# Evaluating Reliability in Medical DNNs: A Critical Analysis of Feature and Confidence-based Out-of-Distribution Detection

Harry Anthony, Konstantinos Kamnitsas

Department of Engineering Science, University of Oxford, Oxford, UK.

✉ harry.anthony@eng.ox.ac.uk

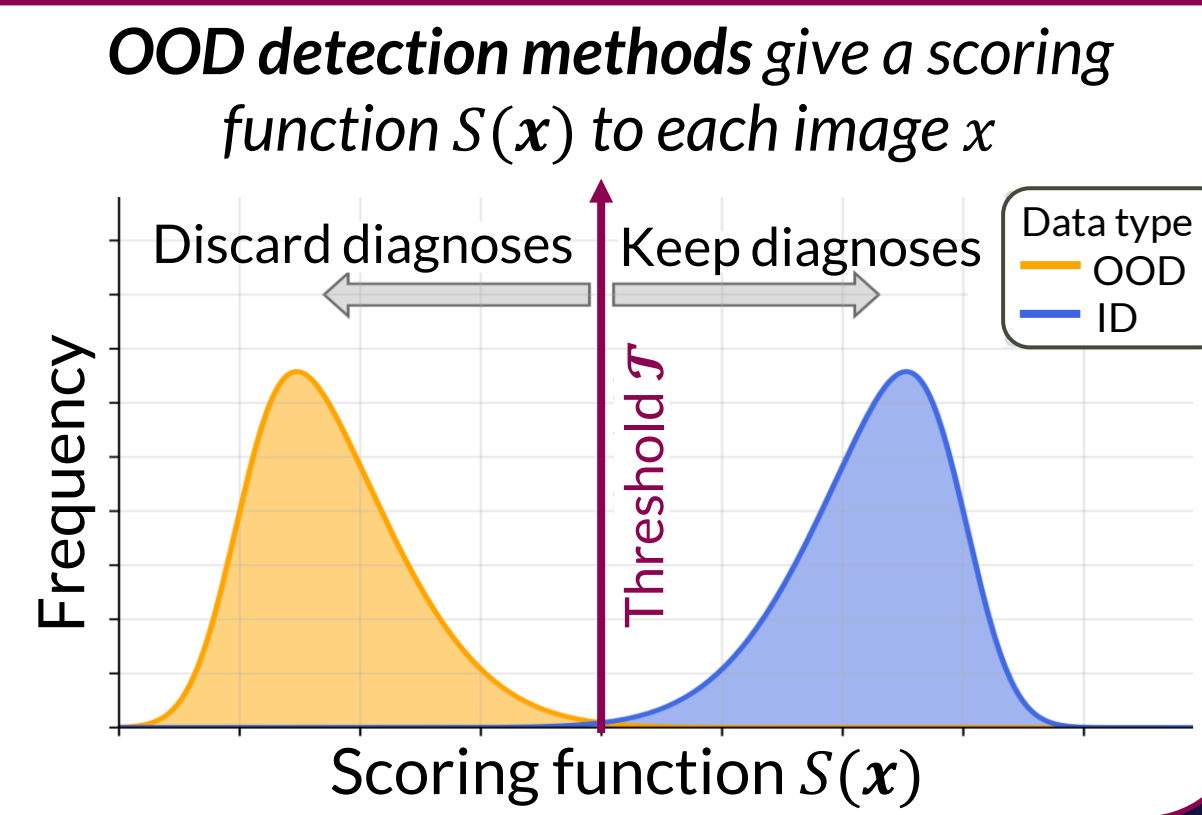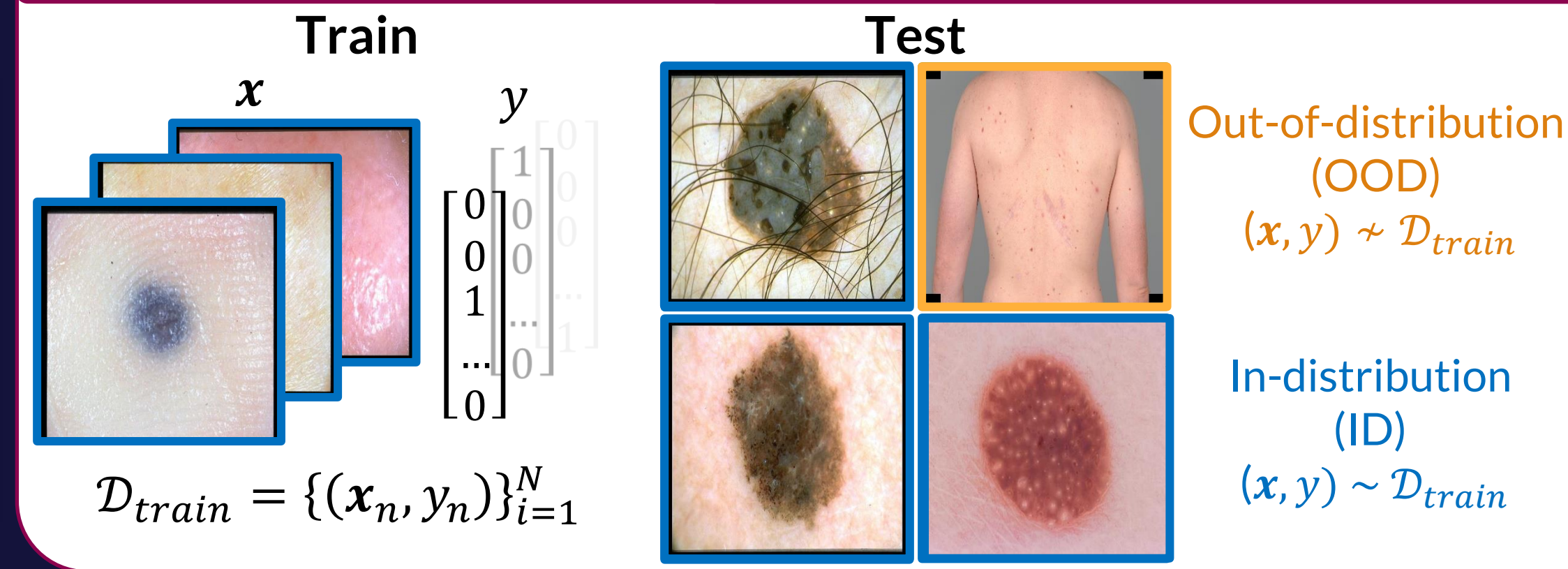○ HarryAnthony
✕ @HarryEJAnthony

