# Outlier Synthesis via Hamiltonian Monte Carlo for Out-of-Distribution Detection

Hengzhuang Li, Teng Zhang

School of Computer Science and Technology,

Huazhong University of Science and Technology, Wuhan, China

# Contents

- Introduction
- Motivation
- Method
- Results
- Conclusion

# Introduction

- **Task:** fine-tuning based out-of-distribution (OOD) detection without access to auxiliary OOD dataset

- **Outlier synthesis:** synthesize virtual outliers which serve as surrogated OOD supervision signals

- **Results:** our proposed framework *HamOS* synthesizes high quality outliers and outperforms previous baselines

# Motivation

- Pixel space outlier synthesis
  - E.g., Dream-OOD [1]
  - To generate pixel space outliers through the generative models, e.g., diffusion model.

- Feature space outlier synthesis
  - E.g., VOS [2], NPOS [3]
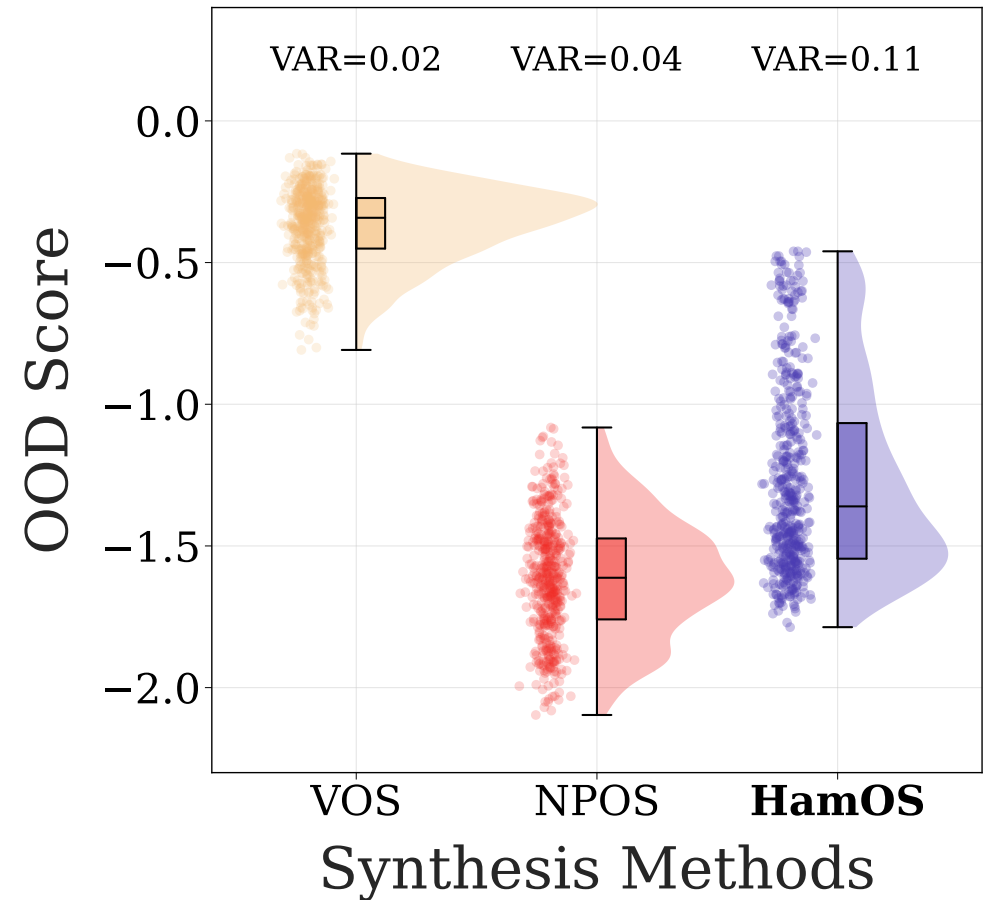  - To generate outliers in the feature space through sampling algorithms, e.g. Gaussian.

