

Artificial Intelligence Safety for Medical Sciences

Varun Ojha

AI Theme Leader (Edge AI Hub)

Senior Lecturer in Artificial Intelligence
School of Computing, Newcastle University
varun.ojha@Newcastle.ac.uk

Date:
21st November 2024

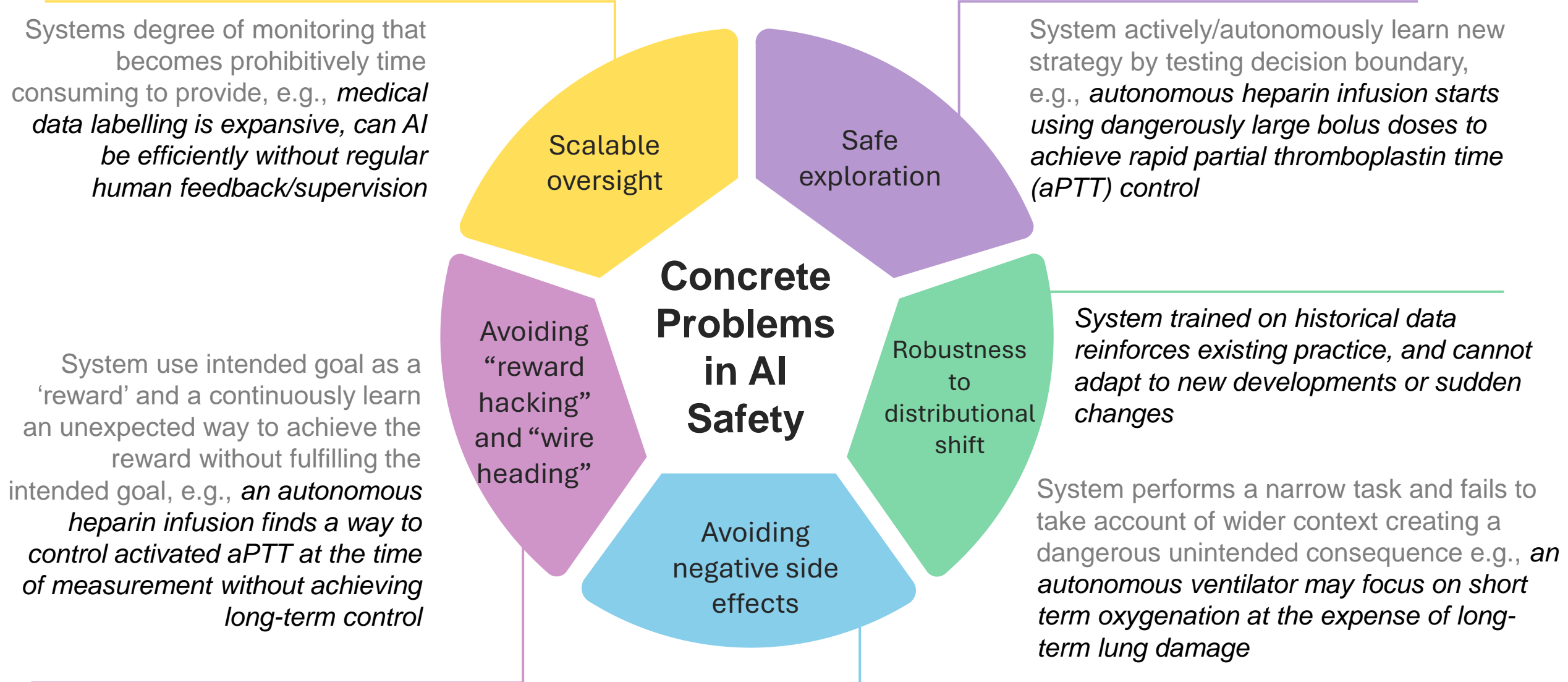


Place:
Newcastle University

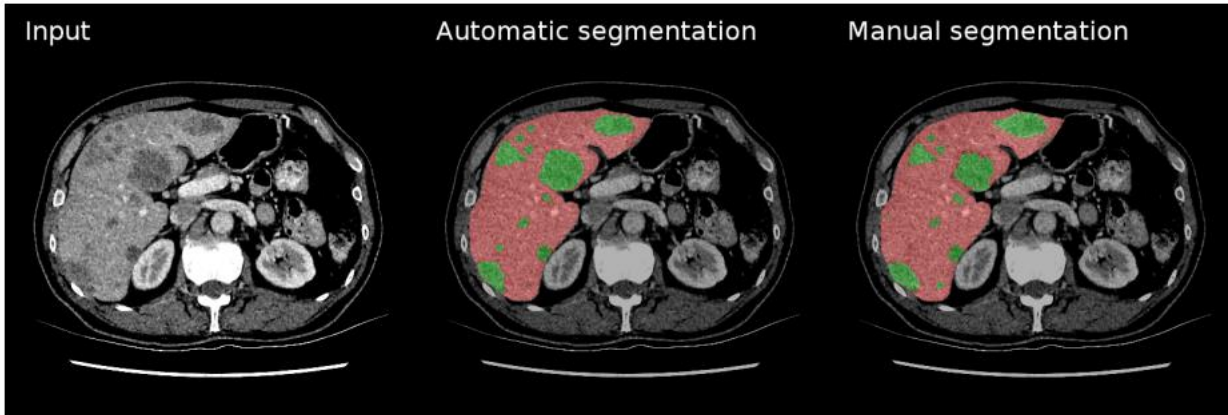
Artificial Intelligence Safety Challenges in Medical Domain