## Introduction

Reliable neural networks must detect inputs that are **out-of-distribution (OOD)**.

⚠ Individual OOD detection methods have strengths and weaknesses.

🗄 We created new OOD counterfactual datasets to analyse these weaknesses.

💡 Combining complementary methods can mitigate against their weaknesses.

Paper available at:

Code available at:

## 1. What is Out-Of-Distribution (OOD) detection?



**Train** $x$ $y$

$\mathcal{D}_{train} = \{(x_n, y_n)\}_{i=1}^{N}$

**Test**

Out-of-distribution (OOD)
$(x, y) \nsim \mathcal{D}_{train}$

In-distribution (ID)
$(x, y) \sim \mathcal{D}_{train}$

OOD detection methods give a scoring function $S(x)$ to each image $x$

Discard diagnoses ← → Keep diagnoses

Threshold $\mathcal{T}$

Data type: OOD / ID

Scoring function $S(x)$

## 2. OOD Detection benchmarks

### a) New out-of-distribution benchmarks

**D7P** — Dermatology dataset

Training data: No rulers (90% ID data)
Training classes: Nevus 59%, Not Nevus 41%

ID Test (10% ID)   OOD Test: Grid rulers

**BreastMNIST** — Ultrasound dataset

Training data: No annotations (90% ID)
Malignant 27%, Normal 17%, Benign 56%
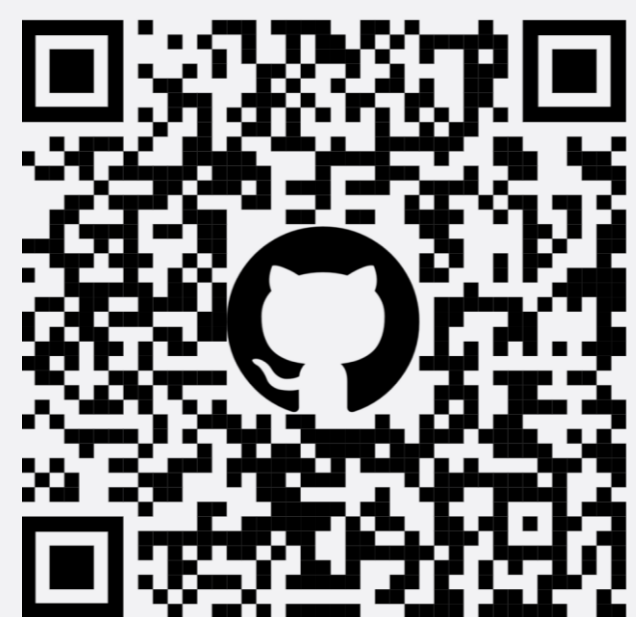
ID Test (10% ID)   OOD: Annotations

