

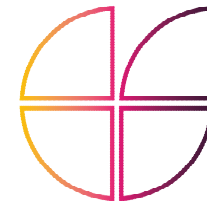
CVPR
JUNE 3-7, 2026



DENVER
COLORADO



IBME
INSTITUTE OF BIOMEDICAL ENGINEERING



DEPARTMENT OF
ENGINEERING
SCIENCE



The Invisible Gorilla Effect In Out-of-distribution Detection

Harry Anthony, Ziyun Liang, Hermione Warr, Konstantinos Kamnitsas

CVPR Highlight 2026

Let's begin with an exercise...

Count the number of basketball passes between **players in white shirts** and ignore passes between in black shirts.

<https://www.youtube.com/watch?v=vJG698U2Mvo>



How many did you count ?

Correct answer is 15

But did you spot the person in a **Gorilla Costume**?



Invisible Gorilla Experiment

In Psychology



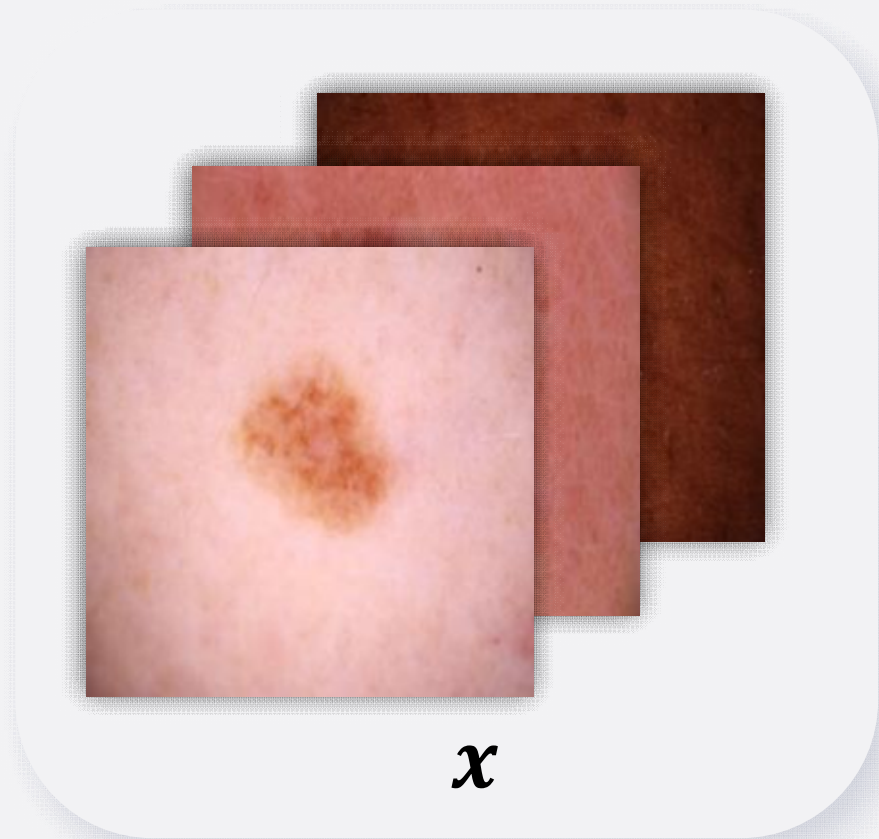
N.Tisza (2023)

Inattentive Blindness

Individuals focusing on a primary task can ignore unexpected stimuli due to a lack of attention.

Unexpected stimuli are **more likely to be seen** when they **closely resemble** the attended target.

