bias.
- **Stage 3 – Group Prototype Refinement:**
Reassign samples and refine prototypes to reduce spurious influence.

Experimental Setup & Datasets

- **Baselines:**
Compared with **19** post-hoc OOD detection methods.
- **Backbones:**
Primarily ResNet-50; **10** different backbones in total.
- **Metrics:** AUROC and FPR@95, AUPR.
- **SP-OOD Benchmarks:**
Waterbirds, CelebA, UrbanCars, Spurious ImageNet, and Animals MetaCoCo.

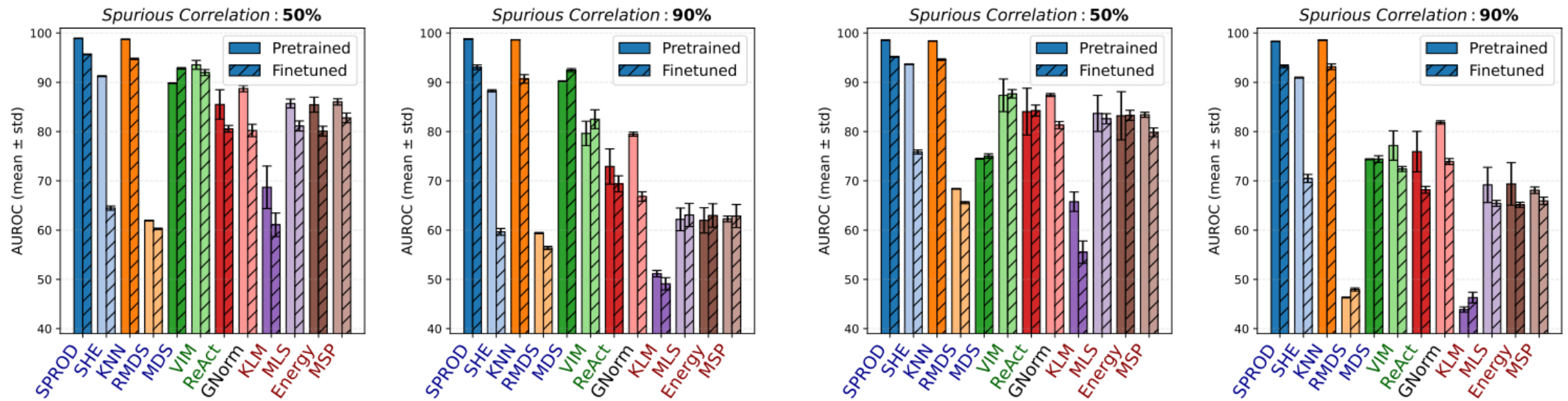
Results: SP-OOD benchmarks

- Comparative performance of **post-hoc OOD detection methods** on SP-OOD benchmarks using a **ResNet-50** backbone.