Amodei et al. (2016). Concrete problems in AI safety

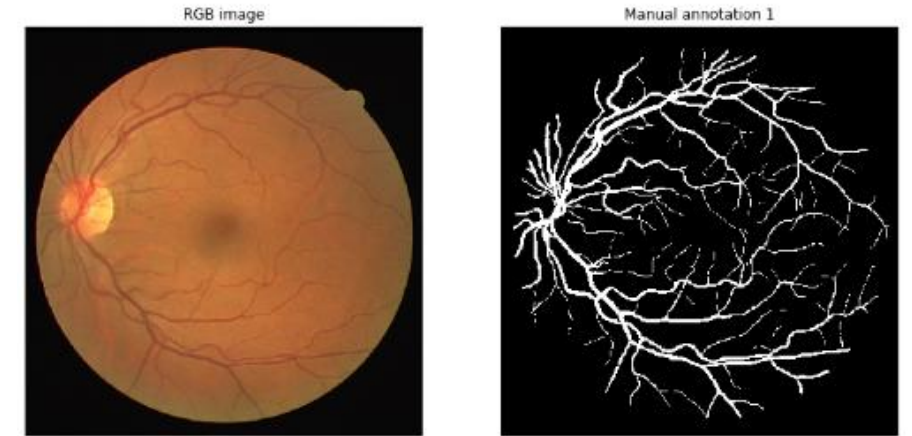
Challen et al. (2019). Artificial intelligence, bias and clinical safety. *BMJ quality & safety*



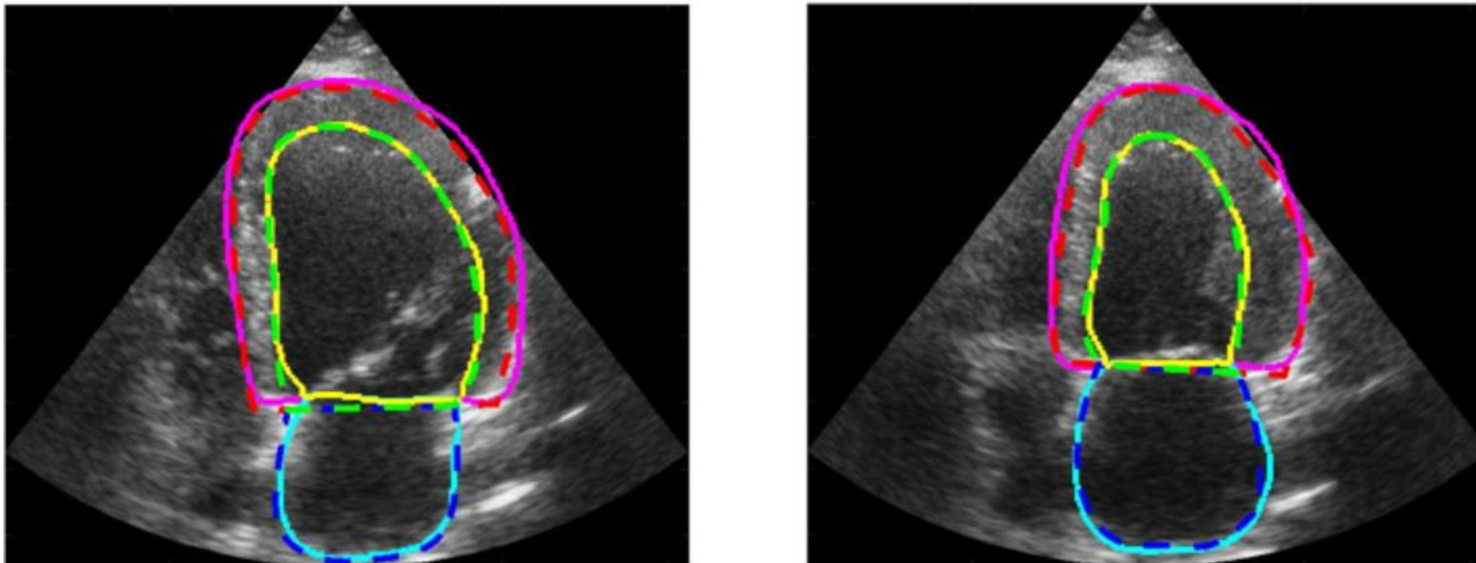
Artificial Intelligence in Medical Domain