Out-of-distribution detection



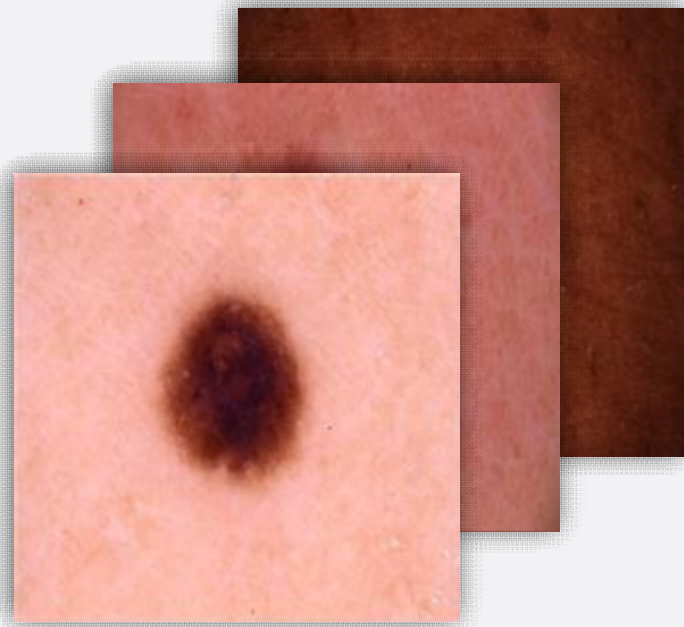
Benign
Malignant
Benign

y

Training
Data

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

Out-of-distribution detection



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

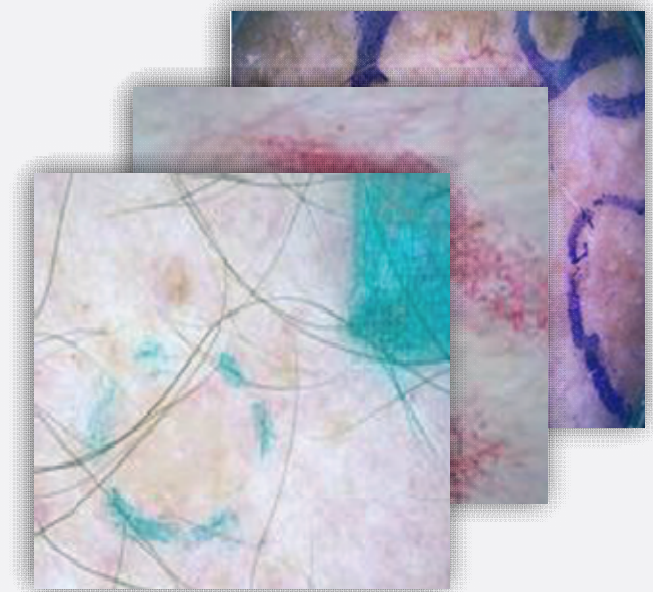


Covariate
Shift

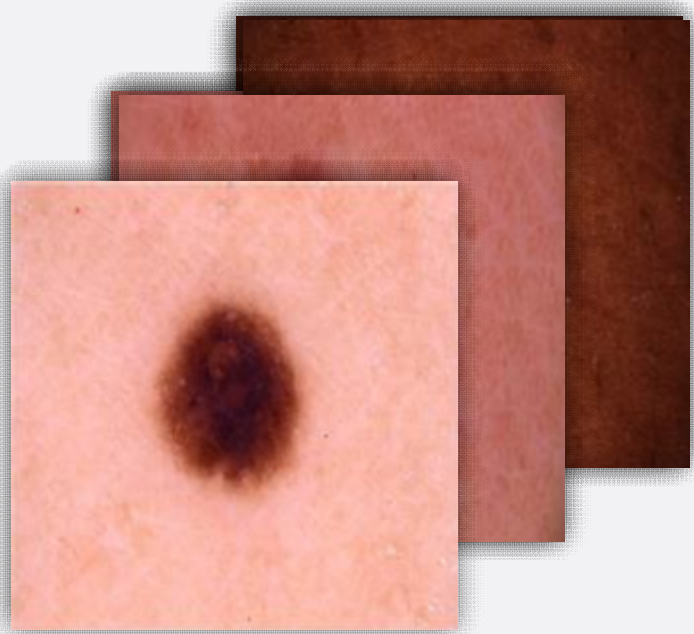
Same Classes

$$p_{ID}(y|x) = p_{OOD}(y|x)$$

$$p_{ID}(x) \neq p_{OOD}(x)$$

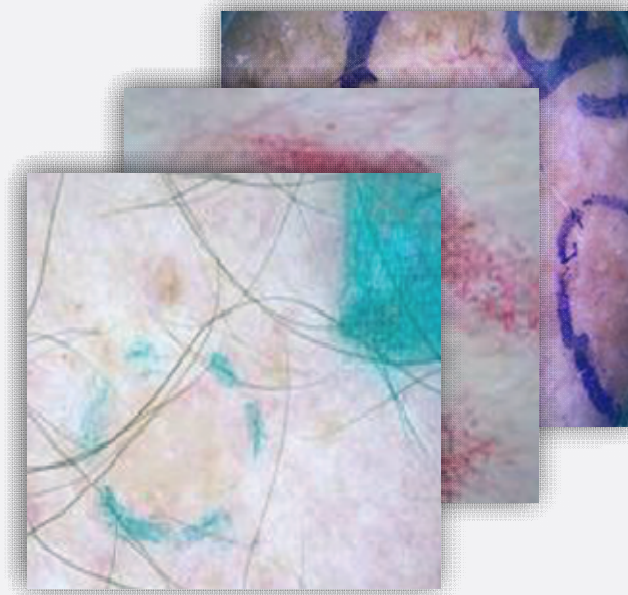


$(x, y) \not\sim \mathcal{D}_{train}$
Out-of-distribution
(OOD)