AUROC↑							FPR@95↓						
Method	WB	CA	UC	AMC	SpI	Avg.	Method	WB	CA	UC	AMC	SpI	Avg.
MSP ^[7]	62.3 \pm 0.6	46.0 \pm 1.4	38.5 \pm 0.3	79.7 \pm 0.4	83.1 \pm 0.3	61.9	MSP ^[7]	87.9 \pm 0.8	98.7 \pm 0.5	97.3 \pm 0.3	83.8 \pm 0.7	74.1 \pm 1.2	88.4
Energy ^[13]	62.0 \pm 2.6	45.4 \pm 3.4	38.4 \pm 2.1	79.9 \pm 0.6	80.6 \pm 0.4	61.3	Energy ^[13]	89.2 \pm 3.2	98.6 \pm 0.7	95.5 \pm 3.1	84.8 \pm 0.8	76.3 \pm 0.9	88.9
MLS ^[44]	62.2 \pm 2.3	45.3 \pm 3.2	38.4 \pm 1.4	80.2 \pm 0.6	81.9 \pm 0.3	61.6	MLS ^[44]	88.1 \pm 2.0	98.8 \pm 0.6	96.7 \pm 2.0	84.4 \pm 0.8	74.6 \pm 0.9	88.5
KLM ^[44]	51.2 \pm 0.7	41.7 \pm 2.5	57.0 \pm 0.2	74.2 \pm 0.6	79.6 \pm 0.8	60.7	KLM ^[44]	89.1 \pm 0.7	98.7 \pm 0.5	97.1 \pm 0.3	80.5 \pm 0.8	76.1 \pm 1.7	88.3
GEN ^[45]	62.3 \pm 0.6	46.0 \pm 1.4	38.5 \pm 0.3	80.2 \pm 0.0	80.8 \pm 0.4	61.6	GEN ^[45]	87.9 \pm 0.8	98.7 \pm 0.5	97.3 \pm 0.3	84.8 \pm 0.1	76.3 \pm 0.7	89.0
GNorm ^[61]	79.5 \pm 0.4	38.0 \pm 1.3	46.6 \pm 0.4	74.2 \pm 0.5	85.2 \pm 0.2	64.7	GNorm ^[61]	84.2 \pm 0.7	98.8 \pm 0.4	97.1 \pm 0.1	84.2 \pm 0.6	54.7 \pm 0.6	83.8
ReAct ^[58]	72.9 \pm 3.6	45.6 \pm 5.3	41.3 \pm 3.1	80.1 \pm 0.6	83.6 \pm 0.7	64.7	ReAct ^[58]	86.9 \pm 7.0	96.3 \pm 2.4	95.5 \pm 3.2	83.9 \pm 0.8	57.5 \pm 1.6	84.0
VIM ^[59]	79.6 \pm 2.5	50.4 \pm 3.1	60.7 \pm 1.7	78.6 \pm 0.6	77.4 \pm 0.9	69.3	VIM ^[59]	61.4 \pm 3.5	96.2 \pm 0.4	69.0 \pm 1.5	86.6 \pm 0.7	79.5 \pm 0.5	78.5
ASH ^[60]	78.5 \pm 3.2	47.3 \pm 2.8	39.6 \pm 1.7	78.0 \pm 0.2	86.6 \pm 0.7	66.0	ASH ^[60]	85.2 \pm 7.0	96.9 \pm 1.4	96.1 \pm 1.5	87.9 \pm 0.4	52.9 \pm 3.1	83.8
MDS ^[11]	90.2 \pm 0.1	57.8 \pm 0.5	91.8 \pm 0.1	62.9 \pm 0.8	58.4 \pm 0.1	72.2	MDS ^[11]	49.2 \pm 0.2	96.0 \pm 0.5	39.0 \pm 0.3	93.0 \pm 0.3	90.5 \pm 0.1	73.5
RMDS ^[46]	59.4 \pm 0.1	33.6 \pm 1.4	47.4 \pm 0.2	81.9 \pm 0.4	68.8 \pm 0.1	58.2	RMDS ^[46]	91.7 \pm 0.2	99.6 \pm 0.1	95.3 \pm 0.1	83.4 \pm 0.9	88.1 \pm 0.1	91.6
KNN ^[47]	98.6 \pm 0.0	54.5 \pm 0.5	91.1 \pm 0.1	79.7 \pm 0.6	77.4 \pm 0.0	80.3	KNN ^[47]	4.8 \pm 0.1	94.4 \pm 1.0	42.5 \pm 0.2	79.9 \pm 1.1	70.4 \pm 0.2	58.4
SHE ^[48]	88.3 \pm 0.2	42.7 \pm 0.6	73.2 \pm 0.1	54.8 \pm 0.7	83.0 \pm 0.1	68.4	SHE ^[48]	33.2 \pm 0.5	96.4 \pm 0.5	76.5 \pm 0.2	93.9 \pm 0.3	52.6 \pm 0.8	70.5
NECO ^[51]	53.5 \pm 1.6	39.5 \pm 3.2	35.1 \pm 1.5	80.2 \pm 0.1	67.2 \pm 0.3	55.1	NECO ^[51]	90.5 \pm 2.2	98.8 \pm 0.6	96.7 \pm 1.9	78.2 \pm 0.0	89.9 \pm 0.8	90.8
NNGuide ^[49]	70.6 \pm 2.9	49.8 \pm 4.2	43.6 \pm 2.1	79.4 \pm 0.0	85.1 \pm 0.8	65.7	NNGuide ^[49]	77.7 \pm 6.0	97.6 \pm 1.2	91.6 \pm 3.6	86.3 \pm 0.1	52.2 \pm 1.4	81.1
Relation ^[50]	80.7 \pm 0.2	60.4 \pm 2.5	96.0 \pm 0.5	74.5 \pm 0.3	81.8 \pm 0.7	78.7	Relation ^[50]	73.8 \pm 0.5	95.4 \pm 0.2	24.2 \pm 1.7	84.6 \pm 0.2	78.0 \pm 1.1	71.2
SCALE ^[52]	89.0 \pm 2.9	44.9 \pm 3.2	54.4 \pm 2.1	78.4 \pm 0.4	86.2 \pm 0.5	70.6	SCALE ^[52]	61.0 \pm 22.5	98.7 \pm 0.6	94.6 \pm 3.1	87.7 \pm 0.3	53.0 \pm 1.5	79.0
fDBD ^[53]	71.1 \pm 0.5	51.3 \pm 1.3	47.4 \pm 0.2	79.9 \pm 0.0	84.2 \pm 0.3	66.8	fDBD ^[53]	85.5 \pm 0.8	98.6 \pm 0.5	96.1 \pm 0.4	85.4 \pm 0.2	70.8 \pm 1.0	87.3
NCI ^[54]	84.0 \pm 0.1	46.4 \pm 2.4	54.8 \pm 0.8	78.5 \pm 0.1	84.9 \pm 0.2	69.7	NCI ^[54]	41.1 \pm 0.1	99.4 \pm 0.3	92.2 \pm 0.6	85.8 \pm 0.3	63.8 \pm 0.9	76.5
SPROD	98.8\pm0.0	61.6\pm0.9	97.4\pm0.0	82.4\pm0.5	85.3\pm0.0	85.1	SPROD	4.7\pm0.1	93.7\pm0.9	19.0\pm0.4	69.5\pm1.2	58.0\pm0.1	49.0