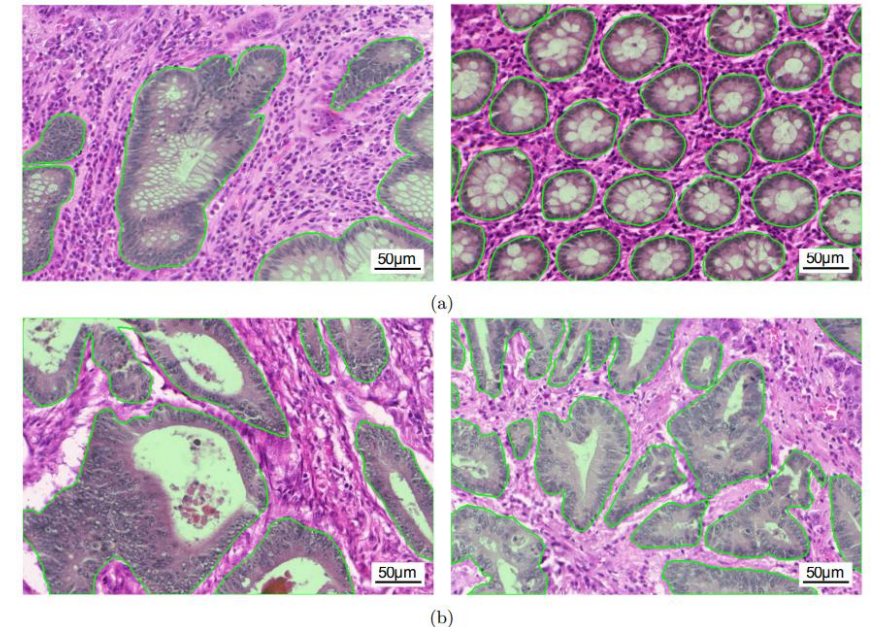
Liver Tumor Segmentation from CT scans



Digital Retinal Images for Vessel Extraction



Cardiac Acquisitions for Multi-structure Ultrasound Segmentation

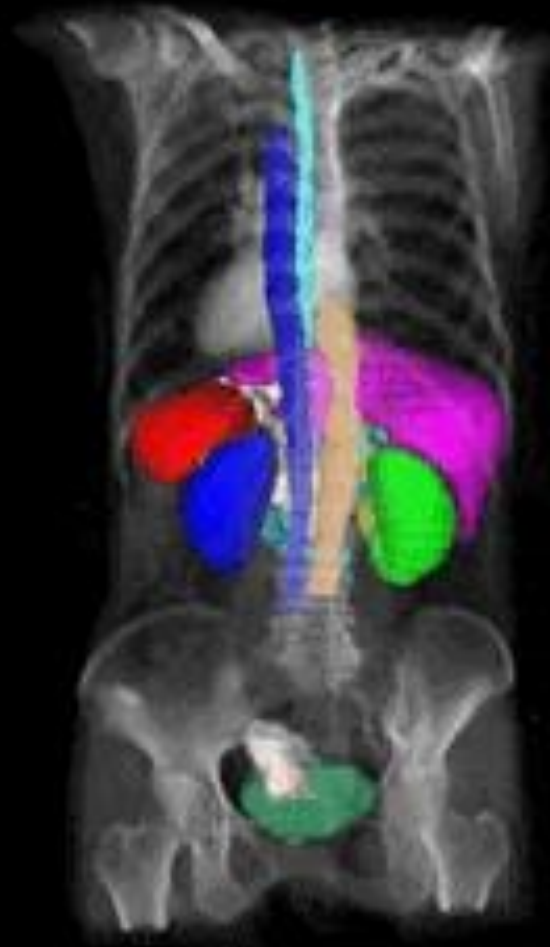


Gland Segmentation in Colon Histology Images

Multi-Modality Abdominal Multi-Organ Segmentation

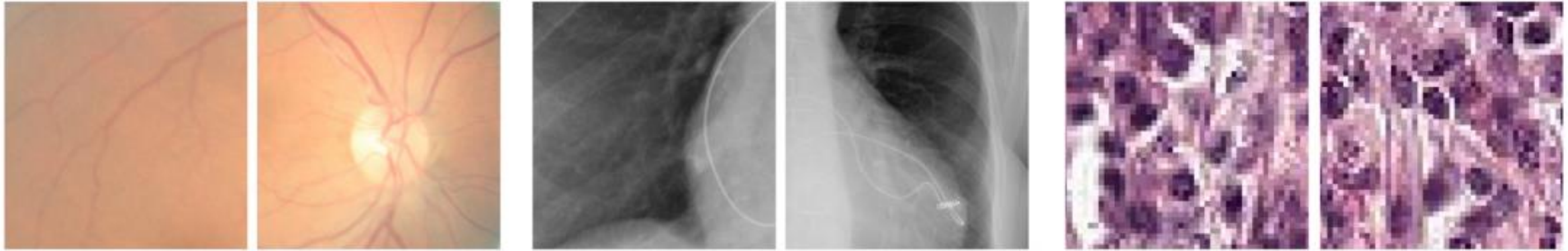
AMOS2022

- Prostate/Uterus
- Bladder
- Duodenum
- Left Adrenal gland
- Right Adrenal gland
- Pancreas
- Postcava
- Aorta
- Stomach
- Liver
- Esophagus
- Gall Bladder
- Left Kidney
- Right Kidney
- Spleen
- Blackground