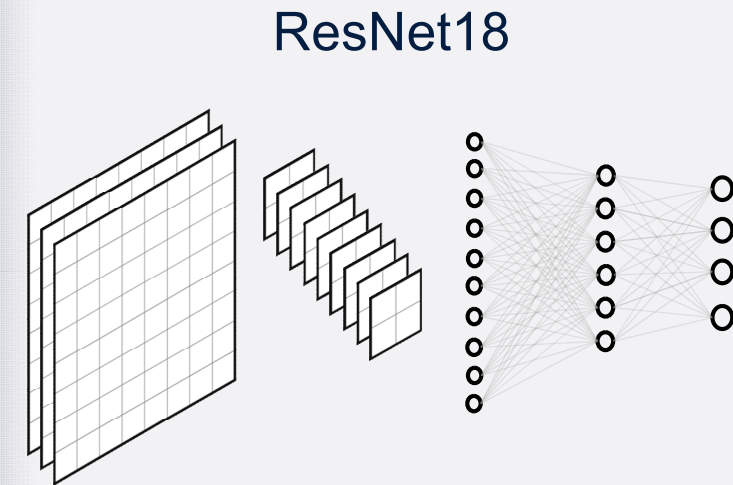
$(x, y) \sim \mathcal{D}_{train}$
In-distribution

90.8%
Diagnostic
Accuracy



$(x, y) \not\sim \mathcal{D}_{train}$
Out-of-distribution

74.3%
Diagnostic
Accuracy





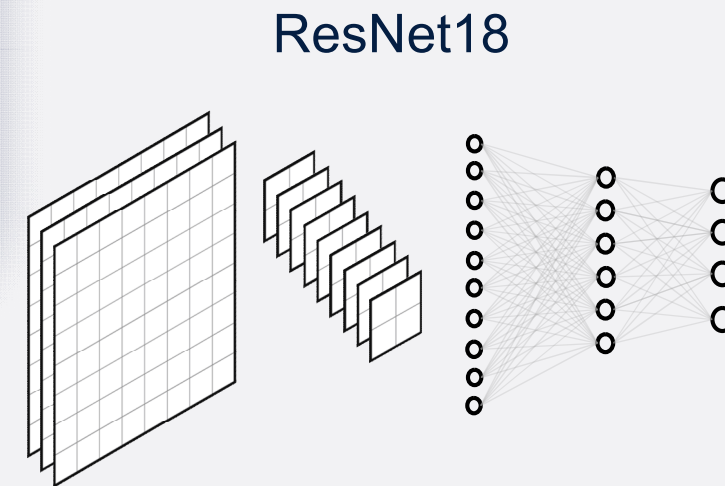
$(x, y) \sim \mathcal{D}_{train}$
In-distribution

90.8%
Diagnostic
Accuracy



$(x, y) \not\sim \mathcal{D}_{train}$
Out-of-distribution

74.3%
Diagnostic
Accuracy



Motivates Out-of-distribution detection methods

Aim of flagging diagnoses on OOD data as unreliable

Finding biases in

Out-of-distribution detection

OOD detection methods that perform well in detecting one OOD artefact can perform poorly in another.

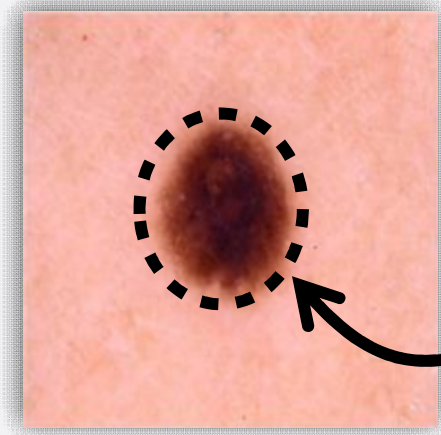
Very limited understanding of what causes these differences.

We introduce an **unreported bias in OOD detection**.