Results: Fine-Tuning and SP Correlation Rate

- **Fine-tuning** the backbone **reduces OOD performance**, especially for feature-based methods.
- **Higher spurious correlation (90% vs. 50%)** further degrades performance, mainly for output-based methods.
- **Smaller backbone (ResNet-18)** performs comparably to ResNet-50.

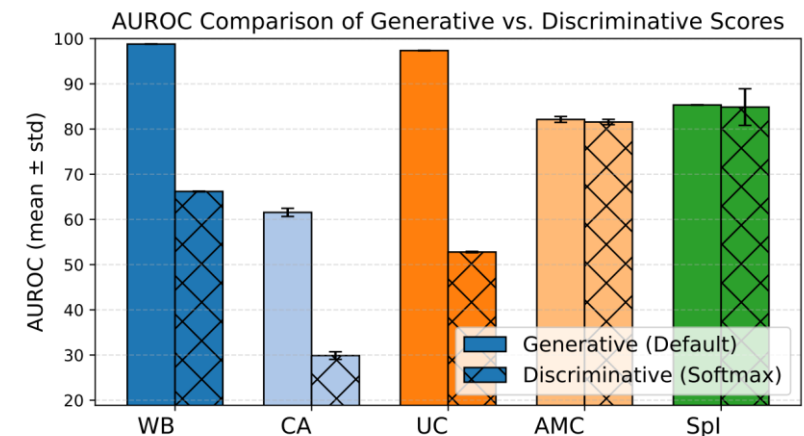
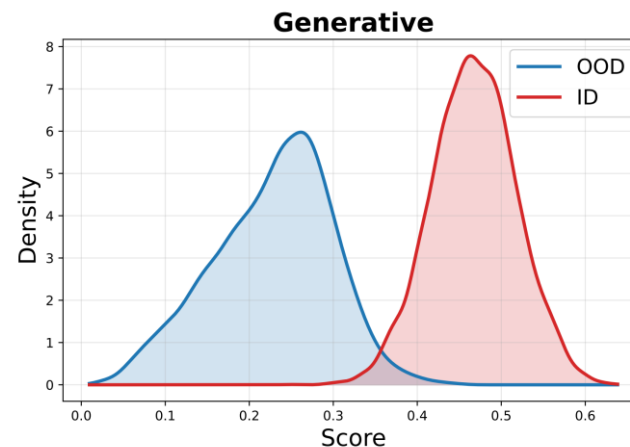
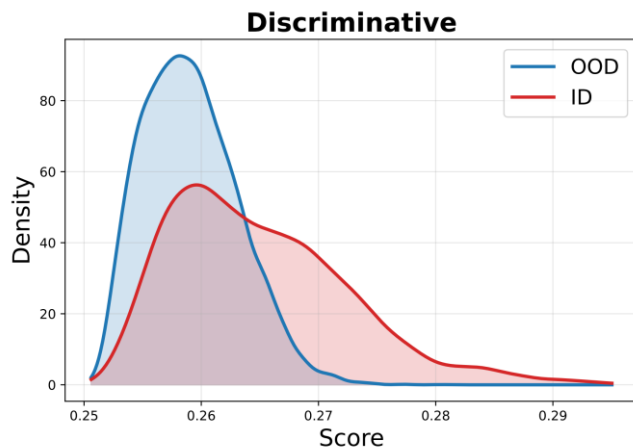


(a) ResNet-50 backbone.

(b) ResNet-18 backbone.

Results: Generative & Discriminative (Ablation)

- **Distance-based vs. softmax-based** scoring, using identical SP-ROD features and prototypes.
- **Generative** scores yield **clearer ID/OOD separation**.
- **Generative** scoring **outperforms** on most SP-OOD datasets. Biggest gap on Waterbirds, CelebA, and UrbanCars.



Results: Low-Shot and Zero-Shot SP-OOD

- **Zero-Shot Comparison (CLIP-based):**
 - MCM: 98.36 CMA: 98.62 SP-PROD: 99.01 (vision-only)
- Text-free SP-PROD outperforms CLIP-based zero-shot methods.
- **Low-Shot Comparison:**
 - SP-PROD maintains high AUROC even with limited samples, unlike KNN.