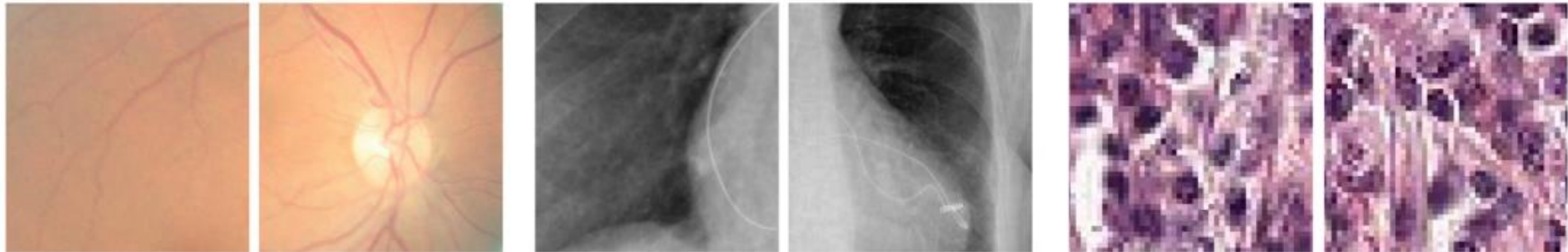


AI Safety Concern in Medical Image Analysis

Original
Image



Adversarially
Modified



Ophthalmology

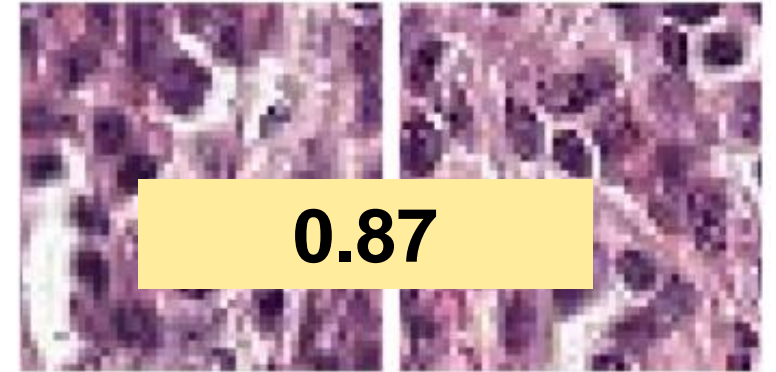
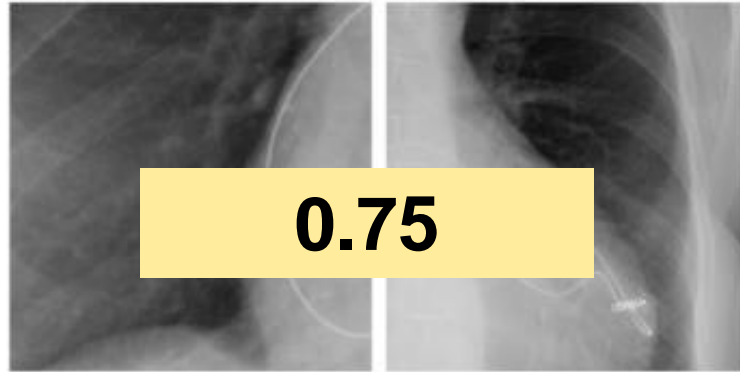
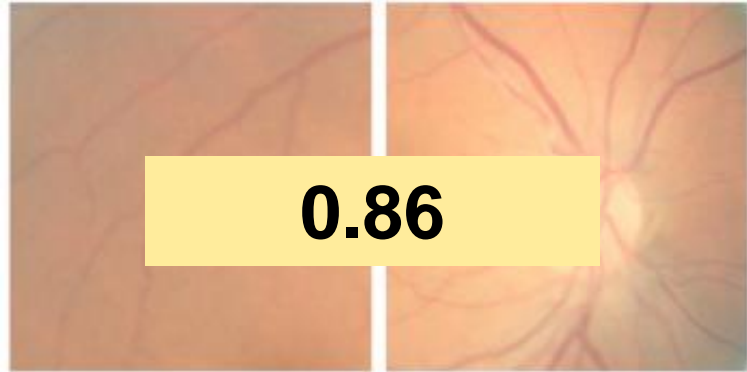
Radiology

Pathology

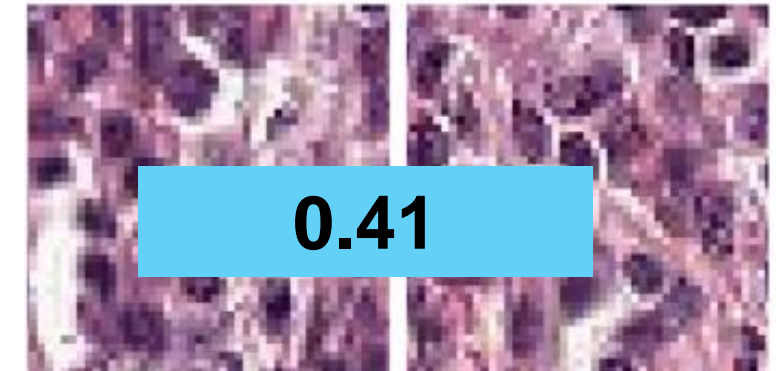
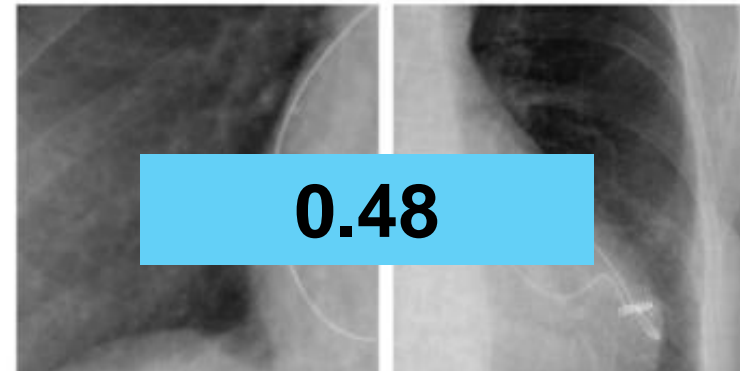
Unperceived changes in images make misclassification