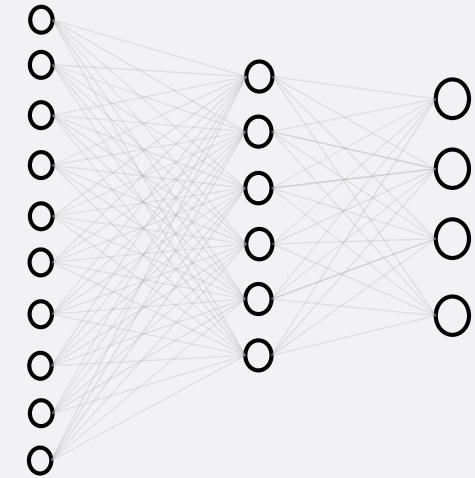
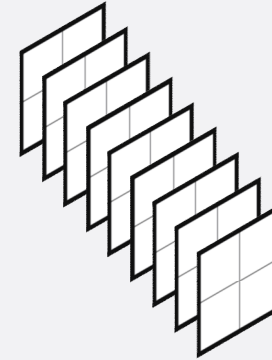
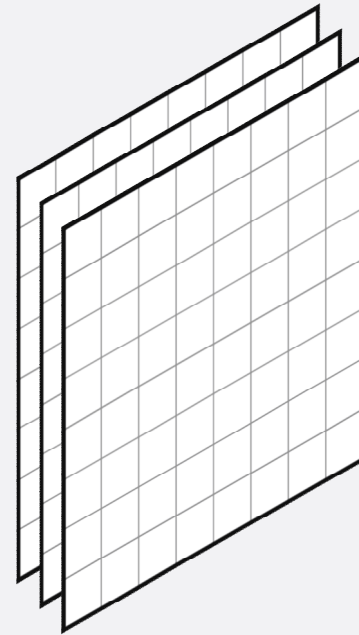
Finding biases in

Out-of-distribution detection

Model's Region of Interest (ROI)

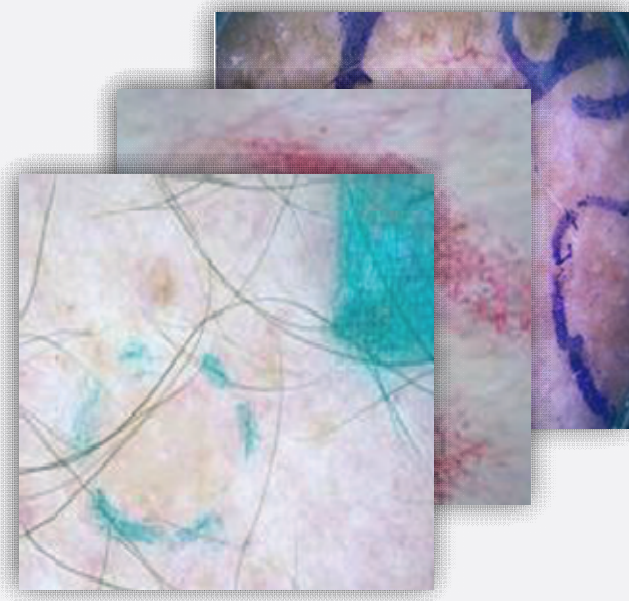


Mean RGB
(176,116,77)

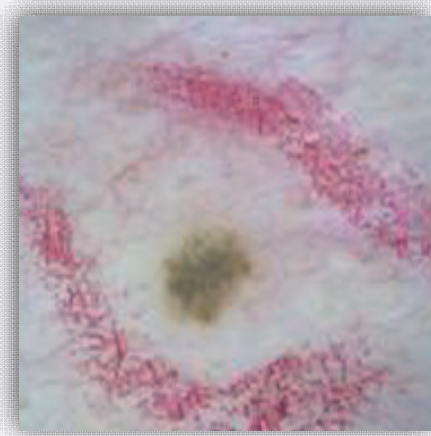


Finding biases in

Out-of-distribution detection



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



Red



Green



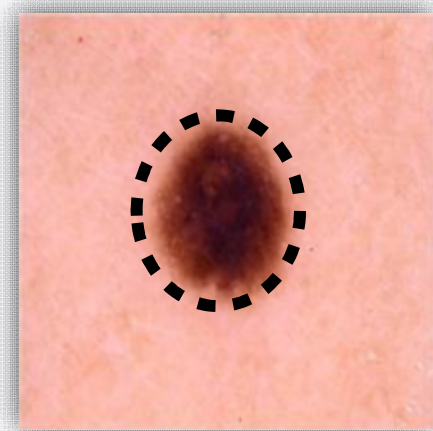
Purple



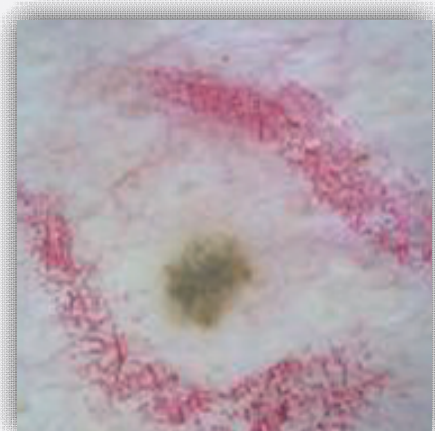
Black

Finding biases in **Out-of-distribution detection**

Model's ROI
(176,116,77)