[1] Du, et al. Dream the impossible: Outlier imagination with diffusion models. In NIPS, 2023.
[2] Du, et al. Vos: Learning what you don't know by virtual outlier synthesis. In ICLR, 2022.
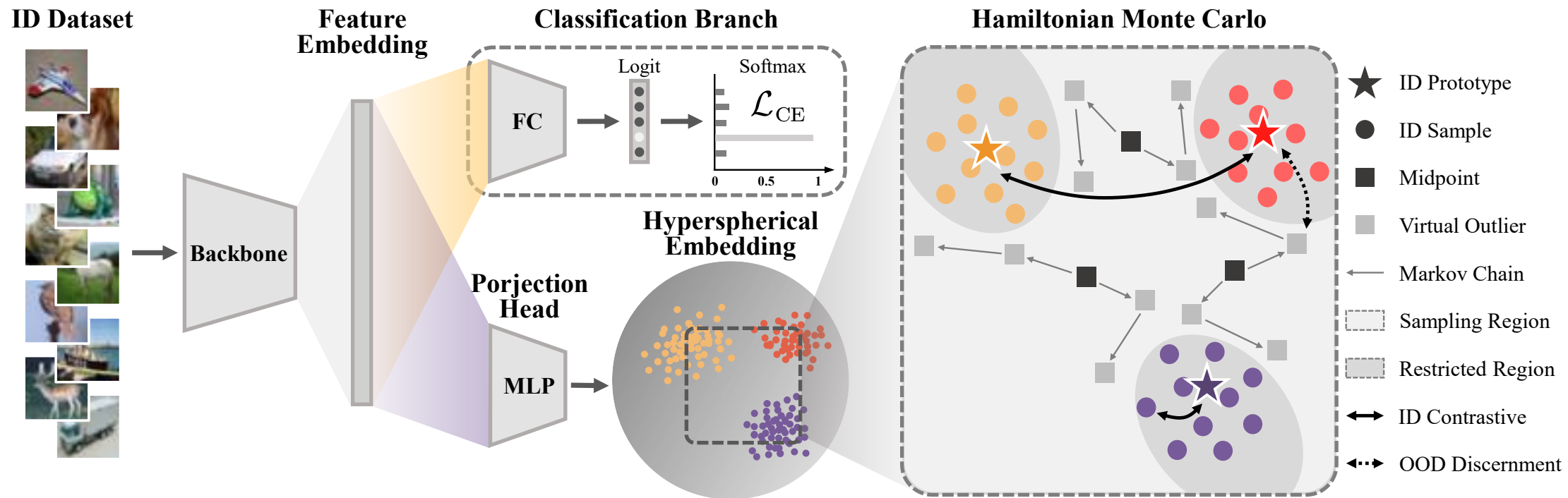[3] Tao, et al. Non-parametric outlier synthesis. In ICLR, 2023.

# Motivation

- ### Our goal
  - To efficiently synthesize dicers and representative outliers based solely on the ID data

- ### Idea
  - Modeling the synthesis process as Markov chain



OOD score distributions

# Method

# Method

- ## Synthesizing Outliers via Hamiltonian Monte Carlo (HMC)
  - ### Estimating OOD density via the distance to the k-th nearest neighbor: we design a quantitative characterizing of the likelihood that a sample is OOD rather than ID.

$$P^{\mathrm{OOD}}(\boldsymbol{z};\, \mathcal{Z}_c) = \|\boldsymbol{z} - \boldsymbol{z}_{c(k)}\|_2 \qquad P^{\mathrm{OOD}}(\boldsymbol{z};\, \mathcal{Z}_u, \mathcal{Z}_v) = \frac{P^{\mathrm{OOD}}(\boldsymbol{z};\, \mathcal{Z}_u) + P^{\mathrm{OOD}}(\boldsymbol{z};\, \mathcal{Z}_v)}{2}$$

$$U^{\mathrm{OOD}}(\boldsymbol{z};\, \mathcal{Z}_u, \mathcal{Z}_v) = -\log P^{\mathrm{OOD}}(\boldsymbol{z};\, \mathcal{Z}_u, \mathcal{Z}_v) = -\log \sum_{i=u,v} P^{\mathrm{OOD}}(\boldsymbol{z};\, \mathcal{Z}_i) + \mathrm{constant}$$

  - ### Synthesizing outliers by OOD density estimation via HMC: we generate virtual outliers along the Markov chains by solving the Hamilton's Equation.

$$H(\boldsymbol{z}, \boldsymbol{q}) = U^{\mathrm{OOD}}(\boldsymbol{z}) + \frac{1}{2}\|\boldsymbol{q}\|_2^2$$

  - ### Rejecting erroneous outliers located within ID clusters: we reject false outliers that conflate with ID embeddings by applying a hard margin according to the ID probability.

$$t_- = -\log\max_c P_c^{\mathrm{ID}}(\mathbf{b}_{u,v}) - \delta \quad \mathbf{b}_{u,v} = \frac{\boldsymbol{\mu}_u + \boldsymbol{\mu}_v}{\|\boldsymbol{\mu}_u + \boldsymbol{\mu}_v\|_2}$$
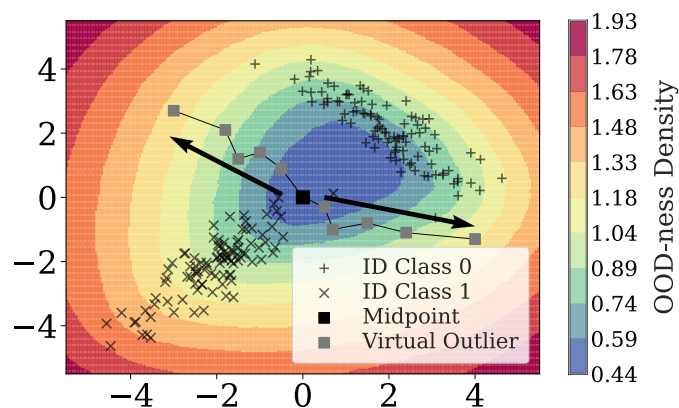
# Method

- ## Training with Synthesized Outliers
  - We fine-tune the model with the OOD discernment loss, the contrastive loss, and the cross-entropy loss to help broaden the gap between ID and OOD data.
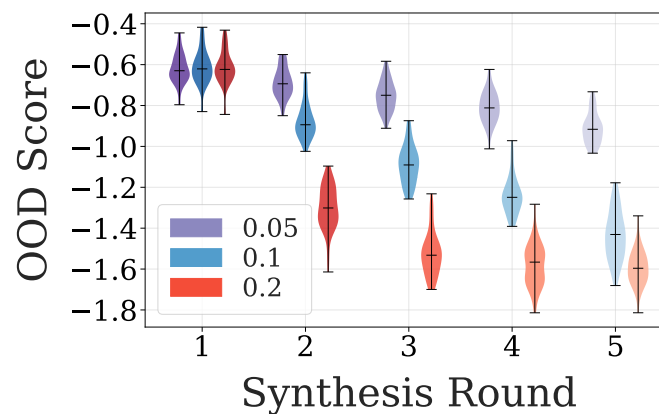
$$\mathcal{L}_{\mathrm{OOD-disc}} = \frac{1}{M}\sum_{i=1}^{M}\frac{1}{C}\sum_{j=1}^{C}\log\frac{\exp(\boldsymbol{z}_i^\top\boldsymbol{\mu}_j/\tau)}{\displaystyle\sum_{l=1}^{C}\exp(\boldsymbol{z}_i^\top\boldsymbol{\mu}_l/\tau)}$$

$$\mathcal{L}_{\mathrm{HamOS}} = \mathcal{L}_{\mathrm{CE}} + \mathcal{L}_{\mathrm{ID-con}} + \lambda_d\,\mathcal{L}_{\mathrm{OOD-disc}}$$
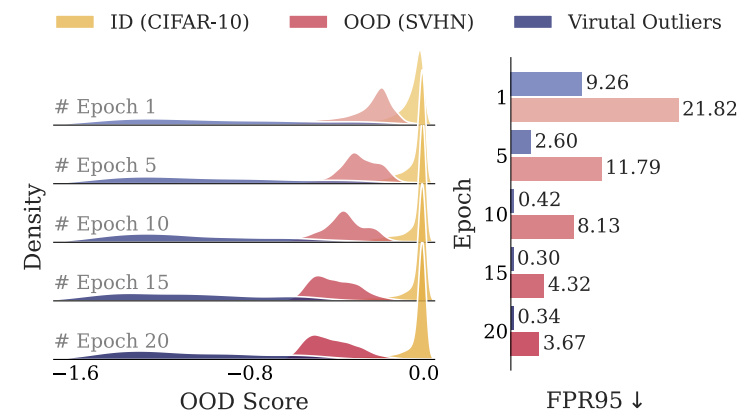
# Method



Depiction of the designed OOD-ness density estimation.

Varied OOD scores of the generated outliers at different synthesis rounds.

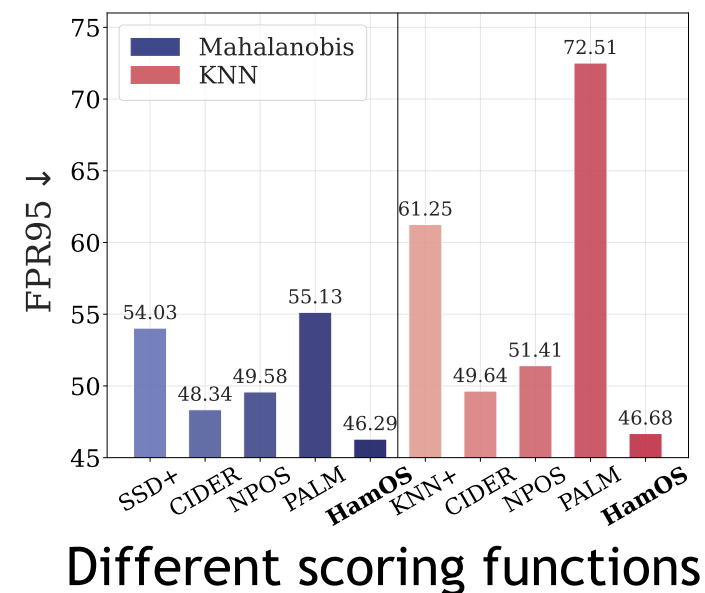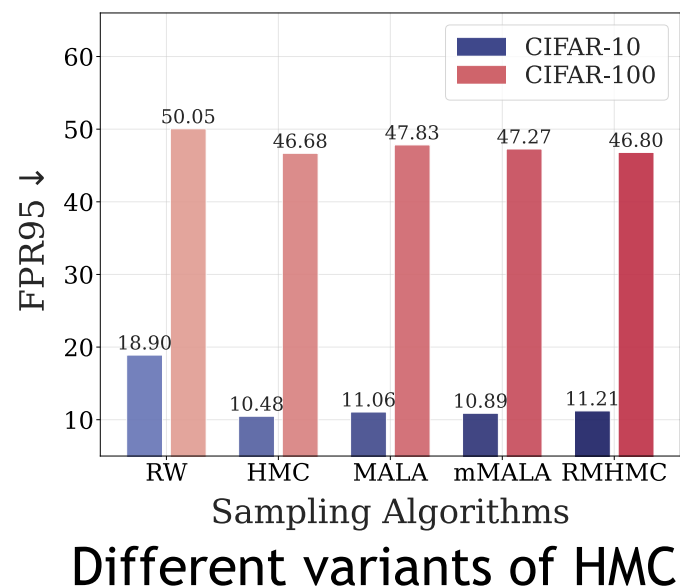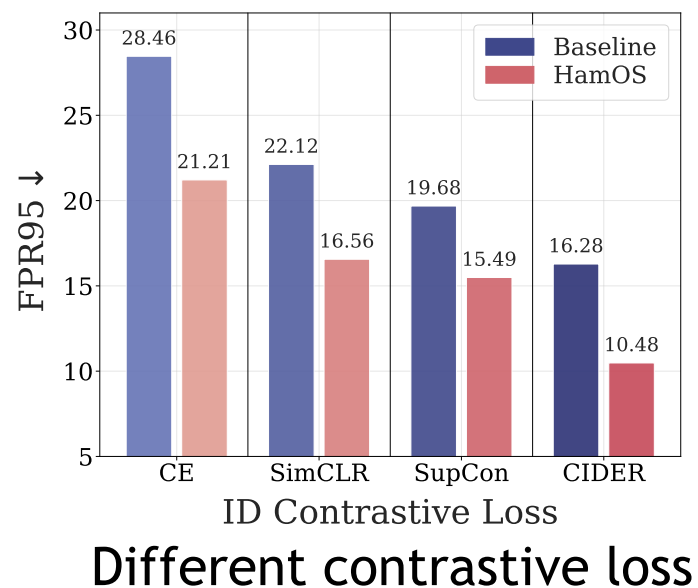OOD performance is improved continuously along the training process.

# Results

## Main results: CIFAR10/100 benchmarks

| Methods | CIFAR-10 | | | | CIFAR-100 | | | |
|---|---|---|---|---|---|---|---|---|
| | FPR95↓ | AUROC↑ | AUPR↑ | ID-ACC↑ | FPR95↓ | AUROC↑ | AUPR↑ | ID-ACC↑ |
| *Post-hoc Methods* | | | | | | | | |
| MSP | $32.17_{\pm6.38}$ | $91.10_{\pm0.71}$ | $81.70_{\pm5.82}$ | $\mathbf{95.17}_{\pm0.16}$ | $59.78_{\pm2.16}$ | $77.25_{\pm1.28}$ | $66.86_{\pm1.58}$ | $\mathbf{76.69}_{\pm0.24}$ |
| ODIN | $58.04_{\pm18.46}$ | $85.70_{\pm4.17}$ | $70.08_{\pm11.84}$ | $\mathbf{95.17}_{\pm0.16}$ | $63.49_{\pm2.51}$ | $78.01_{\pm1.62}$ | $65.20_{\pm2.19}$ | $\mathbf{76.69}_{\pm0.25}$ |
| EBO | $41.85_{\pm13.78}$ | $91.79_{\pm1.54}$ | $79.70_{\pm8.10}$ | $\mathbf{95.17}_{\pm0.16}$ | $60.86_{\pm1.87}$ | $78.32_{\pm1.31}$ | $66.73_{\pm1.35}$ | $\mathbf{76.69}_{\pm0.24}$ |
| KNN | $22.86_{\pm1.12}$ | $92.98_{\pm0.42}$ | $88.74_{\pm0.79}$ | $\mathbf{95.17}_{\pm0.16}$ | $56.96_{\pm2.96}$ | $81.01_{\pm1.19}$ | $70.60_{\pm2.29}$ | $\mathbf{76.69}_{\pm0.24}$ |
| ASH | $54.22_{\pm26.06}$ | $87.37_{\pm6.60}$ | $72.33_{\pm16.40}$ | $95.10_{\pm0.14}$ | $66.84_{\pm0.87}$ | $77.14_{\pm1.12}$ | $62.24_{\pm0.73}$ | $76.20_{\pm0.23}$ |
| Scale | $63.18_{\pm23.64}$ | $77.74_{\pm16.24}$ | $63.03_{\pm20.52}$ | $95.15_{\pm0.16}$ | $69.27_{\pm2.31}$ | $77.25_{\pm1.01}$ | $61.42_{\pm1.42}$ | $\mathbf{76.69}_{\pm0.24}$ |
| Relation | $26.28_{\pm1.63}$ | $92.31_{\pm0.43}$ | $86.75_{\pm0.98}$ | $\mathbf{95.17}_{\pm0.16}$ | $59.64_{\pm2.48}$ | $79.69_{\pm1.08}$ | $68.76_{\pm1.78}$ | $\mathbf{76.69}_{\pm0.24}$ |
| *Regularization-based Methods* | | | | | | | | |
| CSI | $21.21_{\pm1.68}$ | $93.73_{\pm0.33}$ | $89.74_{\pm0.68}$ | $92.03_{\pm0.72}$ | $69.34_{\pm0.86}$ | $73.46_{\pm0.37}$ | $61.57_{\pm0.75}$ | $61.75_{\pm0.15}$ |
| SSD+ | $18.49_{\pm1.20}$ | $94.85_{\pm0.57}$ | $90.88_{\pm0.83}$ | $93.95_{\pm0.57}$ | $54.03_{\pm1.92}$ | $80.64_{\pm0.60}$ | $69.73_{\pm1.09}$ | $75.63_{\pm0.39}$ |
| KNN+ | $19.68_{\pm1.86}$ | $94.41_{\pm0.66}$ | $90.46_{\pm0.66}$ | $93.79_{\pm0.63}$ | $61.25_{\pm0.81}$ | $78.24_{\pm0.93}$ | $66.64_{\pm0.88}$ | $72.18_{\pm0.58}$ |
| VOS | $42.37_{\pm21.13}$ | $91.42_{\pm3.38}$ | $79.16_{\pm11.62}$ | $95.05_{\pm0.05}$ | $58.55_{\pm1.53}$ | $81.40_{\pm0.62}$ | $68.33_{\pm1.61}$ | $74.71_{\pm0.07}$ |
| CIDER | $16.28_{\pm0.68}$ | $95.76_{\pm0.37}$ | $92.36_{\pm0.06}$ | $93.98_{\pm0.16}$ | $49.64_{\pm1.80}$ | $81.77_{\pm0.95}$ | $73.22_{\pm1.12}$ | $75.09_{\pm0.49}$ |
| NPOS | $14.39_{\pm0.87}$ | $96.61_{\pm0.26}$ | $93.35_{\pm0.74}$ | $93.95_{\pm0.13}$ | $51.41_{\pm1.88}$ | $81.02_{\pm0.98}$ | $72.49_{\pm1.54}$ | $74.53_{\pm0.62}$ |
| PALM | $32.25_{\pm4.14}$ | $90.54_{\pm1.46}$ | $84.44_{\pm2.14}$ | $93.93_{\pm0.98}$ | $55.13_{\pm0.97}$ | $79.95_{\pm1.26}$ | $70.21_{\pm1.38}$ | $74.67_{\pm0.36}$ |
| **HamOS(ours)** | $\mathbf{10.48}_{\pm0.76}$ | $\mathbf{97.11}_{\pm0.26}$ | $\mathbf{94.94}_{\pm0.86}$ | $94.67_{\pm0.15}$ | $\mathbf{46.68}_{\pm1.44}$ | $\mathbf{83.64}_{\pm0.64}$ | $\mathbf{75.52}_{\pm1.30}$ | $76.12_{\pm0.14}$ |

# Results

## Ablation study



Different contrastive loss

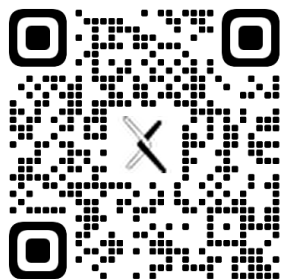Different variants of HMC

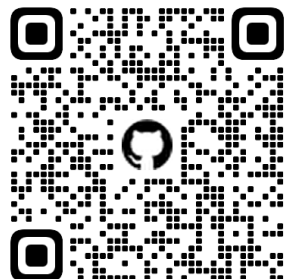Different scoring functions

# Conclusion

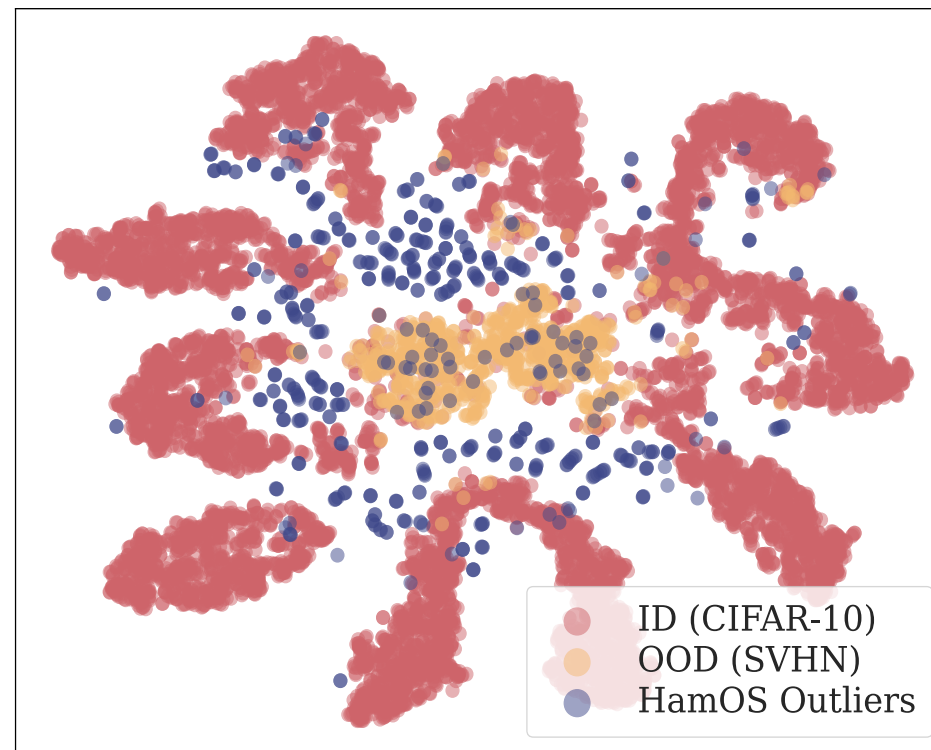- We propose a novel framework HamOS to synthesize virtual outliers for OOD detection

## *Please Contact Us !*

arXiv

Github

WeChat



Feature visualization via t-SNE