### b) New Counterfactual Datasets

*Dataset was created with inter-image interpolation, using a patch from the same image to remove OOD artefacts.*
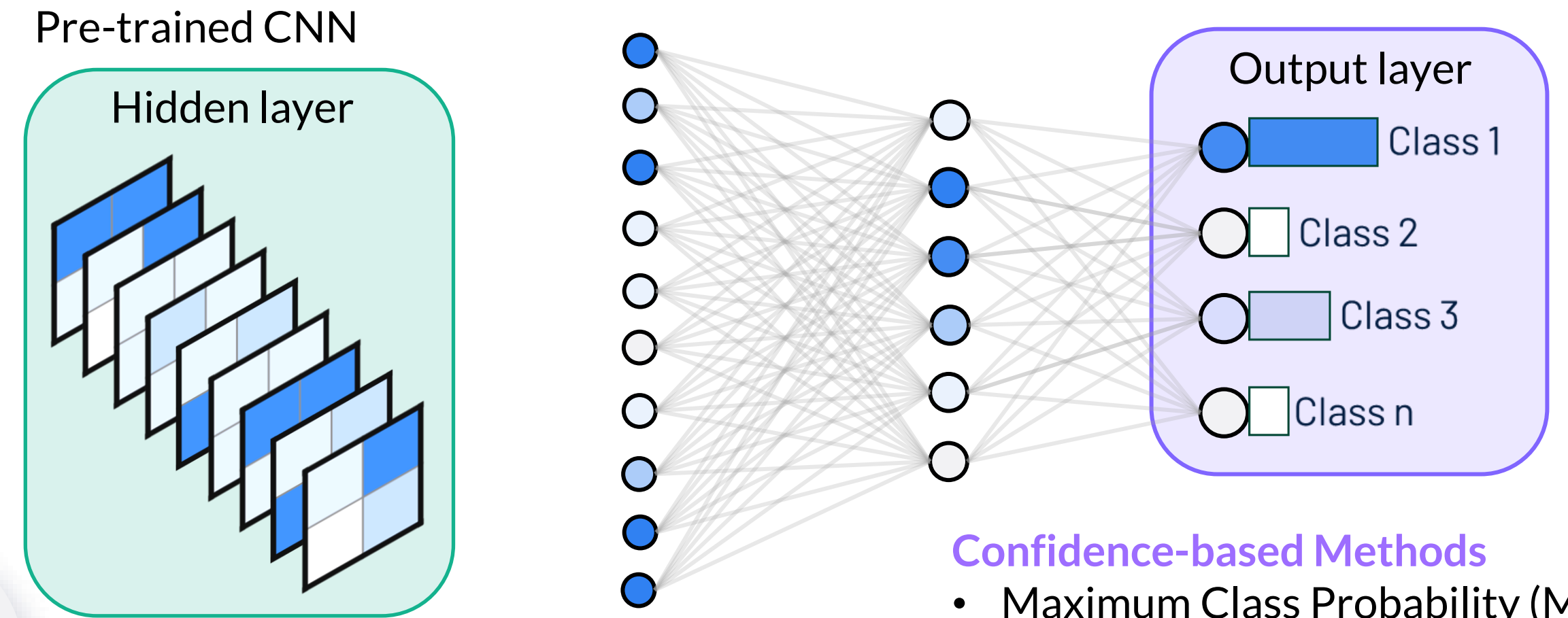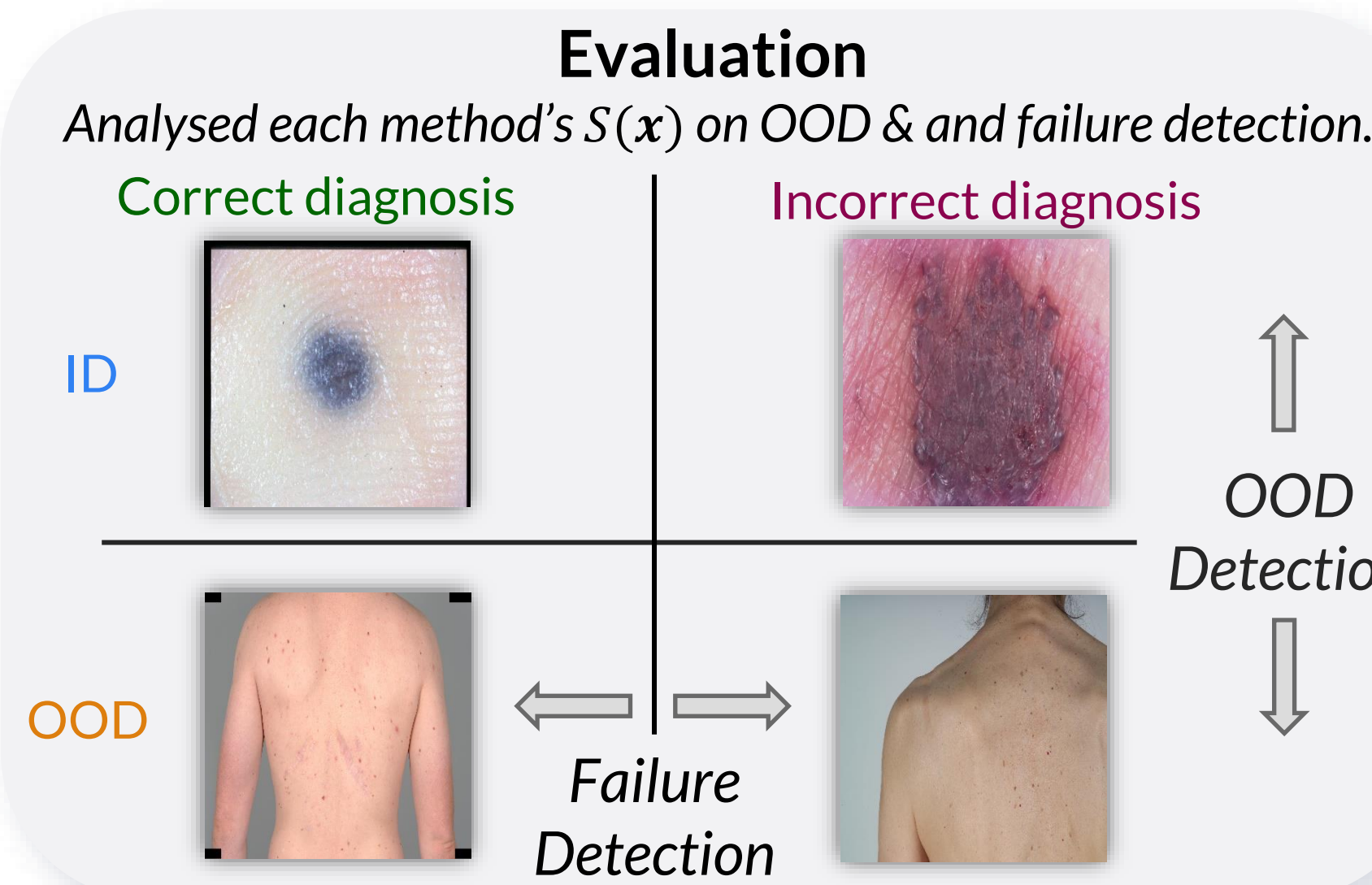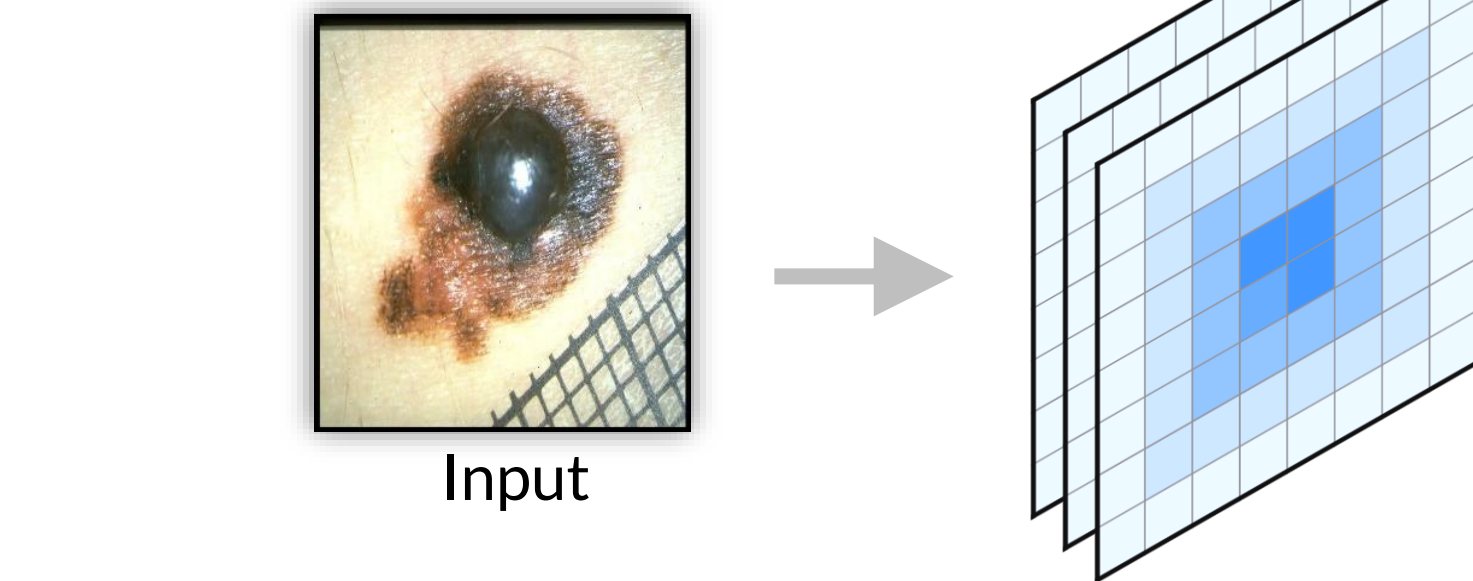
Original   Mask   Interpolation   Result