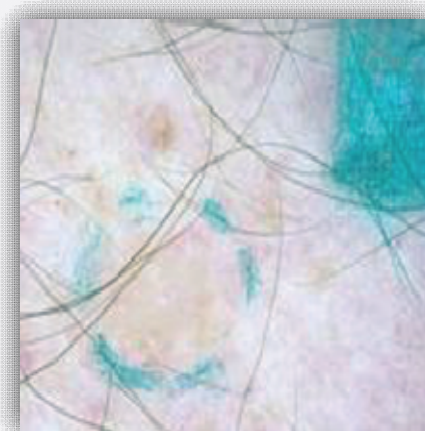


Similar to
ROI



Red

Dissimilar
to ROI



Green



Purple



Black

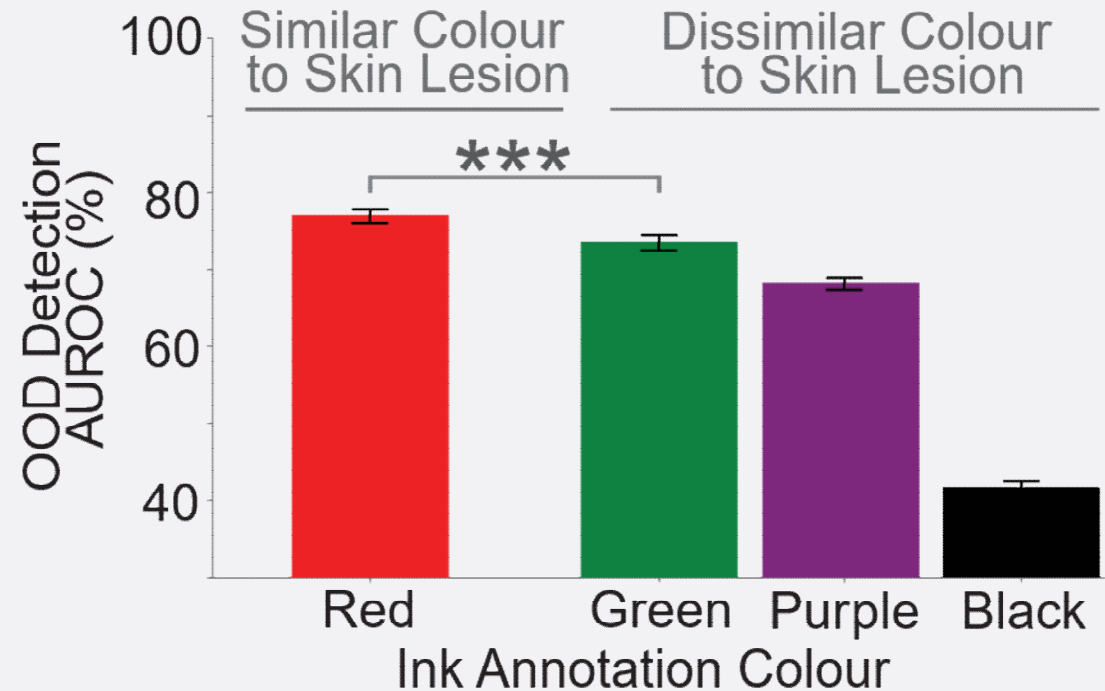


Euclidean RGB Distance

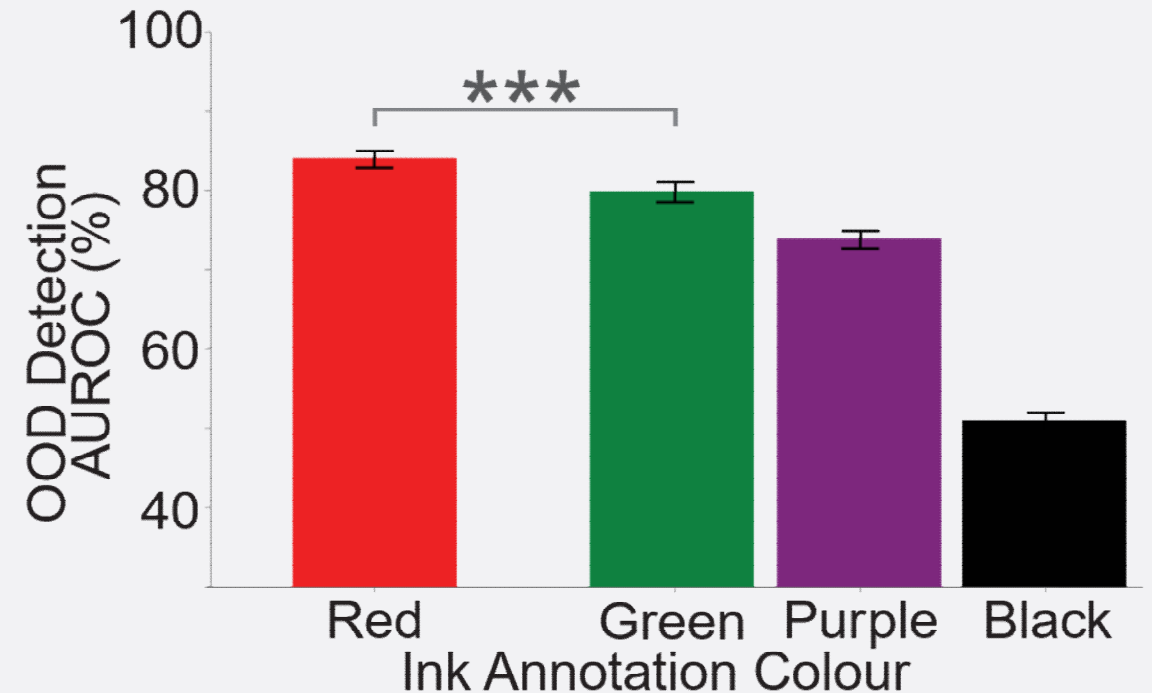
Finding biases in

Out-of-distribution detection

OOD Detection Method: Mahalanobis Score