AI Safety Concern in Medical Image Analysis

Original Image



Adversarially Modified



Ophthalmology

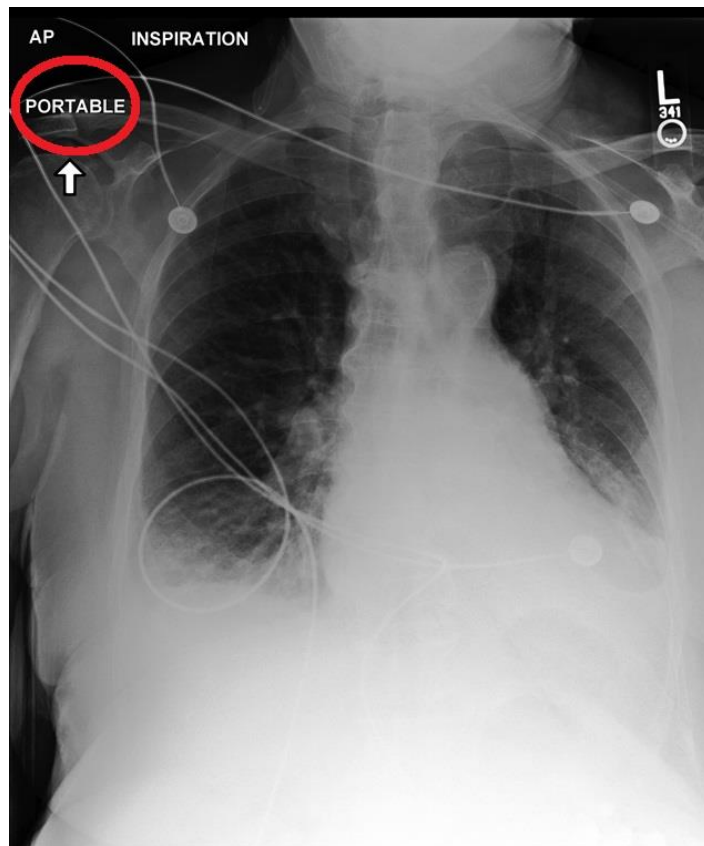
Radiology

Pathology

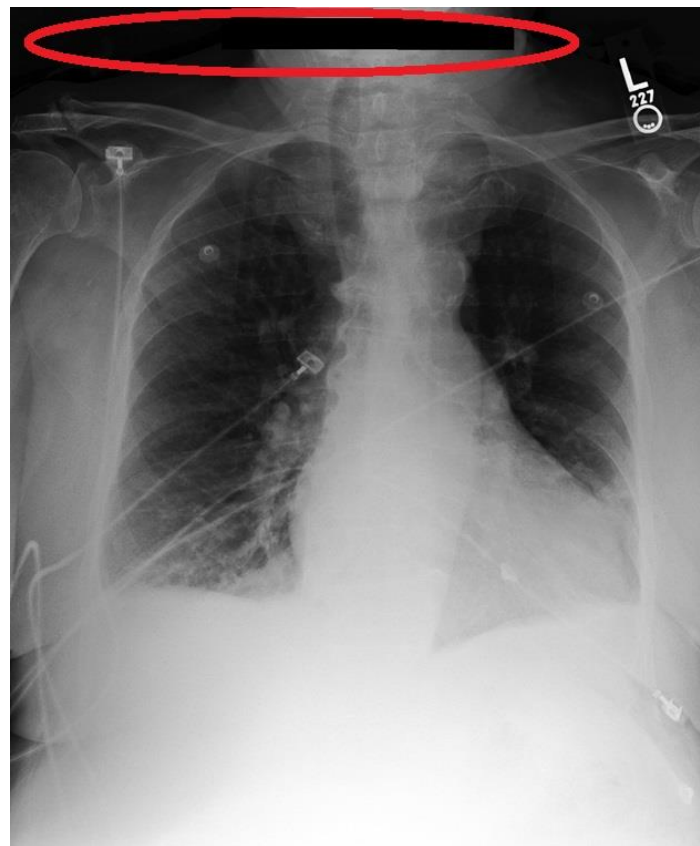
Accuracy of detection decreases even on an unperceived modification

Spurious Correlation in Medical Training Data

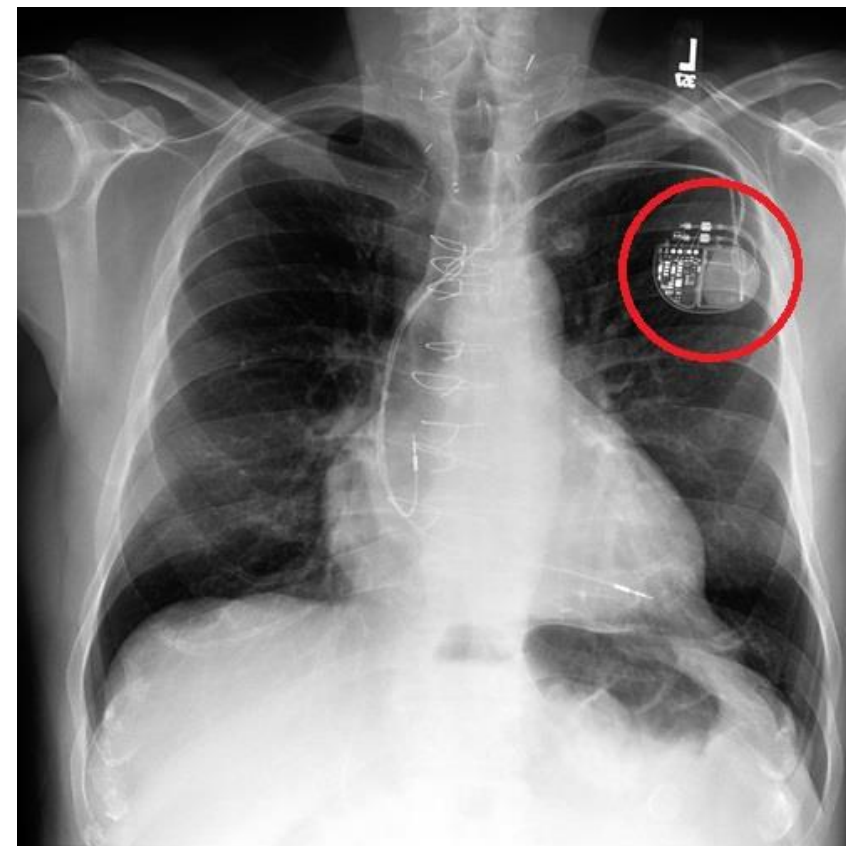
Spurious correlation occur in medical training data where diagnosis results are affected by variables (e.g., Hospital tags, Stripes, Medical devices) that are not related to the diagnostic information being predicted. This phenomenon leads to misleading interpretations.



Hospital tags



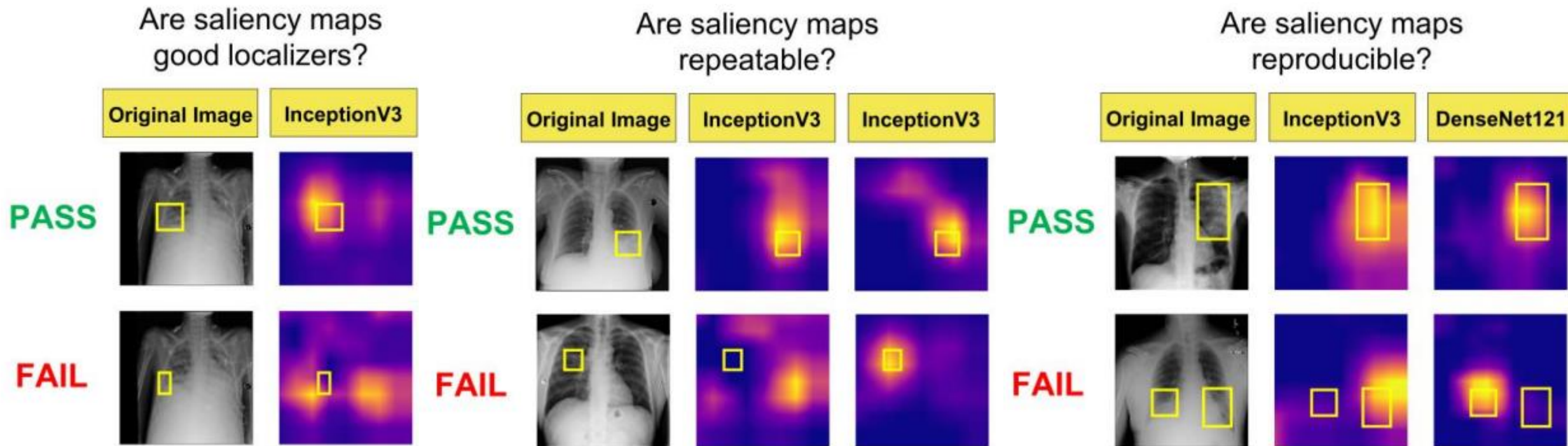
Stripes



Medical devices

Artificial Intelligence (Un)Trustworthiness Performance

Saliency Map is way to explain the prediction of AI models on Medical Diagnosis



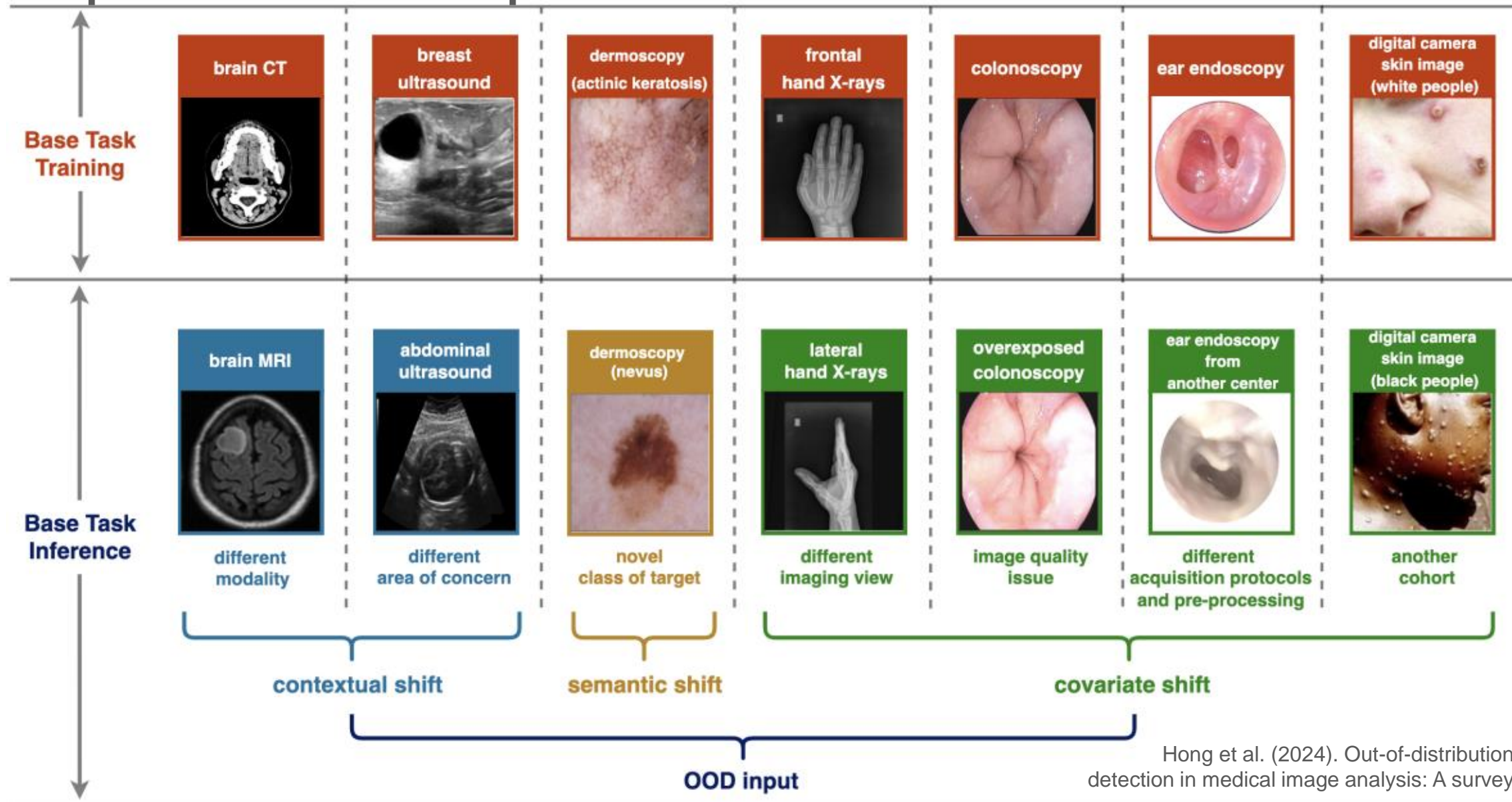
PASS : Only in most ideal cases AI works

FAIL : In most cases AI fails to locate, repeat, and reproduce results

Arun et al. (2021). Assessing the trustworthiness of saliency maps for localizing abnormalities in medical imaging. Radiology: Artificial Intelligence

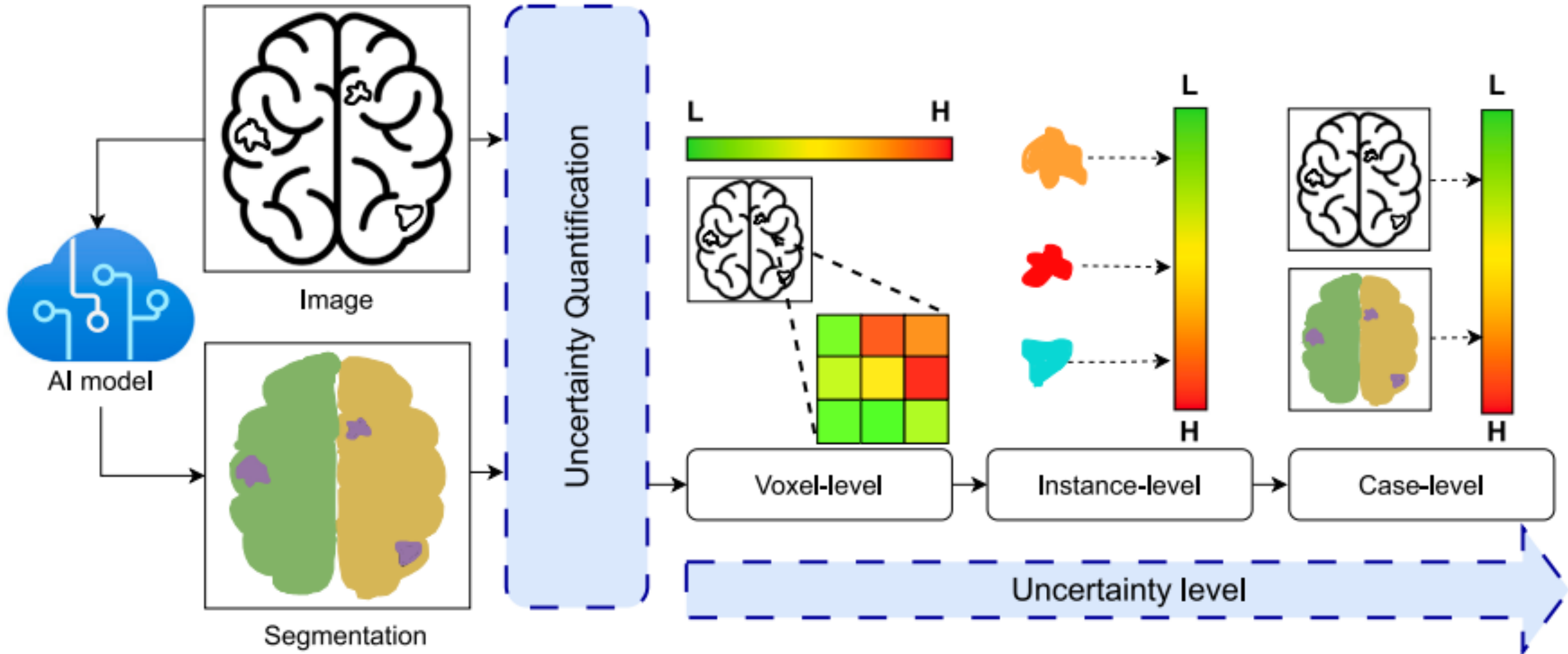
Risk of Training Data Error

Examples of out of distribution problems



Risk of Training Error and AI (Un)Trustworthiness/Uncertainty?

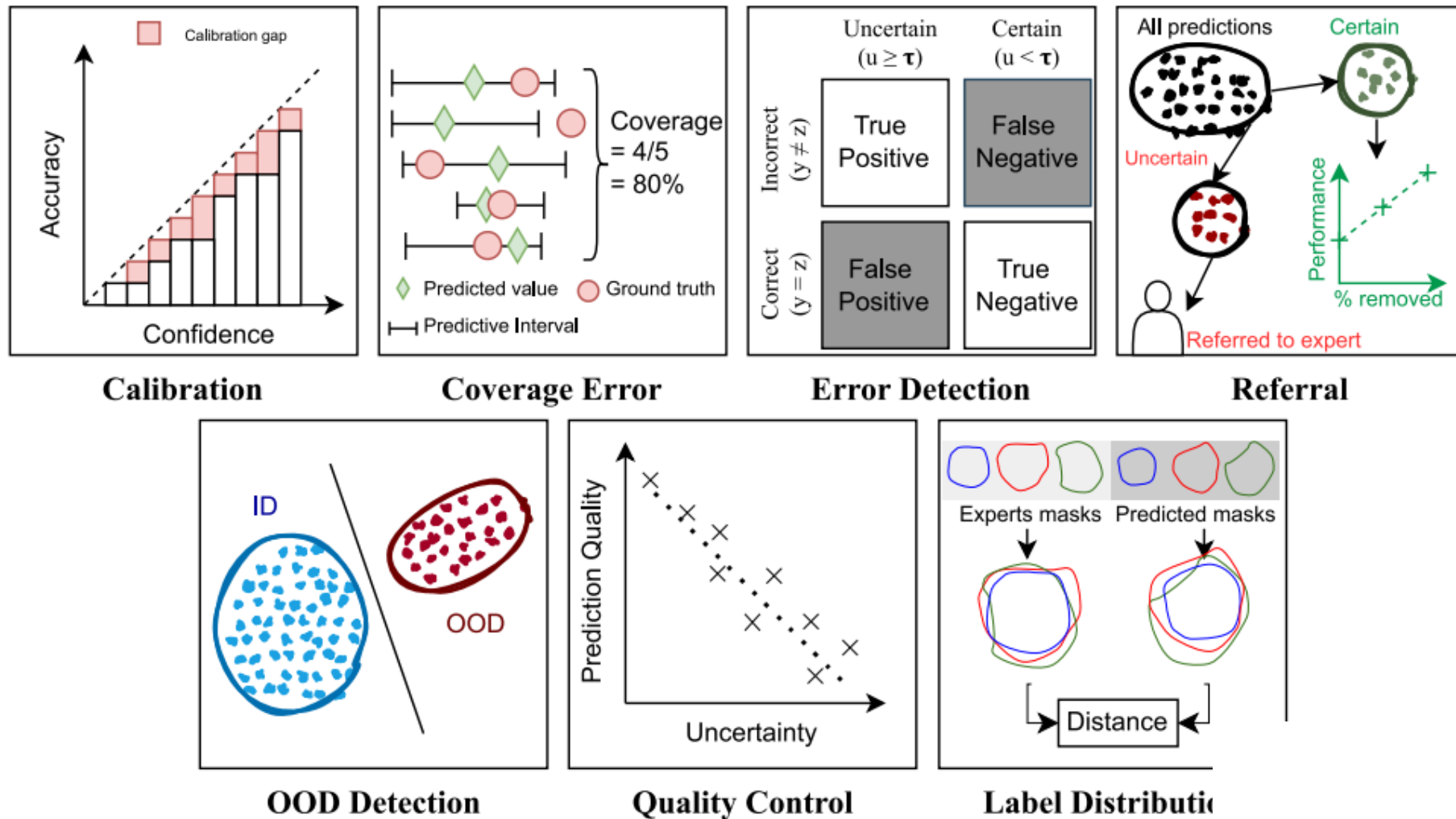
Trained AI Models inference suffers from inference/test images being out of training distribution



Mitigating Risk of AI Uncertainty / Scalable Oversight

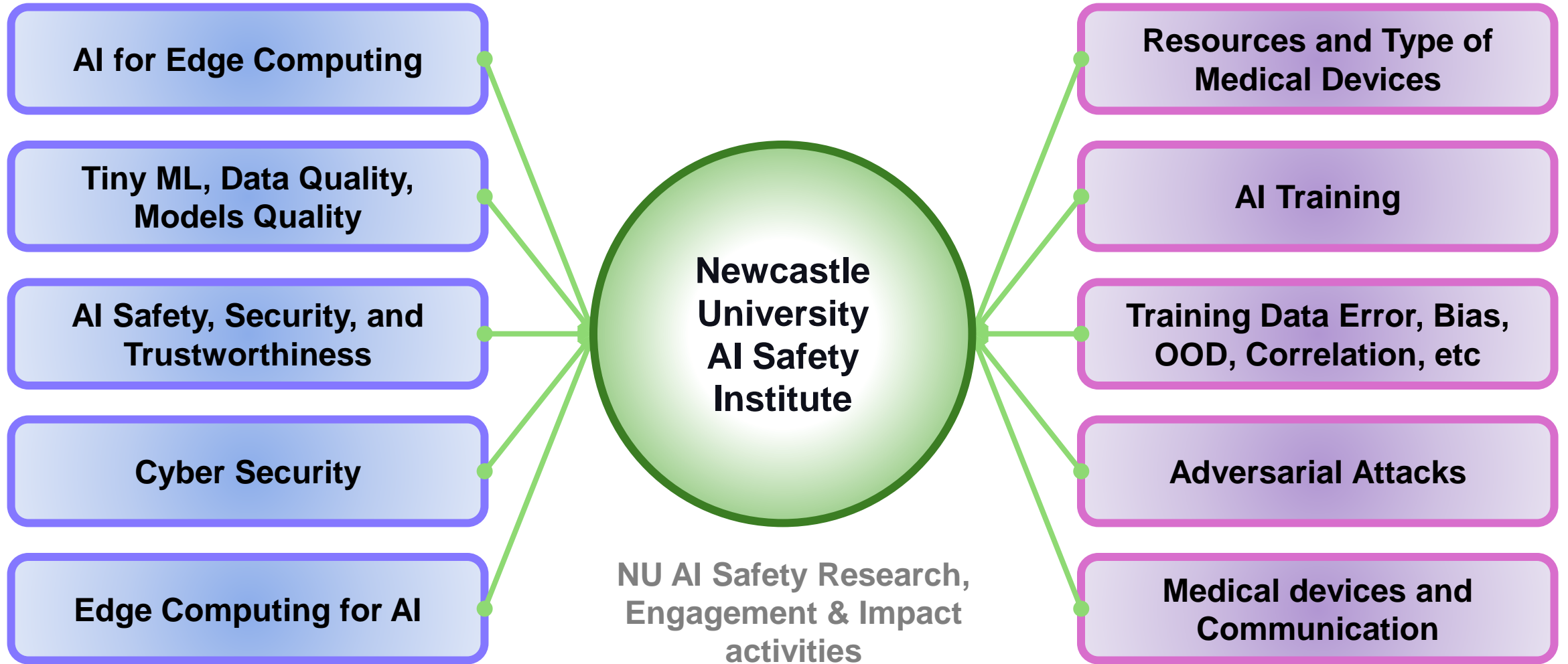
Trained AI Models inference suffers from inference/test images being out of training distribution

Designing AI to actively bring (Observer) **Human in the Loop** and return their feedback to make AI systems better



Edge AI Hub → NU AI Safety Institute → Faculty of Medical Sciences

NU AI Safety Institute will deliver for FMS via Edge AI Hub and University Wide Research Skills



Edge AI Hub AI Safety Research, Skills and Leadership

FMS's AI Safety Research Needs



**AI Safety
Institute**

Powered by



National Edge AI Hub

Get in touch



Address

Urban Sciences Building, 1 Science Square,
Newcastle upon Tyne NE4 5TG, UK



Email

varun.ojha@newcastle.ac.uk



Web

<https://ojhavk.github.io/>

