

Know Your Space: Inlier and Outlier Construction for Calibrating Medical OOD Detectors



Vivek
Narayanaswamy
LLNL



Yamen
Mubarka
LLNL



Rushil
Anirudh
LLNL



Deepta
Rajan
Microsoft



Andreas
Spanias
ASU



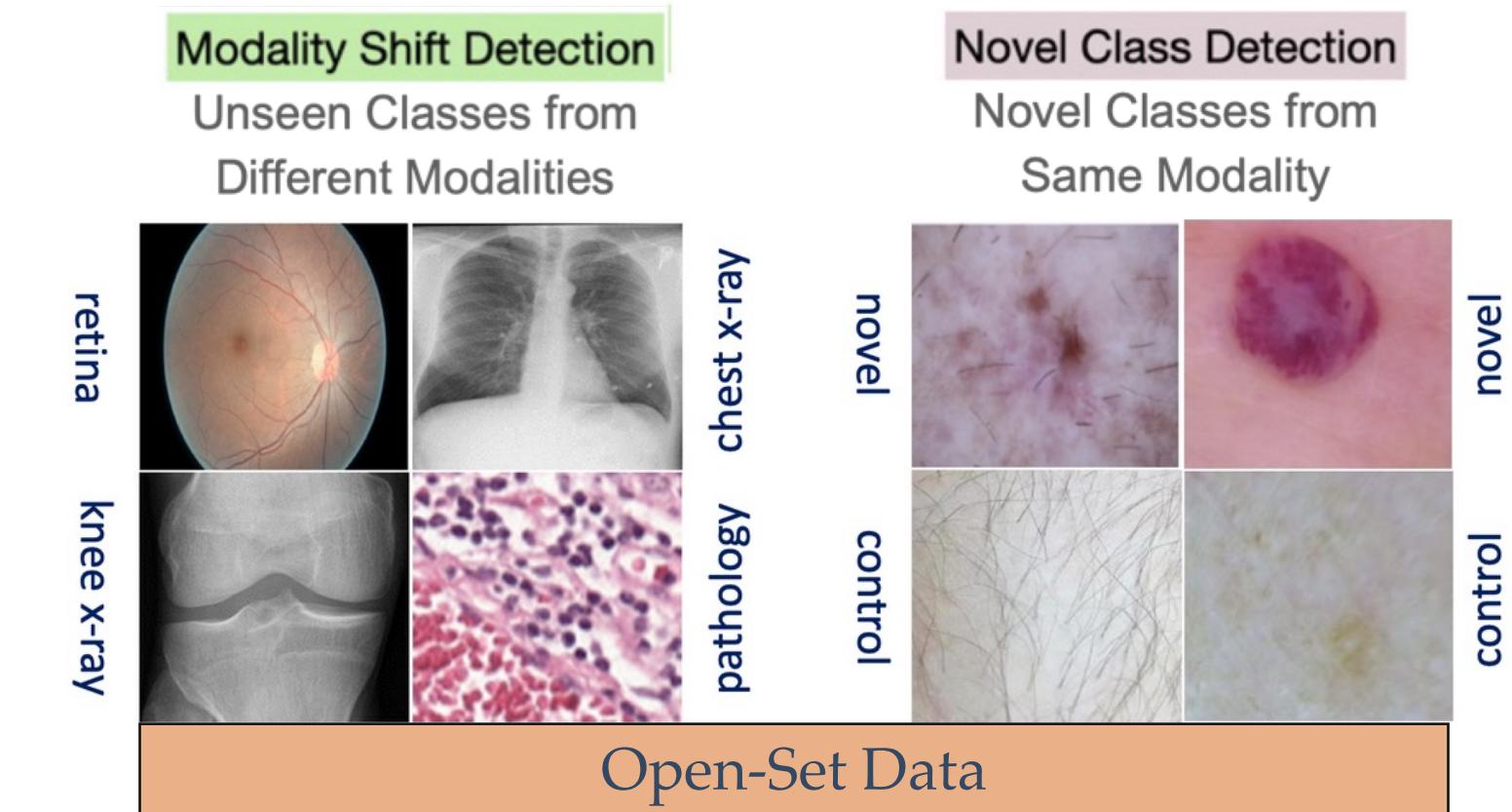
Jay
Thiagarajan
LLNL

Safe Deployment of AI Models Requires us to Monitor Predictions and Detect Unexpected Model Functionality

Ability to flag **open-set** samples with diverse semantic characteristics w.r.t the training data is a critical aspect of safety in medical AI

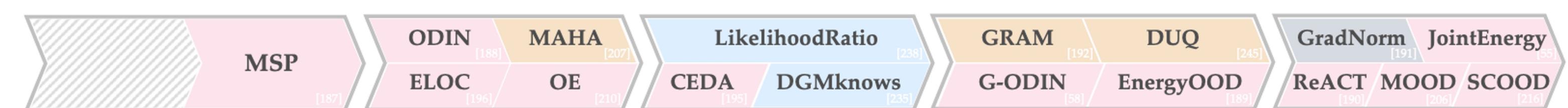
$$\gamma_{\text{OOD}} \neq \gamma_{\text{ID}}$$

Unseen semantic concepts



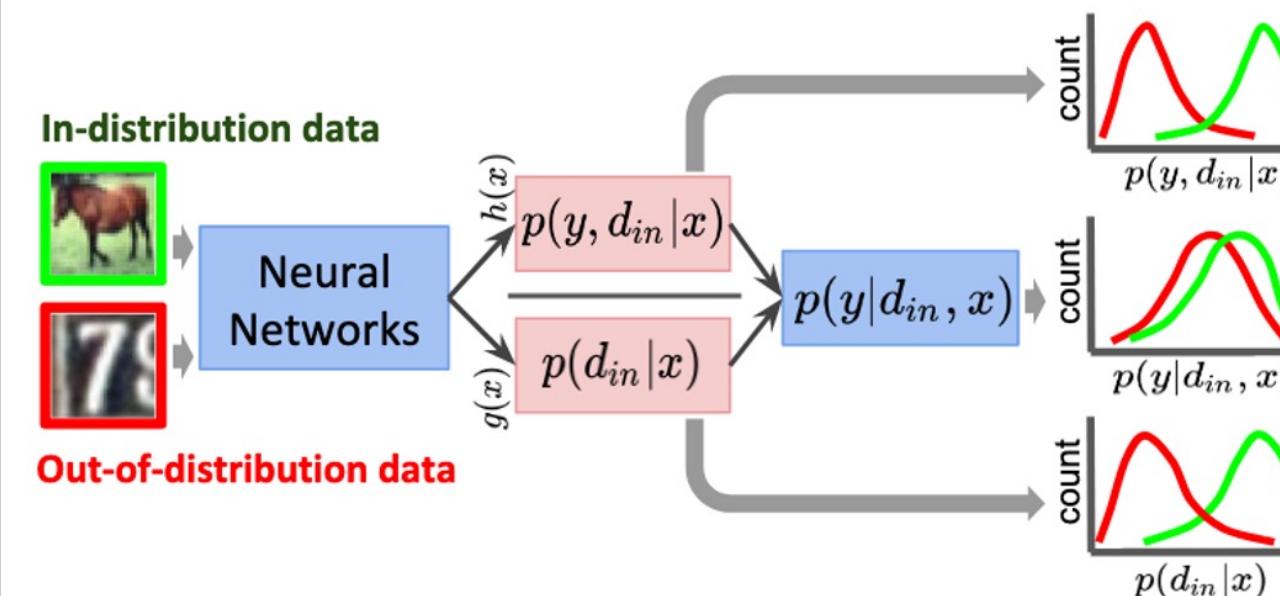
Flag as OOD and defer to experts

Constructing model scores for open-set detection is an active area of research in the vision community



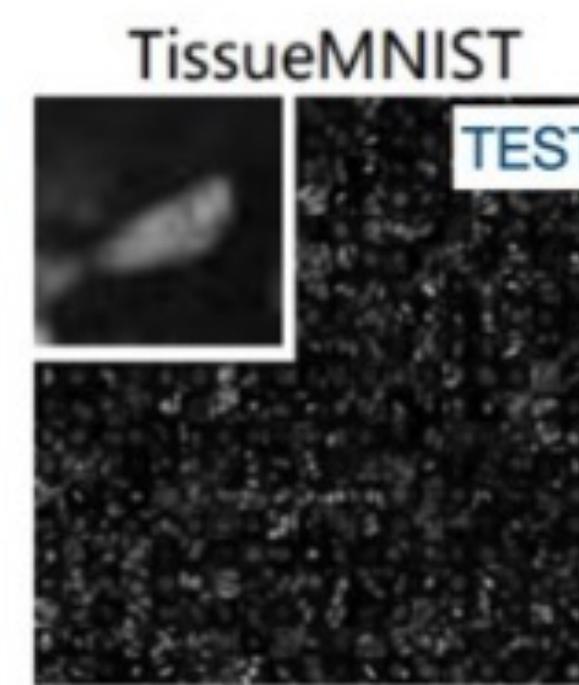
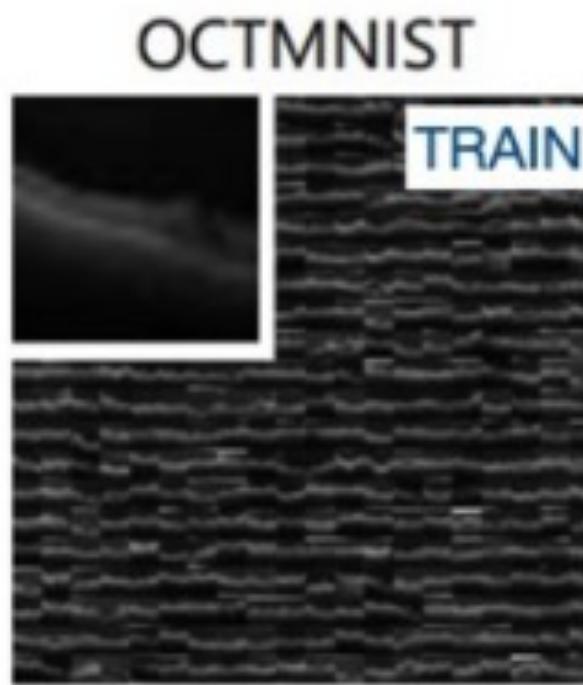
Key question: Can Off-the-shelf OOD Detectors Flag Open-Set Medical Image Data?

Generalized-ODIN



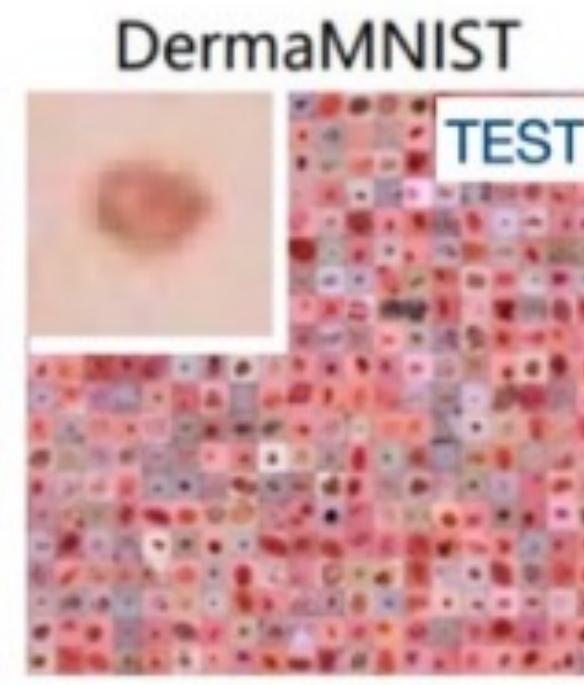
Despite achieving high accuracies on test data, conventional OOD detectors struggle with open data settings in medical imaging!

Modality Shift Detection



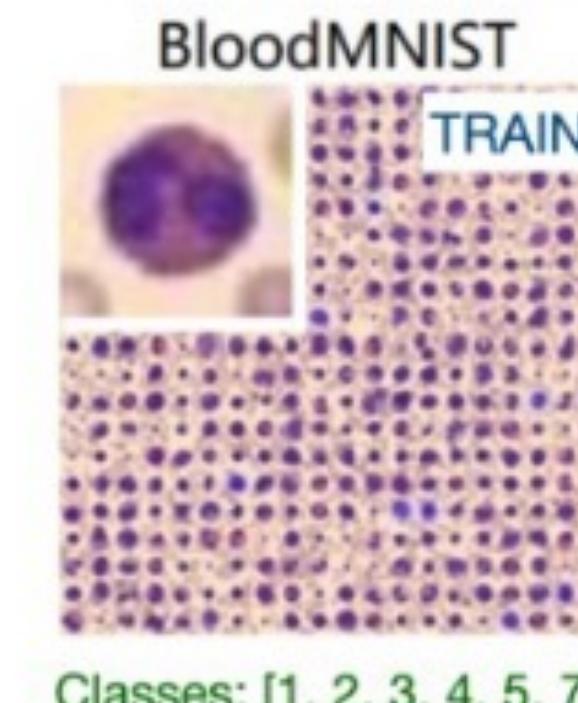
Balanced Accuracy: 95.9 ✓
OOD Rejection AUROC: 49.7 ✗

Modality Shift Detection



Balanced Accuracy: 99.3 ✓
OOD Rejection AUROC: 75.4 ✗

Novel Class Detection



Balanced Accuracy: 96.2 ✓
OOD Rejection AUROC: 53.9 ✗

A Potential Fix: Explicitly Calibrate OOD Detectors during Predictive Model Training

In this work, we consider the popular energy-based OOD detectors

$$G(\mathbf{x}; \theta, \tau) = \begin{cases} \text{outlier}, & \text{if } -E(\mathbf{x}; \theta) \leq \tau \\ \text{inlier}, & \text{if } -E(\mathbf{x}; \theta) > \tau \end{cases}$$

$$\min_{\theta} \quad \mathbb{E}_{(\mathbf{x}, y) \in \mathcal{D}} \mathcal{L}_{CE}(\mathcal{F}_{\theta}(\mathbf{x}), y) + \alpha \mathbb{E}_{\tilde{\mathbf{x}} \in \mathcal{D}_{in}} \mathcal{L}_{ID}(E(\tilde{\mathbf{x}}); \theta) + \beta \mathbb{E}_{\bar{\mathbf{x}} \in \mathcal{D}_{out}} \mathcal{L}_{OOD}(E(\bar{\mathbf{x}}); \theta)$$

Gibbs Distribution

$$p(y | \mathbf{x}) = \frac{e^{-E(\mathbf{x}, y)/T}}{\int_{y'} e^{-E(\mathbf{x}, y')/T}} = \frac{e^{-E(\mathbf{x}, y)/T}}{e^{-E(\mathbf{x})/T}}$$

Softmax Distribution

$$p(y | \mathbf{x}) = \frac{e^{f_y(\mathbf{x})/T}}{\sum_i^K e^{f_i(\mathbf{x})/T}}$$

Energy

$$E(\mathbf{x}; f) = -T \cdot \log \sum_i^K e^{f_i(\mathbf{x})/T}.$$

Temperature

Logit of Class i

A Potential Fix: Explicitly Calibrate OOD Detectors during Predictive Model Training

In this work, we consider the popular energy-based OOD detectors

$$G(\mathbf{x}; \theta, \tau) = \begin{cases} \text{outlier}, & \text{if } -E(\mathbf{x}; \theta) \leq \tau \\ \text{inlier}, & \text{if } -E(\mathbf{x}; \theta) > \tau \end{cases}$$

OOD detector calibration is implemented using margin-based loss functions

$$\min_{\theta} \quad \mathbb{E}_{(\mathbf{x}, y) \in \mathcal{D}} \mathcal{L}_{CE}(\mathcal{F}_{\theta}(\mathbf{x}), y) + \alpha \mathbb{E}_{\tilde{\mathbf{x}} \in \mathcal{D}_{in}} \mathcal{L}_{ID}(E(\tilde{\mathbf{x}}); \theta) + \beta \mathbb{E}_{\bar{\mathbf{x}} \in \mathcal{D}_{out}} \mathcal{L}_{OOD}(E(\bar{\mathbf{x}}); \theta)$$

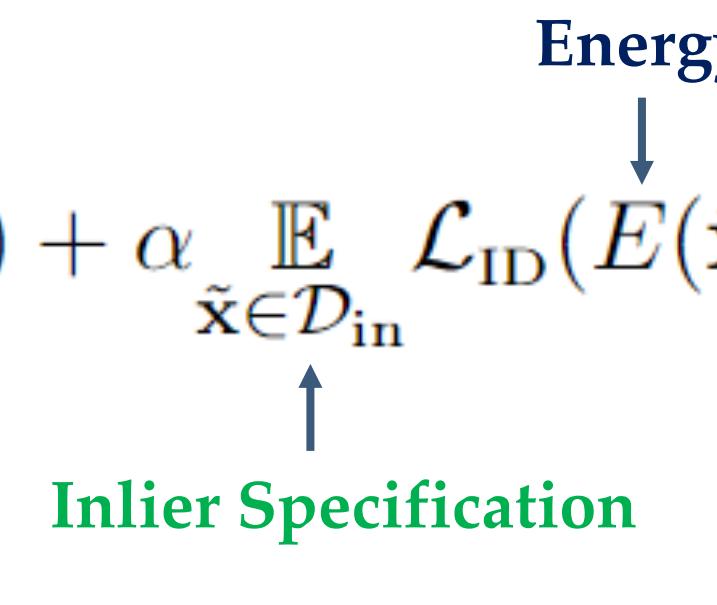
$$\mathcal{L}_{ID} = \mathbb{E}_{\mathbf{t}_k \sim \mathcal{T}} \left[\max \left(0, E(h(\mathbf{x}) = \mathbf{t}_k) - m_{ID} \right) \right]^2$$

$$\mathcal{L}_{OOD} = \mathbb{E}_{\bar{\mathbf{x}} \sim \mathcal{D}_{out}} \left[\max \left(0, m_{OOD} - E(\mathbf{x} = \bar{\mathbf{x}}) \right) \right]^2.$$

A Potential Fix: Explicitly Calibrate OOD Detectors during Predictive Model Training

In this work, we consider the popular energy-based OOD detectors

$$G(\mathbf{x}; \theta, \tau) = \begin{cases} \text{outlier}, & \text{if } -E(\mathbf{x}; \theta) \leq \tau \\ \text{inlier}, & \text{if } -E(\mathbf{x}; \theta) > \tau \end{cases}$$

$$\min_{\theta} \mathbb{E}_{(\mathbf{x}, y) \in \mathcal{D}} \mathcal{L}_{CE}(\mathcal{F}_{\theta}(\mathbf{x}), y) + \alpha \mathbb{E}_{\tilde{\mathbf{x}} \in \mathcal{D}_{in}} \mathcal{L}_{ID}(E(\tilde{\mathbf{x}}); \theta) + \beta \mathbb{E}_{\bar{\mathbf{x}} \in \mathcal{D}_{out}} \mathcal{L}_{OOD}(E(\bar{\mathbf{x}}); \theta)$$


Held-out calibration set from the train data can be used to specify inliers – **challenging in small data scenarios**

How is this optimization carried out in practice?

A representative OOD dataset is curated for specifying the outlier regimes – **non-trivial in medical imaging**

OOD calibration must not compromise the accuracy of the trained detector – **avoid over-conservative models**

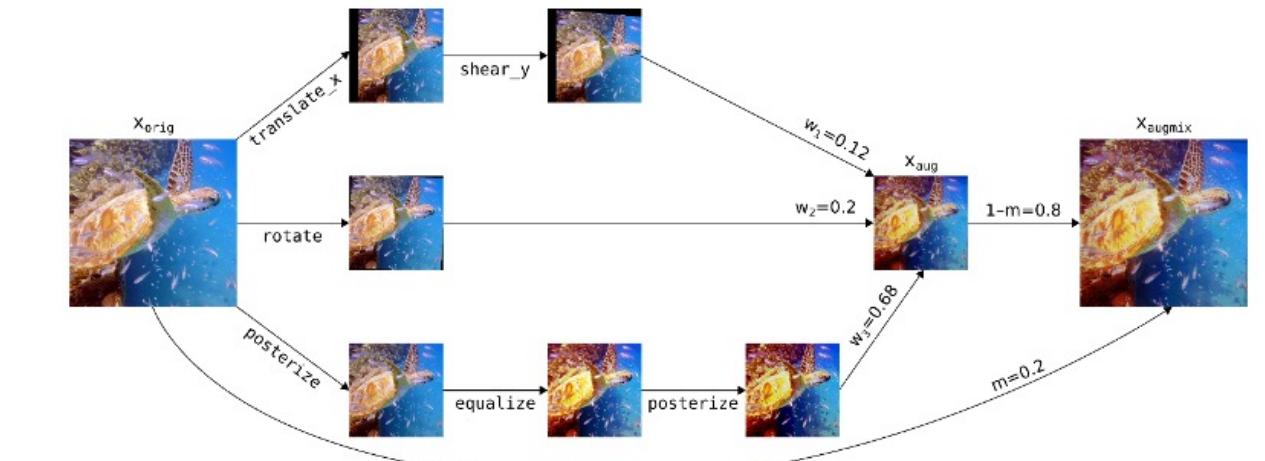
An Alternative Approach: Using Synthetic Data Augmentations to Specify Inliers and Outliers

Inlier Specification



Geometric Transforms

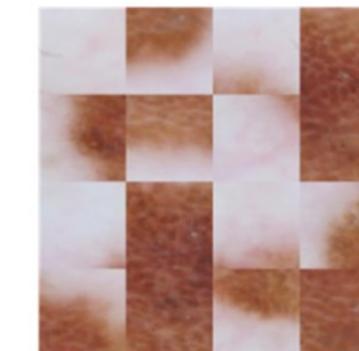
Pixel-Space Augmentations



Compositional (e.g., AugMix))

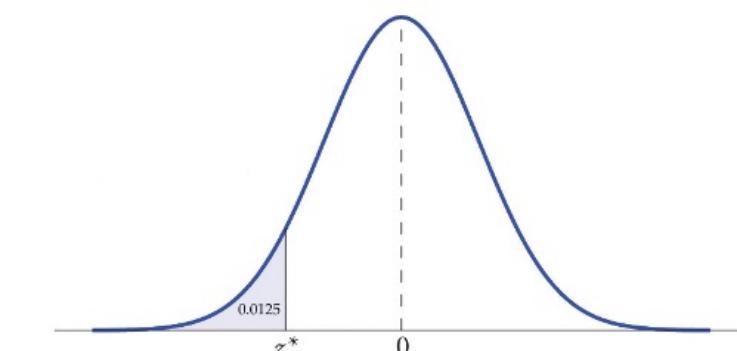
Outlier Specification

Outlier Exposure with Representative Datasets



Pixel-Space Augmentation
(Negative Data Augmentation)

Hard to Define for Medical Imaging Models



Latent-Space Augmentation
(Virtual Outlier Synthesis)

Hendrycks, Dan, et al. "Augmix: A simple data processing method to improve robustness and uncertainty." *arXiv preprint arXiv:1912.02781* (2019).

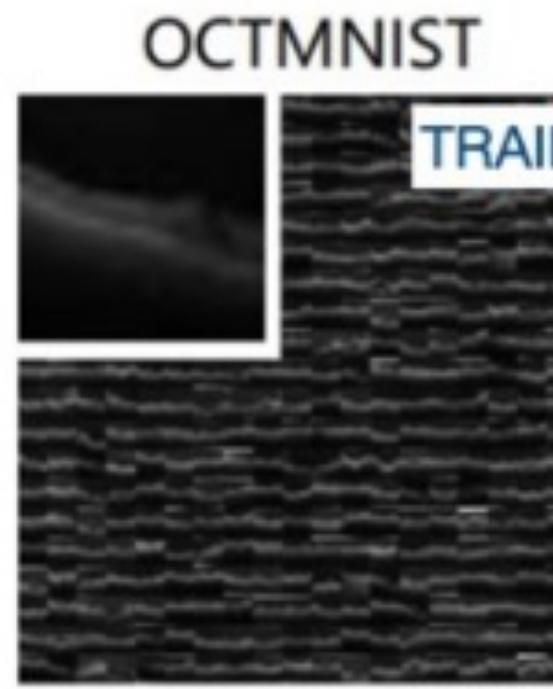
Hendrycks, Dan, Mantas Mazeika, and Thomas Dietterich. "Deep anomaly detection with outlier exposure." *arXiv preprint arXiv:1812.04606* (2018).

Sinha, A et al. Negative data augmentation. In International Conference on Learning Representations, 2021

Xuefeng Du, Zhaoning Wang, Mu Cai, and Yixuan Li. Vos: Learning what you don't know by virtual outlier synthesis. *arXiv preprint arXiv:2202.01197*, 2022

How Effective are the Calibrated OOD Detectors?

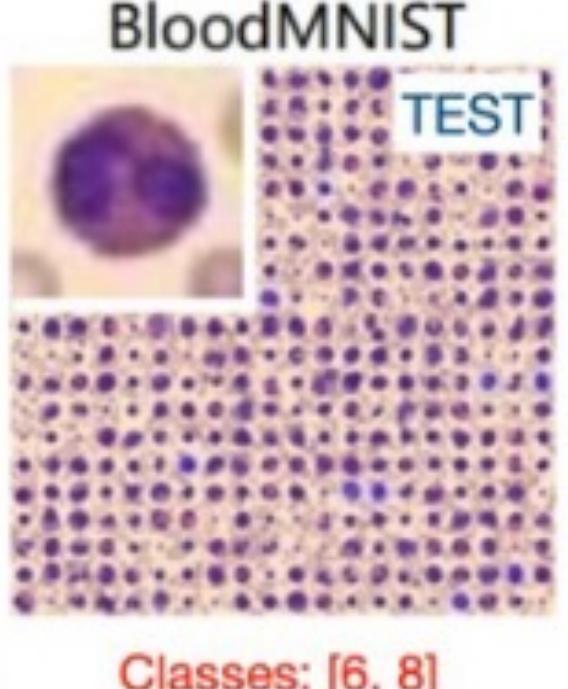
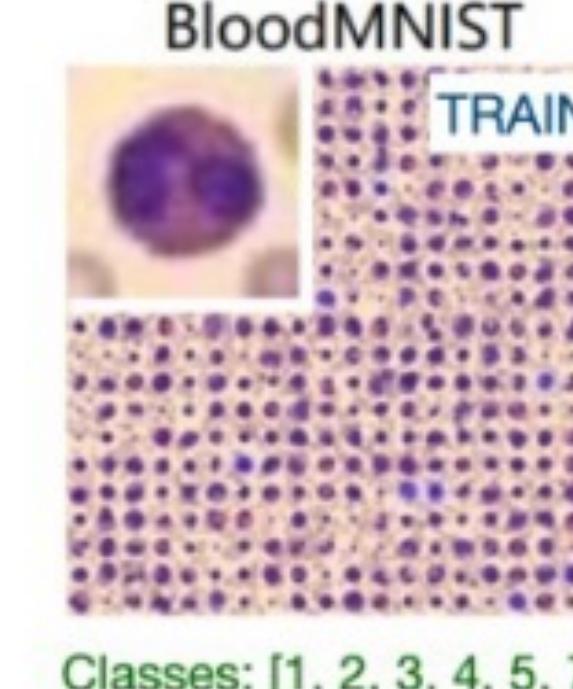
Modality Shift Detection



Modality Shift Detection



Novel Class Detection



G-ODIN

Balanced Accuracy: 95.9

OOD Rejection AUROC: 49.7 ✗

Balanced Accuracy: 99.3

OOD Rejection AUROC: 75.4 ✗

Balanced Accuracy: 96.2

OOD Rejection AUROC: 53.9 ✗

Augmix
+
vos

Balanced Accuracy: 95.3

OOD Rejection AUROC: 62.0 ✗

Balanced Accuracy: 99.2

OOD Rejection AUROC: 39.8 ✗

Balanced Accuracy: 96.5

OOD Rejection AUROC: 38.2 ✗

Augmix
+
NDA

Balanced Accuracy: 95.5

OOD Rejection AUROC: 90.9 ✓

Balanced Accuracy: 99.2

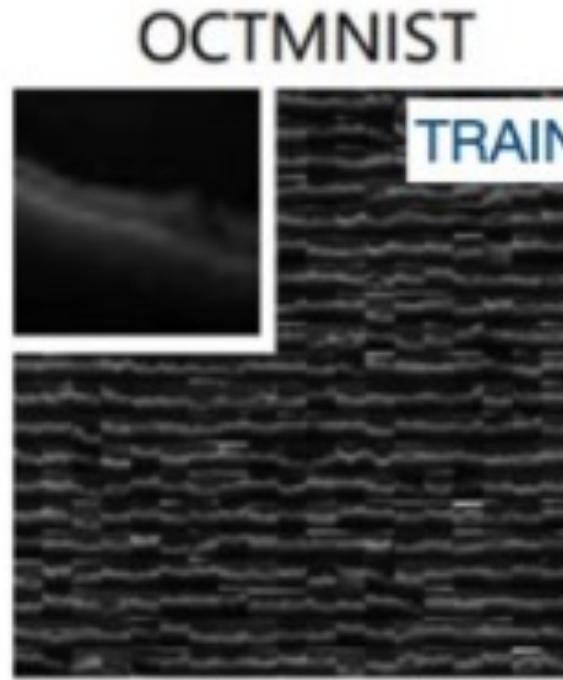
OOD Rejection AUROC: 43.5 ✗

Balanced Accuracy: 96.9

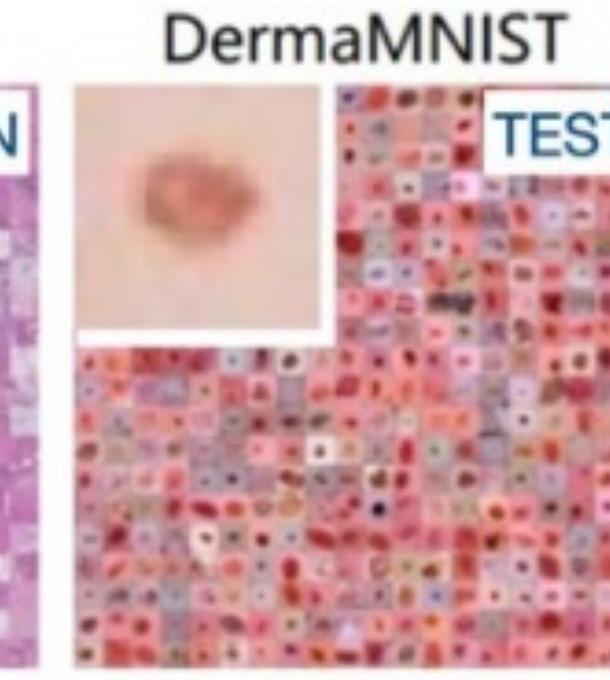
OOD Rejection AUROC: 53.5 ✗

How Effective are the Calibrated OOD Detectors?

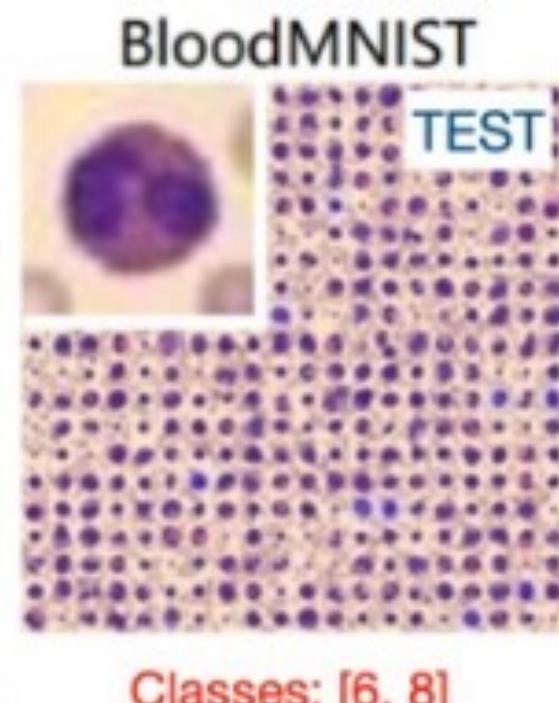
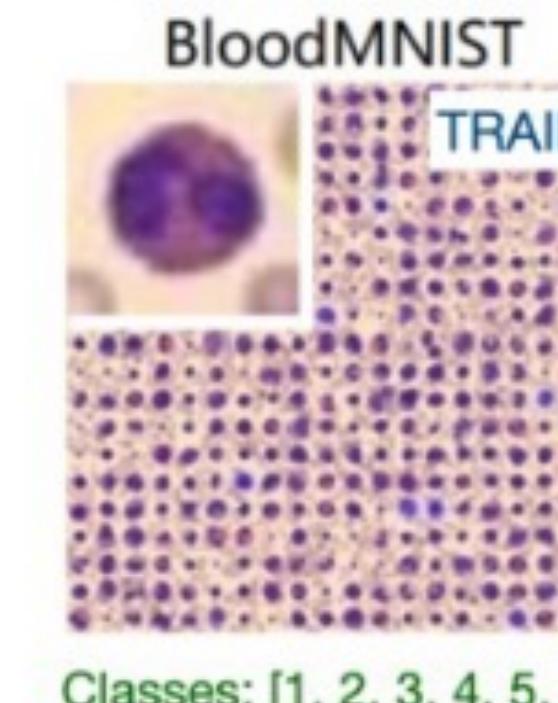
Modality Shift Detection



Modality Shift Detection



Novel Class Detection



G-ODIN

Balanced Accuracy: 95.9

OOD Rejection AUROC: 49.7 ✗

Balanced Accuracy: 99.3

OOD Rejection AUROC: 75.4 ✗

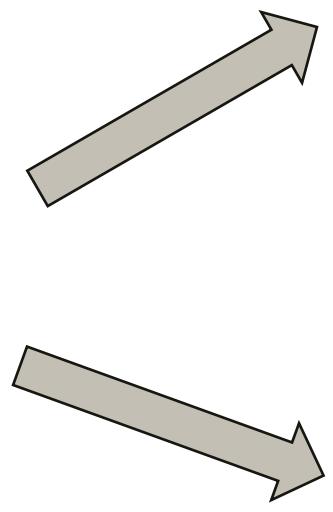
Balanced Accuracy: 96.2

OOD Rejection AUROC: 53.9 ✗

Surprisingly, OOD detectors calibrated using state-of-the-art approaches from vision literature do not perform consistently on both modality shifts and novel class scenarios

Hypothesis: Space in which inliers and Outliers are Specified Plays a Critical Role in Calibrating Medical OOD Detectors

With no outlier exposure, feature updates in a deep network are concentrated in the subspaces pertinent to the ID data

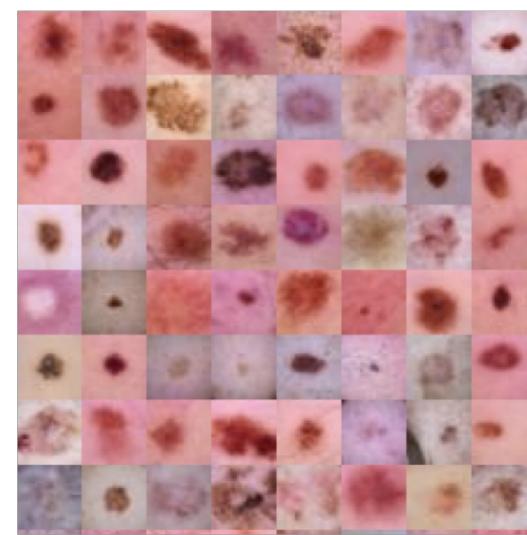
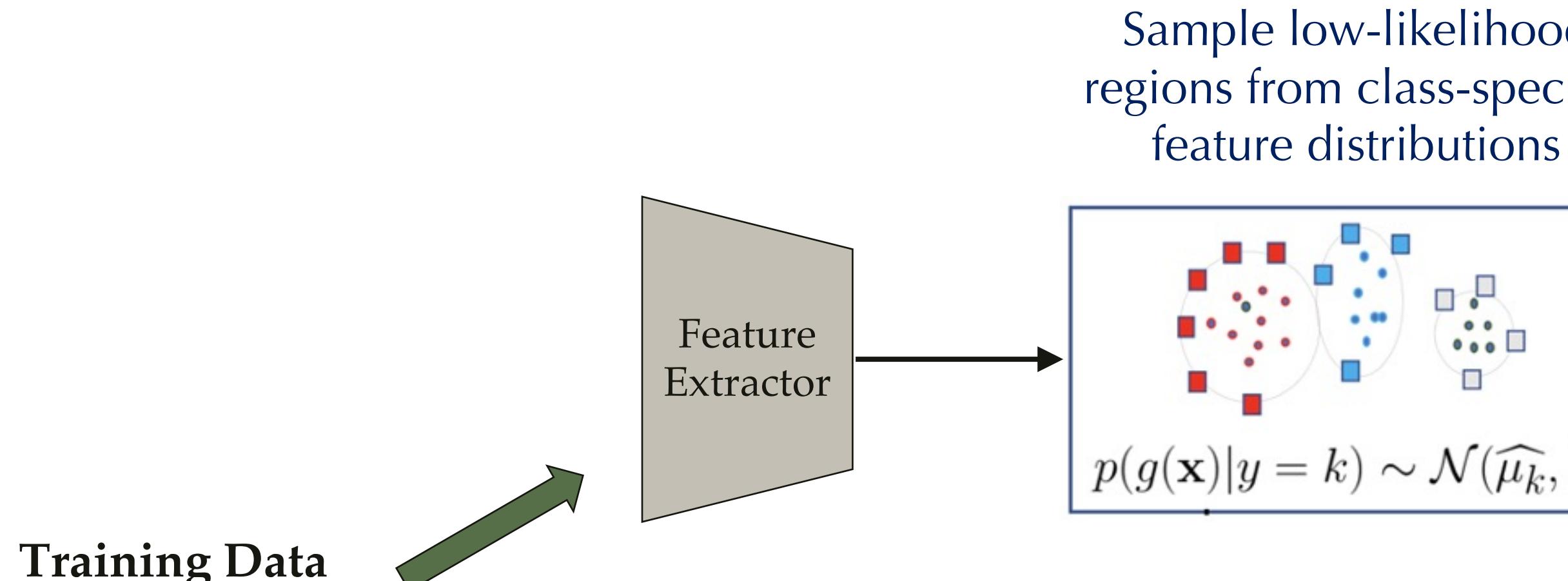


Inlier specification is used to expand model generalization and identify the optimal subspace for ID data – Protect against shortcut learning

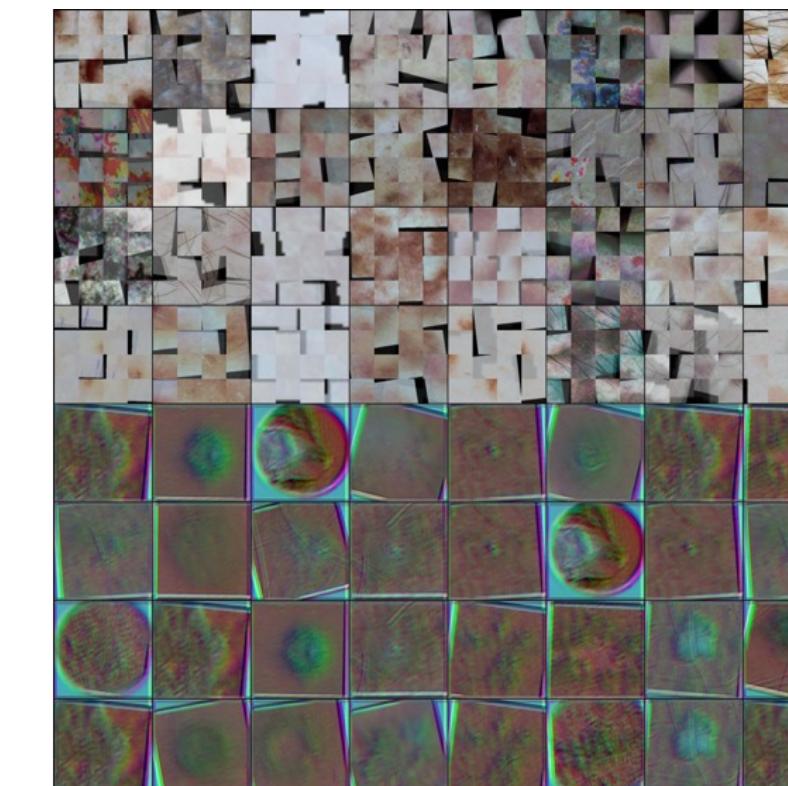
Outlier specification is required to ensure that the subspaces for outlier data do not overlap with the ID subspaces – Avoid over-generalization

Contrary to existing works, we advocate for the use of the latent-space for inlier specification and pixel-space for outlier specification to perform OOD calibration, while not compromising on the test accuracy

Proposed Approach for Calibrating Energy-based OOD Detectors in Medical Imaging Models



Synthesize high-severity compositional image manipulations
(e.g., Augmix, RandConv)



$$-\log \sum_{k=1}^K \exp \left(-E(h(\mathbf{x}) = \hat{\mu}_k, k) + \frac{1}{2} (\mathbf{t}_k - \hat{\mu}_k)^\top \hat{\Sigma}^{-1} (\mathbf{t}_k + \hat{\mu}_k) \right)$$

Lower bound on energy for the virtual inliers

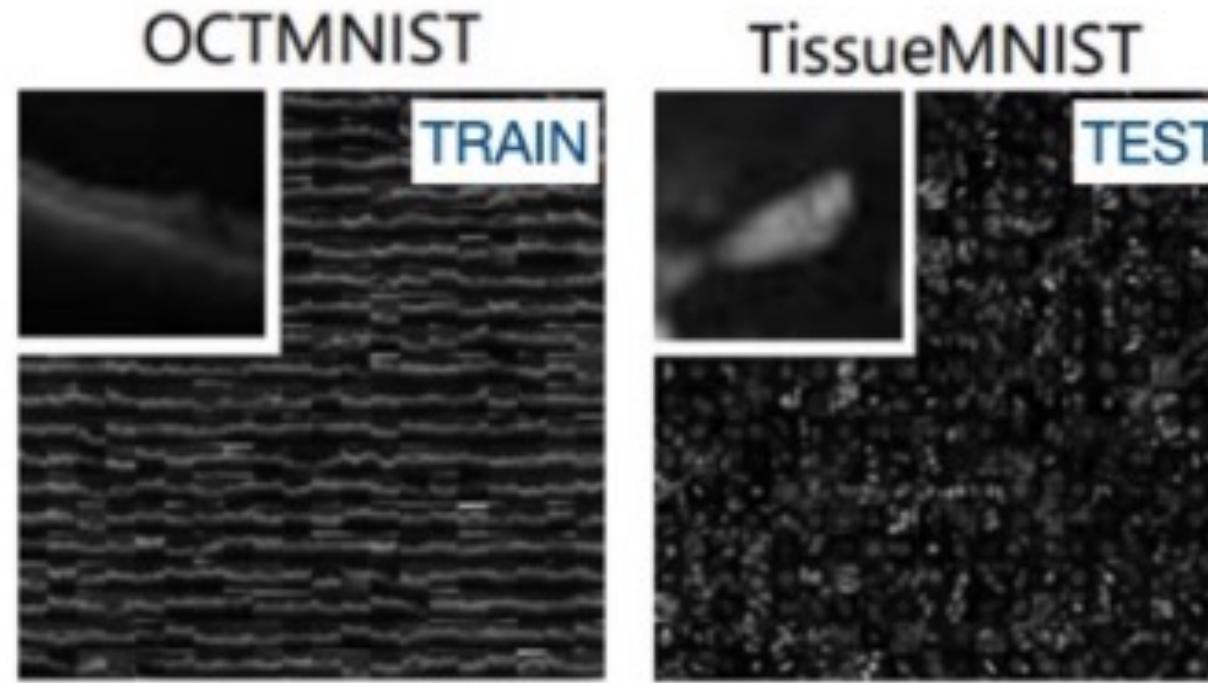
$$E(h(\mathbf{x}) = \mathbf{t}_k)$$

push the tail samples closer to the class-specific prototypes

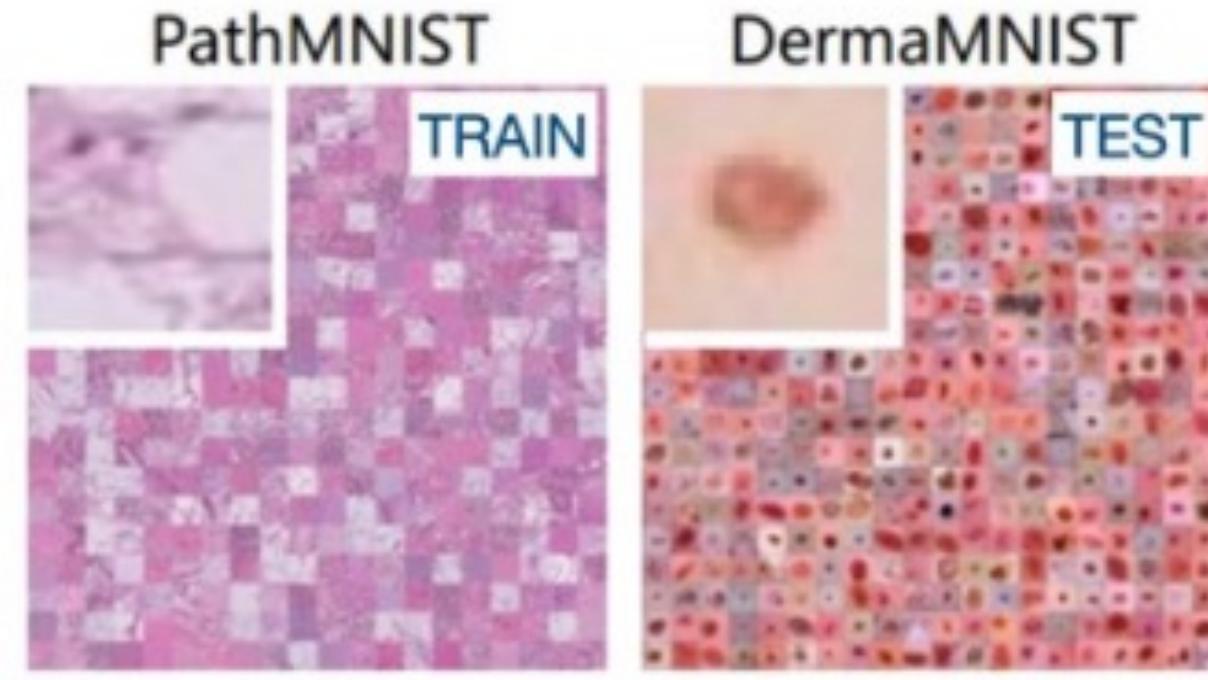
Highly diverse set of pixel-space outliers to ensure the OOD subspace does not overlap with the ID subspace in the feature space of the deep network

Using the Combination of Latent-Space Inliers and Synthetic Pixel-Space Outliers Leads to Powerful OOD Detectors

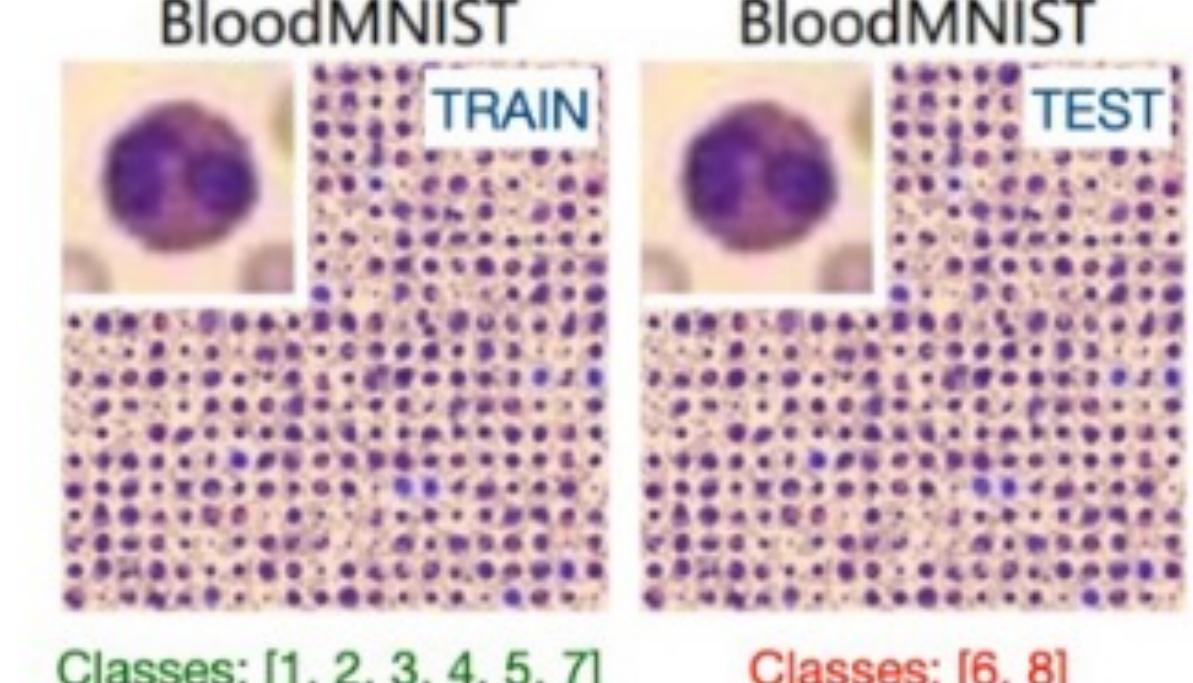
Modality Shift Detection



Modality Shift Detection



Novel Class Detection



G-ODIN | OOD Rejection AUROC: 49.7 ✗

Augmix
+
VOS | OOD Rejection AUROC: 62.0 ✗

Augmix
+
NDA | OOD Rejection AUROC: 90.9 ✓

Ours | OOD Rejection AUROC: 97.5 ✓

OOD Rejection AUROC: 75.4 ✗

OOD Rejection AUROC: 39.8 ✗

OOD Rejection AUROC: 43.5 ✗

OOD Rejection AUROC: 98.1 ✓

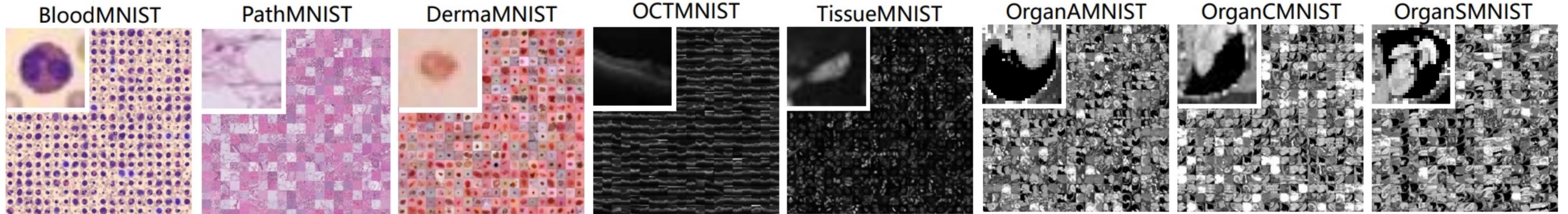
OOD Rejection AUROC: 53.9 ✗

OOD Rejection AUROC: 38.2 ✗

OOD Rejection AUROC: 53.5 ✗

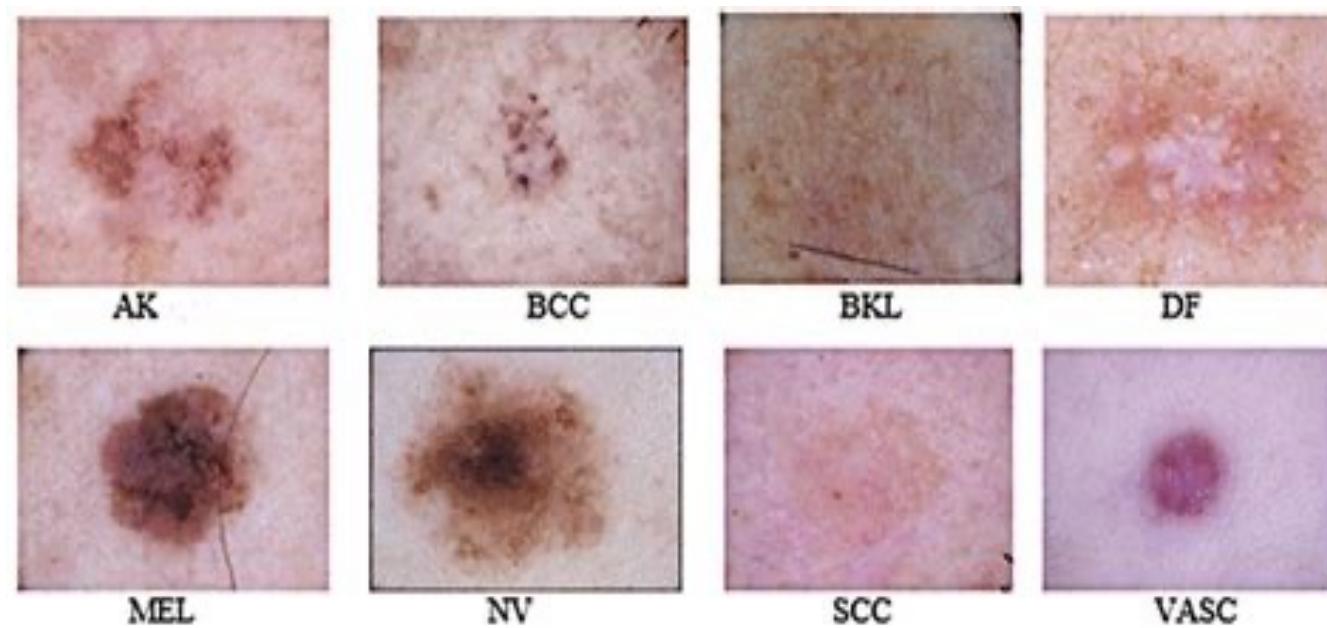
OOD Rejection AUROC: 89.1 ✓

With a Large Suite of Medical Imaging Benchmarks, We Systematically Evaluate our Proposed Approach



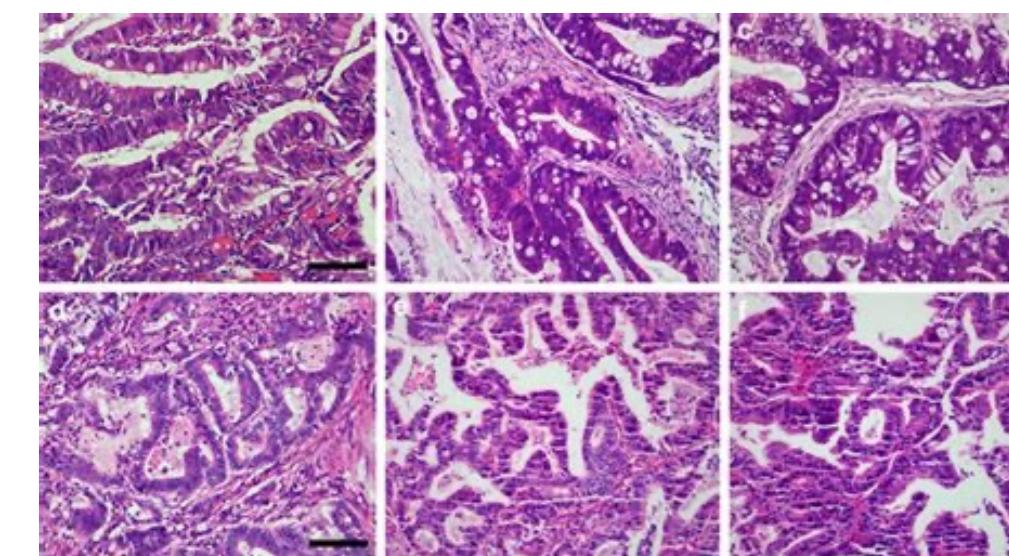
- **Modality shifts:** For each case, samples from all other datasets
- **Semantic shifts:** Held-out classes from the training dataset
- **Architecture:** WideResNet 40-2

ISIC 2019 Skin Lesion



- **Modality shifts:** CXR, WILDS, Retina images
- **Semantic shifts:** Held-out classes from train dataset, Clin Skin (control group), Derm Skin
- **Architecture:** ResNet 50

Colorectal Cancer



- **Modality shifts:** Clin Skin, Derm Skin, CXR, WILDS, Retina images
- **Semantic shifts:** Held-out classes from train dataset
- **Architecture:** ResNet 50

Evaluation

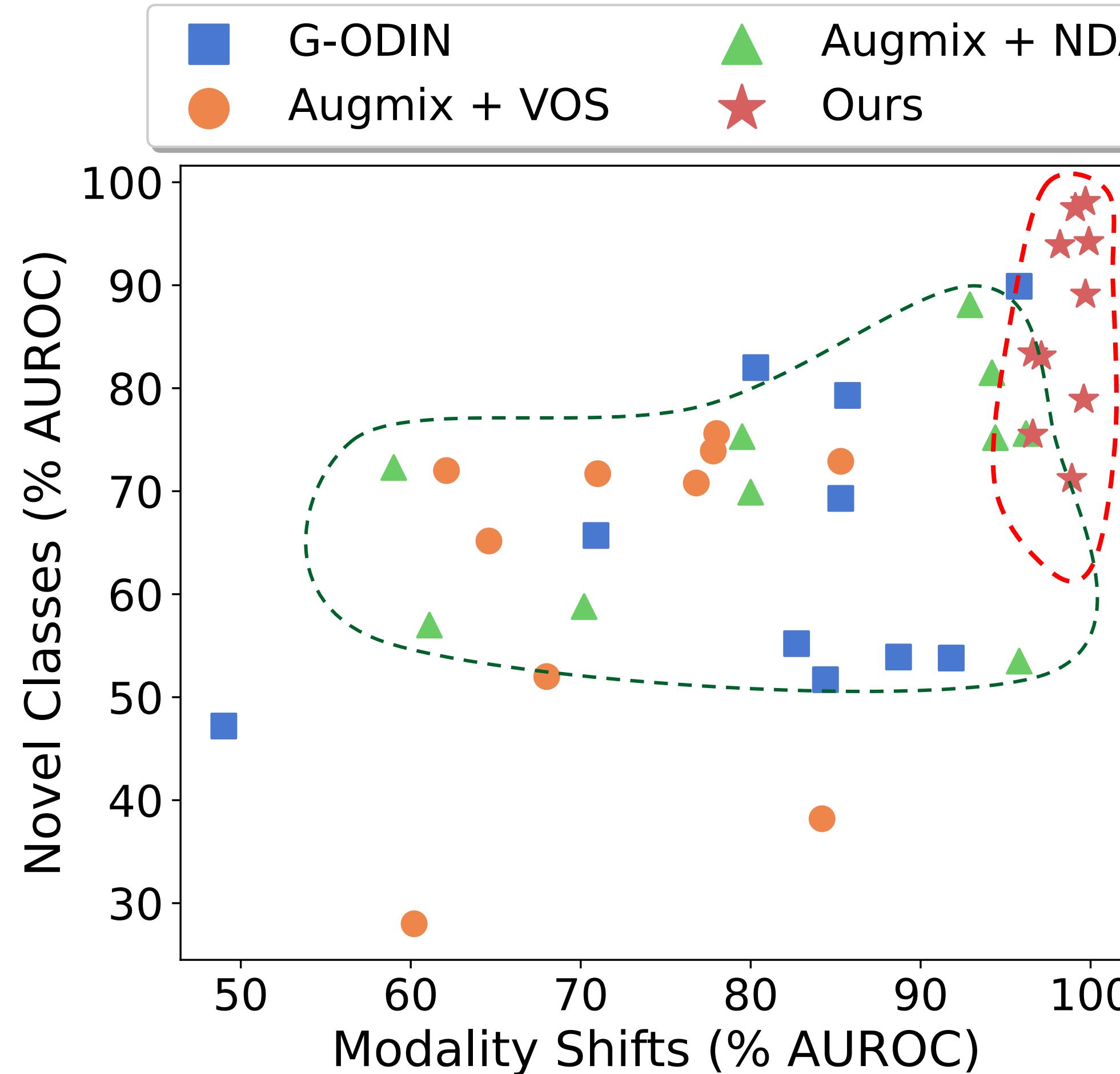
Balanced Accuracy

Modality Shift Detection

Novel Class Detection

AUROC

Strikingly, our Calibration Approach Leads to Significantly Superior Detection Performance in all Cases



Across all benchmarks, our approach achieves large gains over existing baselines in both modality shifts and novel classes

	G-ODIN	Augmix + VOS	Augmix + NDA	Ours
Modality Shifts	81.5	72.8	82.3	98.5
Novel Classes	64.8	62.1	70.7	86.49

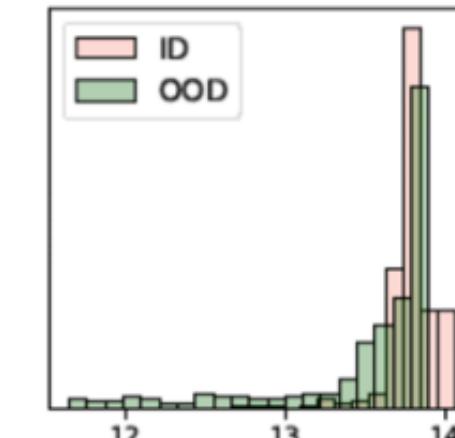
Existing baselines tend to produce large variances in AUROC scores across datasets

Visualization of the Energy Scores for ID and OOD Data Clearly Reveals the Benefits of the Proposed Calibration Protocol

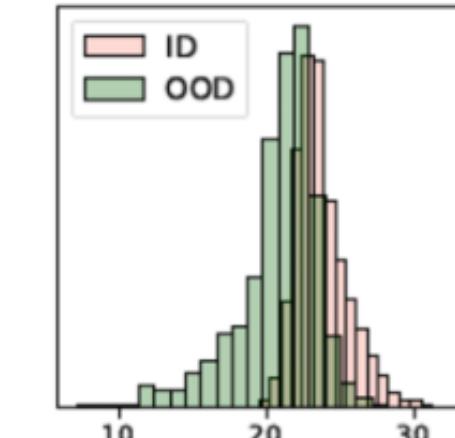
ID: BloodMNIST

Modality Shift: Derma MNIST

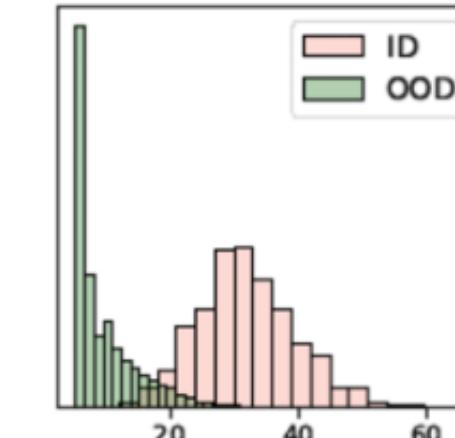
Augmix + VOS



Augmix + NDA



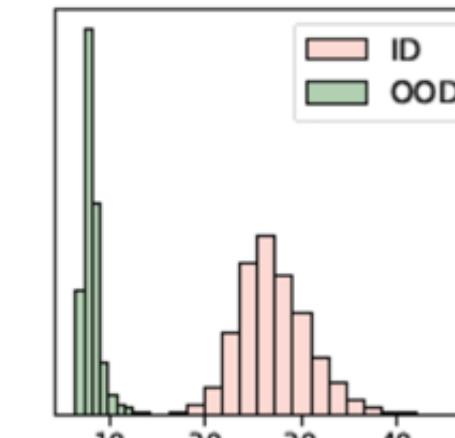
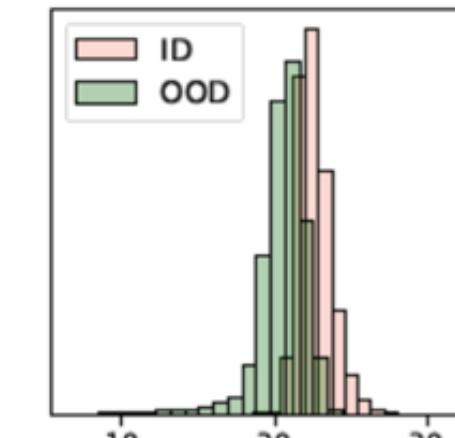
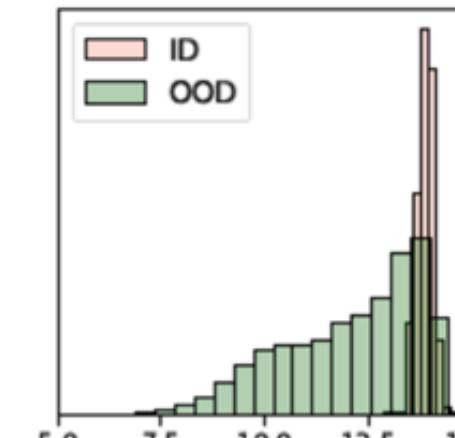
Ours



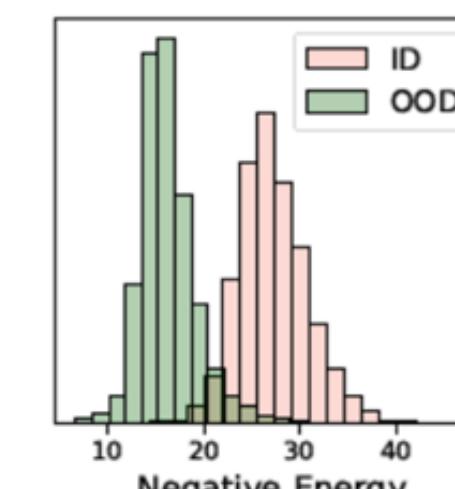
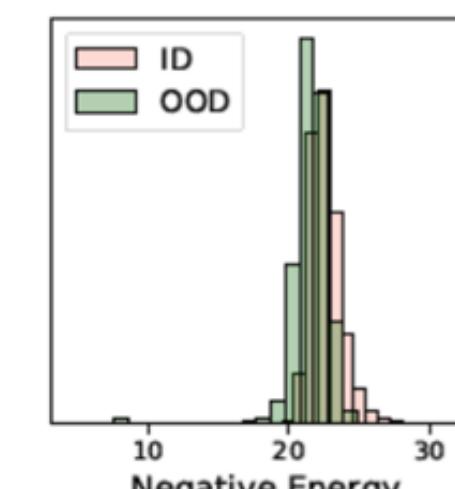
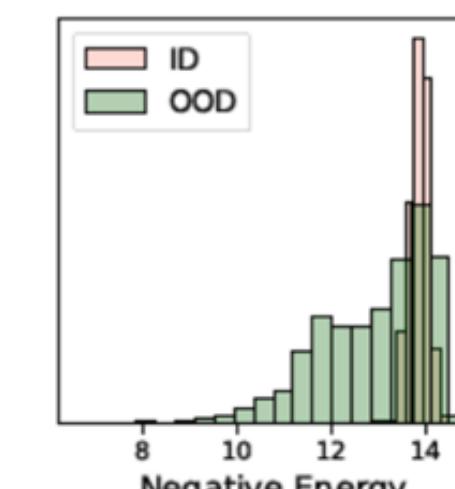
ID: OrganAMNIST

Novel Classes

Modality Shift: Tissue MNIST

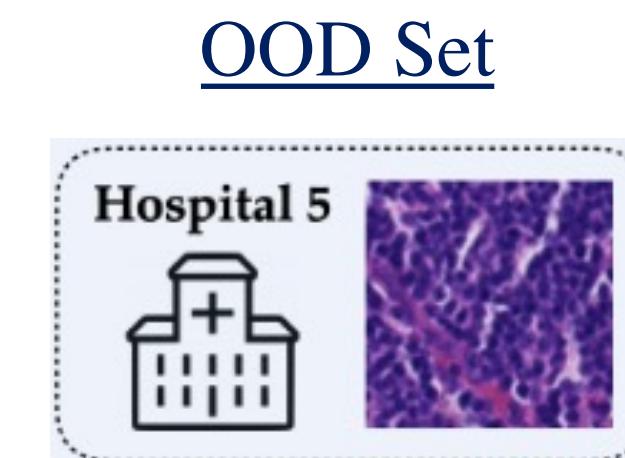
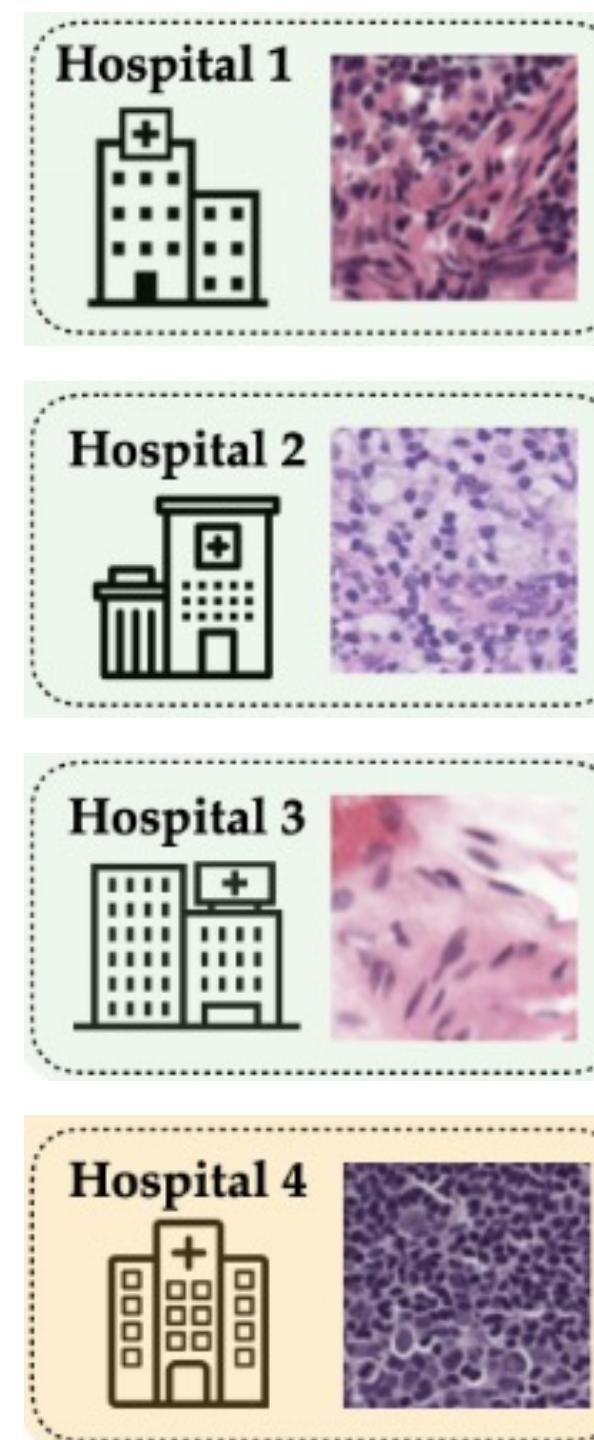


Novel Classes



Interestingly, our Approach can Even Detect Nuanced Covariate Shifts Arising from Different Hospitals

Observed Data



	AUROC (%)
Augmix + VOS	79.5
Augmix + NDA	36.5
Ours	87.4

While latent-space outliers are superior to pixel-space outliers synthesized via negative data augmentation, our approach performs the best

Summary: A New State-of-the-Art Baseline for Medical OOD Detection that can be used with any Imaging Modality or Model Architecture

Calibrating OOD detectors is significantly challenging with medical imaging data, and existing solutions from the vision literature do not work effectively!

We find that the choice of space for synthesizing augmentations is critical when calibrating OOD detectors for open-set data

We advocate for the use of virtual inliers from the classifier's latent-space and diverse pixel-space outliers with energy-based training

Using a large suite of medical imaging benchmarks, we show state-of-the-art open-set recognition performance (both modality shifts and novel classes), as well as in detecting covariate shifts arising from different hospital data.

Know Your Space - Inlier and Outlier Construction for Calibrating Medical OOD Detectors

Vivek Narayanaswamy *¹ Yamen Mubarka *¹ Rushil Anirudh¹ Deepta Rajan² Andreas Spanias³
Jayaraman J. Thiagarajan¹

¹Lawrence Livermore National Laboratory, CA, USA ²Microsoft, WA, USA ³Arizona State University, AZ, USA

[Paper](#) [Slides](#) [Poster](#) [Github](#)

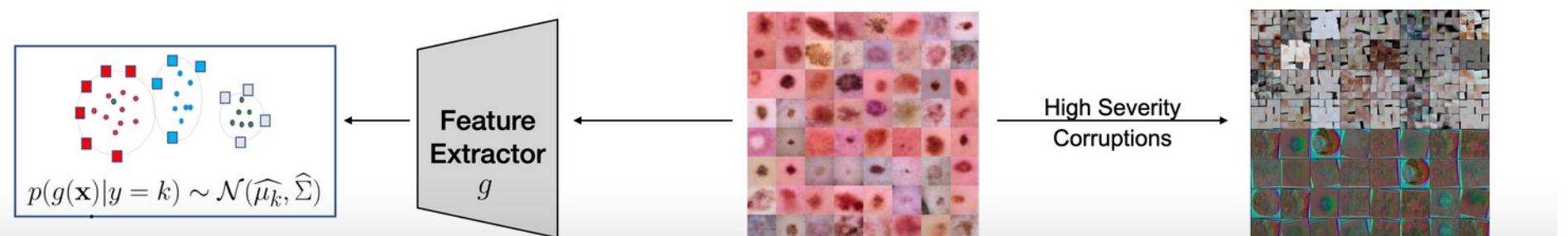
Summary

Our primary focus is on developing well-calibrated out-of-distribution (OOD) detectors to ensure the safe deployment of medical image classifiers. The use of synthetic augmentations has become common for specifying regimes of data inliers and outliers. However, our research findings highlight the substantial influence of both the synthesis space and the type of augmentation on the performance of OOD detectors. After conducting an extensive study using medical imaging benchmarks and open-set recognition settings, we recommend employing a combination of virtual inliers in the classifier's latent space and diverse synthetic outliers in the pixel space. This approach proves highly effective in producing OOD detectors with superior performance.

Video

To be posted

Method



Codes



Website



For questions contact
narayanaswam1@llnl.gov