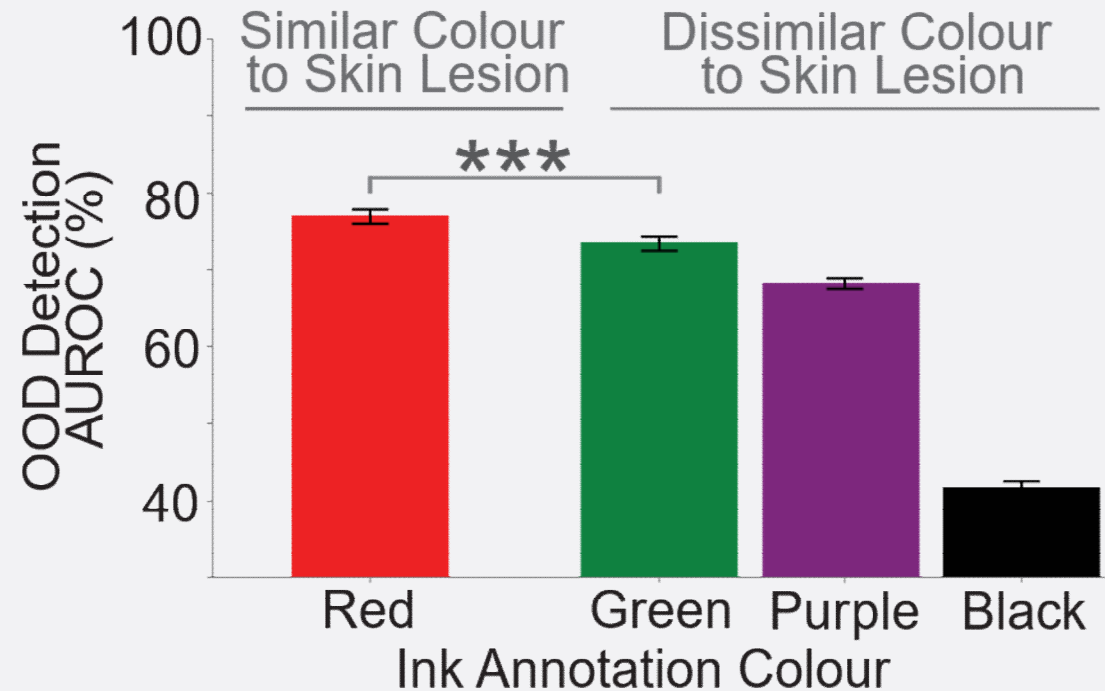


OOD Method: Flow-based Model (RealNVP)

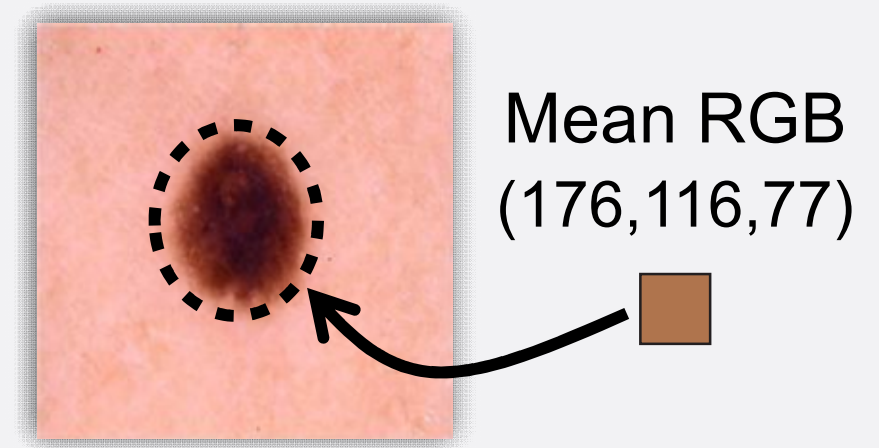


Finding biases in **Out-of-distribution detection**

OOD Detection Method: Mahalanobis Score



OOD detection performance **improves** when the OOD artefact is **visually similar** to model's Region of Interest.