🗄 **Access this new data**
- Annotations for 2 new OOD benchmarks
- Pixel-wise artefact masks
- 478 image counterfactual datasets

## 3. Out-of-distribution Detection Methods: Confidence-based & Feature-based

*We analysed Post-hoc methods which are applied to pre-trained models*

Input

Pre-trained CNN
Hidden layer

Output layer
Class 1, Class 2, Class 3, Class n

**Evaluation**
*Analysed each method's $S(x)$ on OOD & failure detection.*

Correct diagnosis | Incorrect diagnosis
ID
OOD

OOD Detection
Failure Detection

**Feature-based Methods**
- Mahalanobis Score
- Multi-Branch Mahal. (MBM)
- RMS
- GRAM Matrices
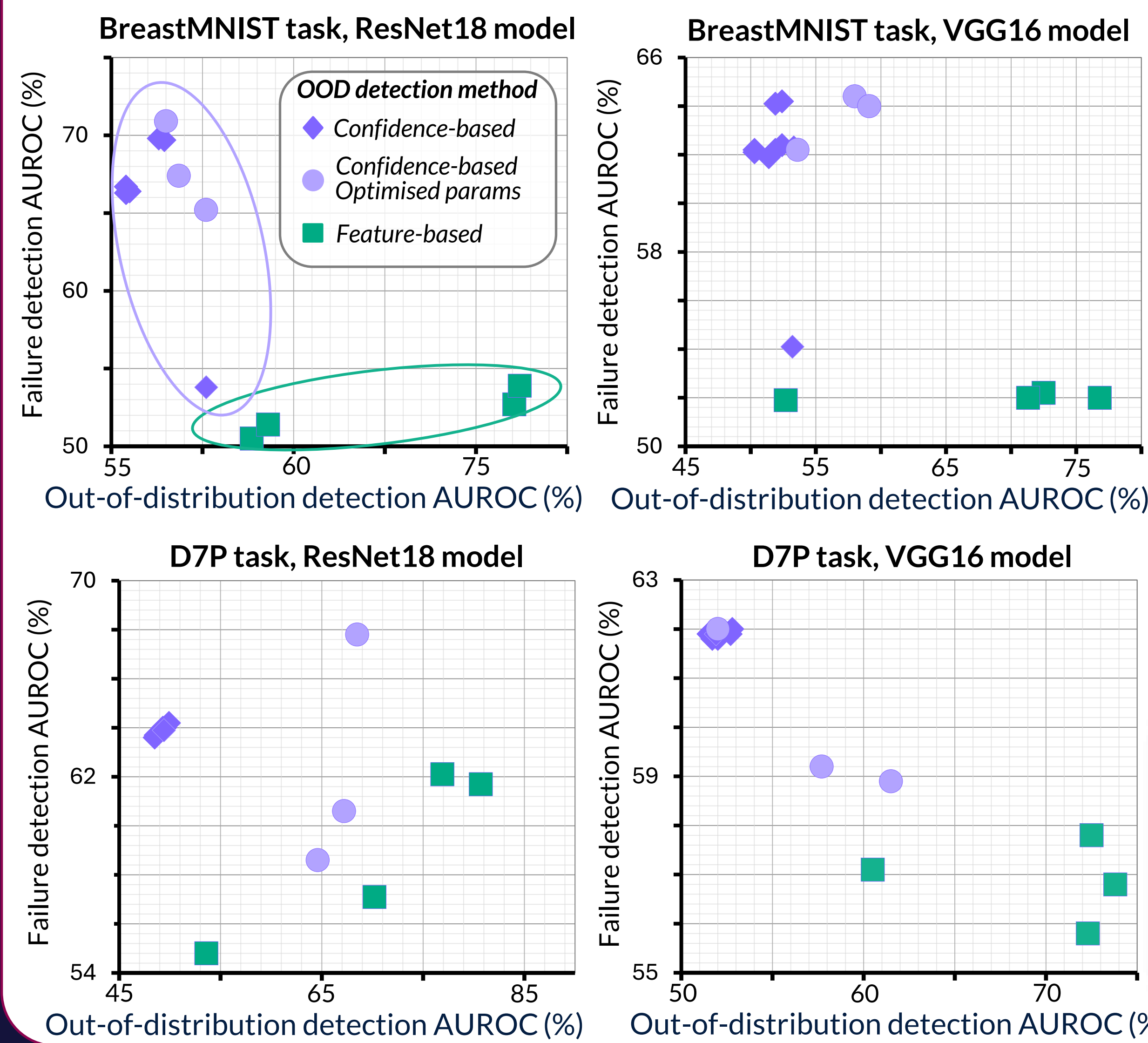
**Confidence-based Methods**
- Maximum Class Probability (MCP)
- Maximum Logit Score (MLS)
- Shannon Entropy (SE)
- Energy Score
- MC Dropout – MCP
- MC Dropout – Predictive Entropy
- MC Dropout – Mutual Information
- Deep Ensembles – MCP
- GradNorm
+ with optimised hyperparameters
- ODIN
- ReAct
- DICE

## 4. Experiments on OOD detection tasks

### a) OOD detection and failure detection evaluation



BreastMNIST task, ResNet18 model

OOD detection method:
- Confidence-based
- Confidence-based Optimised params
- Feature-based

Failure detection AUROC (%) vs Out-of-distribution detection AUROC (%)

BreastMNIST task, VGG16 model

D7P task, ResNet18 model

D7P task, VGG16 model

### b) Why do confidence-based methods have poor OOD detection?

**i** Synthetic image without artefact
Image | LRP heatmap
Softmax dist. — 0.99 Nevus, 0.01 Not Nevus → Correct diagnosis

**ii** Original Image with artefact
Image | LRP heatmap
Softmax dist. — 0.01 Nevus, 0.99 Not Nevus → Very confident mistake

→ OOD artefacts can lead to high confidence predictions which confidence-based methods won't detect!
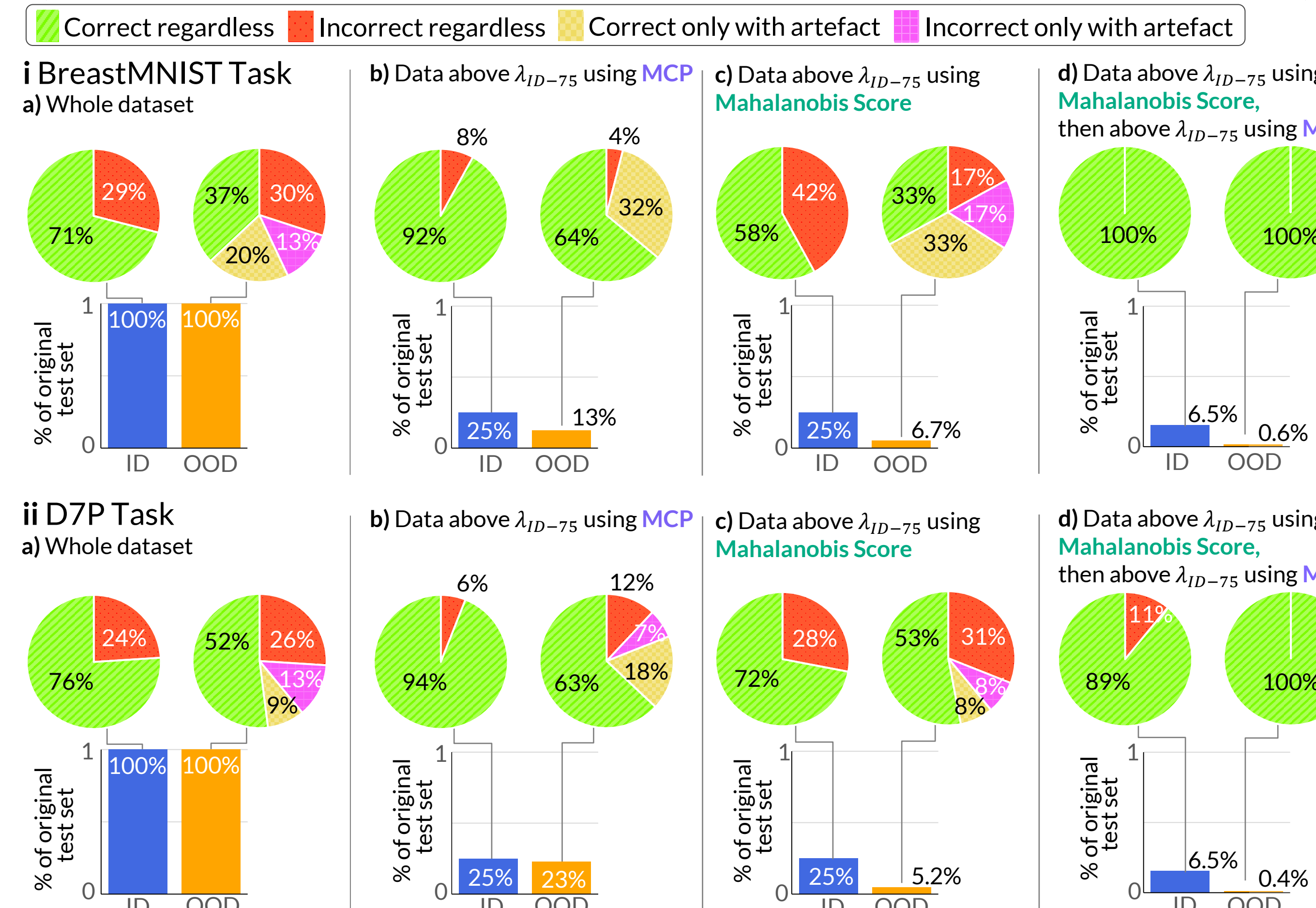
🔑 **Key Takeaways**
- Out-of-distribution detection ≠ Failure detection
- **Confidence-based** methods typically outperform **feature-based** methods at **Failure Detection**
- **Feature-based** methods typically outperform **confidence-based** methods at **OOD Detection**

## 5. Combining OOD detection methods to mitigate their weaknesses

### a) Method for studying diagnoses above $S(x)$ threshold

1. Sort all OOD images into one of four categories

| Label | Image **with** artefact | Counterfactual Image **without** artefact |
|---|---|---|
| Correct regardless | Correct | Correct |
| Incorrect regardless | Incorrect | Incorrect |
| Correct only with artefact | Correct | Incorrect |
| Incorrect only with artefact | Incorrect | Correct |

2. Calculate Scoring function for OOD method.
3. Calculate 75 percentile on ID test data $\lambda_{ID\_75}$.
4. Remove the diagnoses below the threshold.

Frequency — Discard ← → Keep diagnoses — $\lambda_{ID\_75}$ — Scoring function $S(x)$

### b) Results for a confidence and a feature-based method (and a combination)

Legend: Correct regardless | Incorrect regardless | Correct only with artefact | Incorrect only with artefact

**i BreastMNIST Task**

a) Whole dataset: ID 71% / 29%; OOD 37% / 30% / 20% / 13%
% of original test set: ID 100%, OOD 100%

b) Data above $\lambda_{ID\_75}$ using **MCP**: 92% / 8%; 64% / 32% / 4%
% of original test set: ID 25%, OOD 13%

c) Data above $\lambda_{ID\_75}$ using **Mahalanobis Score**: 58% / 42%; 33% / 33% / 17% / 17%
% of original test set: ID 25%, OOD 6.7%

d) Data above $\lambda_{ID\_75}$ using **Mahalanobis Score, then above $\lambda_{ID\_75}$ using MCP**: 100%; 100%
% of original test set: ID 6.5%, OOD 0.6%

**ii D7P Task**

a) Whole dataset: 76% / 24%; 52% / 26% / 13% / 9%
% of original test set: ID 100%, OOD 100%

b) Data above $\lambda_{ID\_75}$ using **MCP**: 94% / 6%; 63% / 18% / 12%
% of original test set: ID 25%, OOD 23%

c) Data above $\lambda_{ID\_75}$ using **Mahalanobis Score**: 72% / 28%; 53% / 31% / 8%
% of original test set: ID 25%, OOD 5.2%

d) Data above $\lambda_{ID\_75}$ using **Mahalanobis Score, then above $\lambda_{ID\_75}$ using MCP**: 89% / 11%; 100%
% of original test set: ID 6.5%, OOD 0.4%

### c) In-depth insights

- **Correct only with artefact** diagnoses can inflate failure detection AUROC.

ID — 92% — Accuracy: 92%
OOD — 64% / 32% — Accuracy: 64+32=96% > 92%; Correct regardless: 64% < 92%

→ *Could give false sense of security*

- OOD detection methods with high OOD AUROC, but which also cause a higher risk of incorrect diagnoses, may not make neural nets more trustworthy.
- Using multiple detectors leads to more discarded diagnoses, but the remaining diagnoses are more trustworthy.

🔑 **Key Takeaways**
- **Combining** confidence & feature-based methods in a pipeline can **mitigate their respective weaknesses.**