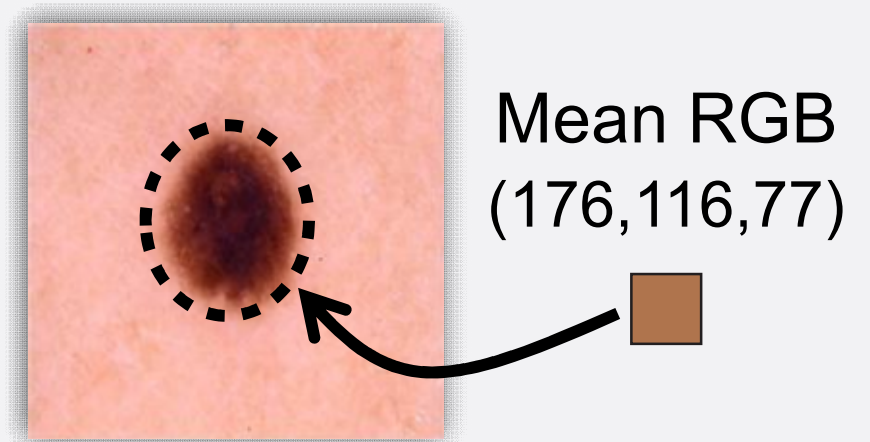


Finding biases in **Out-of-distribution detection**

Unexpected stimuli are **more likely to be seen** when they **closely resemble** the attended target.



OOD detection performance **improves** when the OOD artefact is **visually similar** to model's Region of Interest.



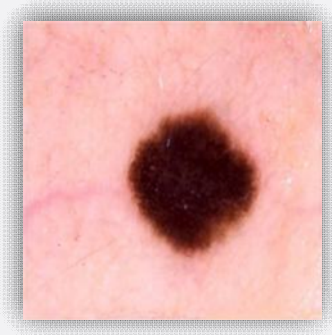
Finding biases in

Out-of-distribution detection

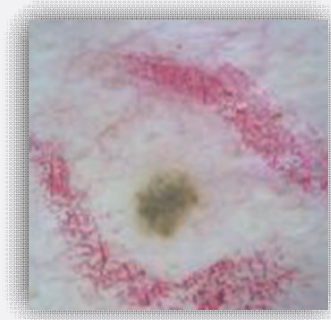
The Invisible Gorilla Effect

OOD detection performance **improves**
when the OOD artefact is
visually similar
to model's Region of Interest.

Community Assumptions



Training
Data



Near
OOD



Far
OOD



The prevailing assumption:

The more similar OOD data is to the training data, the harder it is to detect

Community Assumptions

Our findings challenge that assumption.

We reveal a bias case where the **opposite is true**.

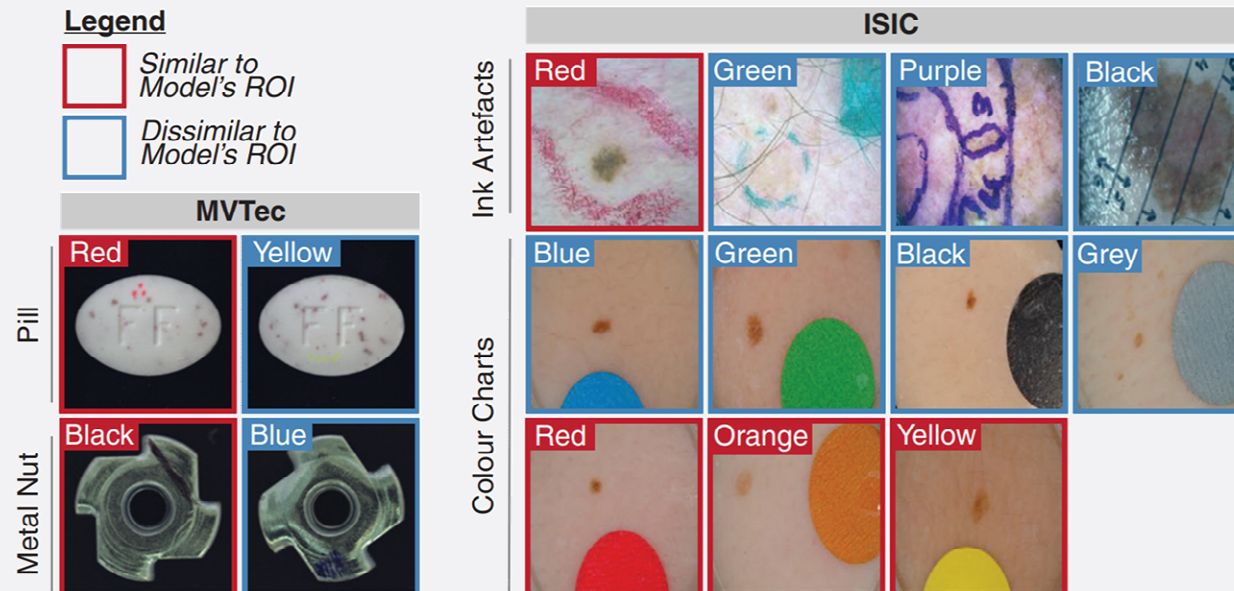
Suggests that OOD task difficulty is not a simple function of global similarity to the training data, but depends on:

- What the model attends to
- How that interacts with the OOD detection method

Contributions of this work

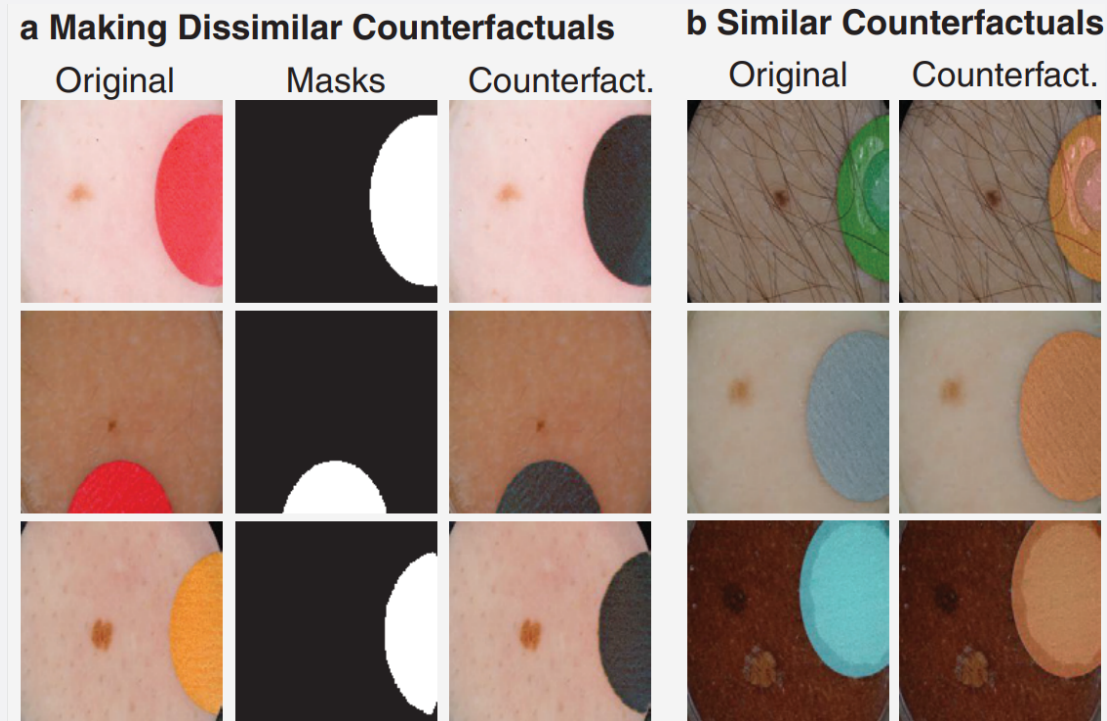
- 1 Create Benchmarks and Conduct a Large-Scale Evaluation
 - 7 **Benchmarks** made public

11,355 images



Contributions of this work

- 1 Create Benchmarks and Conduct a Large-Scale Evaluation
 - 7 **Benchmarks made public**
 - Colour-swapped Counterfactuals



Contributions of this work

- 1 Create Benchmarks and Conduct a Large-Scale Evaluation
 - 7 **Benchmarks made public**
 - Colour-swapped Counterfactuals
 - 3 Model Architectures
 - 40 OOD methods
 - 3,795 Hyperparameters
 - 25 Random Seeds

Contributions of this work

- 1 Create Benchmarks and Conduct a Large-Scale Evaluation
- 2 Identify Trends
- 3 Hypothesise the Causes for the Effect.
- 4 Use Insights to Inform Mitigation Strategies

Thank you for listening

I welcome any questions

Paper available at:



Data available at:



Full Presentation at:

