

Interpretable Medical AI with Vision-Language Alignment

Shantanu Ghosh, PhD Prospectus

Jan 28, 2026



**BATMAN
LAB**

The goal: **Detect the Systematic Mistakes**

Why?

To design better mitigation strategies(data collection, reweighting, architecture/training changes etc) to enhance **robustness**.

Dissecting Systematic Mistakes?



Class: **Waterbirds**
Background: **Water**



Class: **Waterbirds**
Background: **Land**

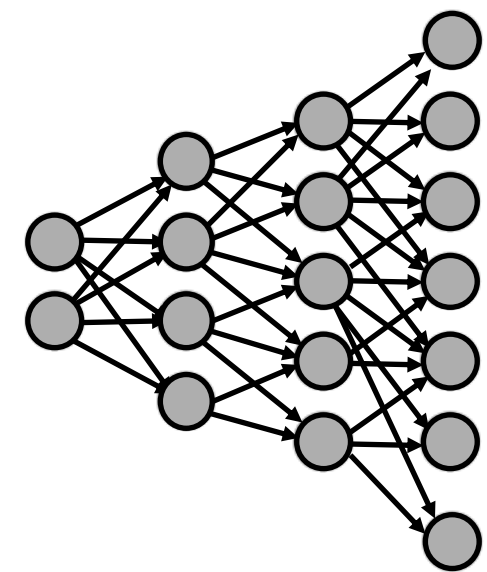


Class: **Landbirds**
Background: **Water**

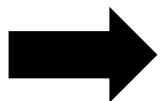


Class: **Landbirds**
Background: **Land**

Dissecting Systematic Mistakes?



ResNet50
Mean Accuracy:
88.6%



Class: **Waterbirds**
Background: **Water**
Accuracy: **94.2%**



Class: **Waterbirds**
Background: **Land**
Accuracy: **68.8%**

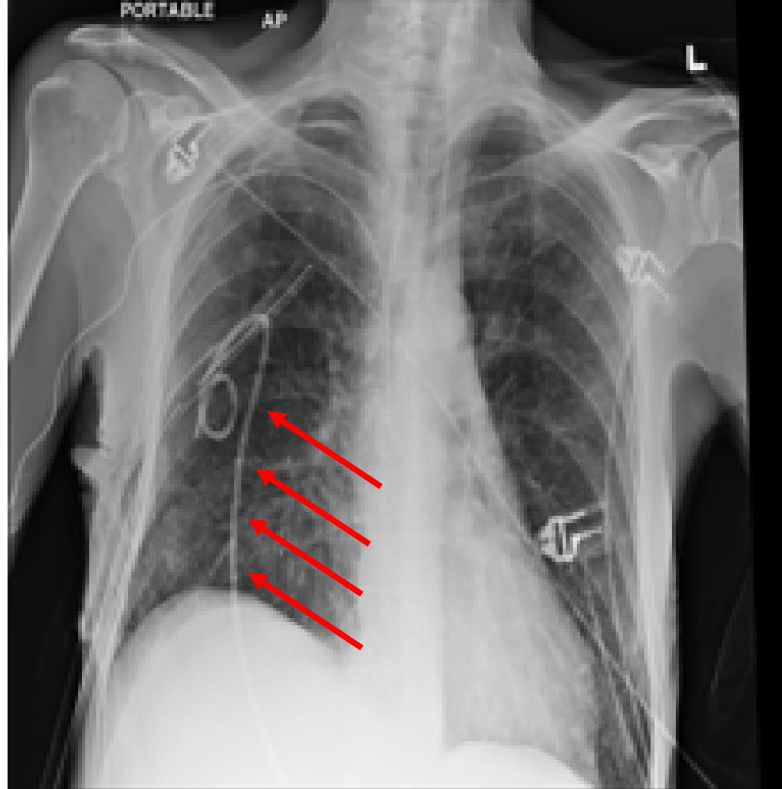
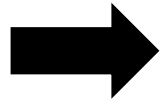
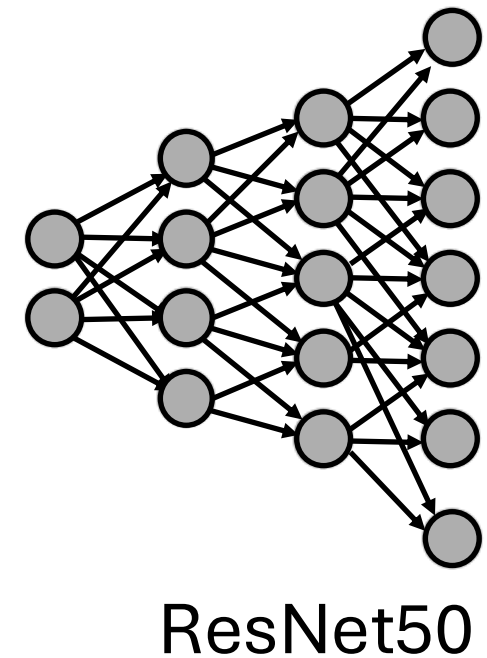


Class: **Landbirds**
Background: **Water**
Accuracy: **80.2%**



Class: **Landbirds**
Background: **Land**
Accuracy: **99.6%**

Dissecting Systematic Mistakes?



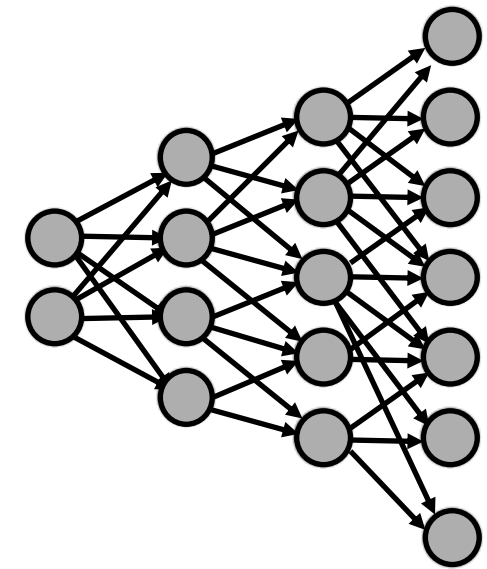
Class: **Pneumothorax**
Correlation: **Chest tube**



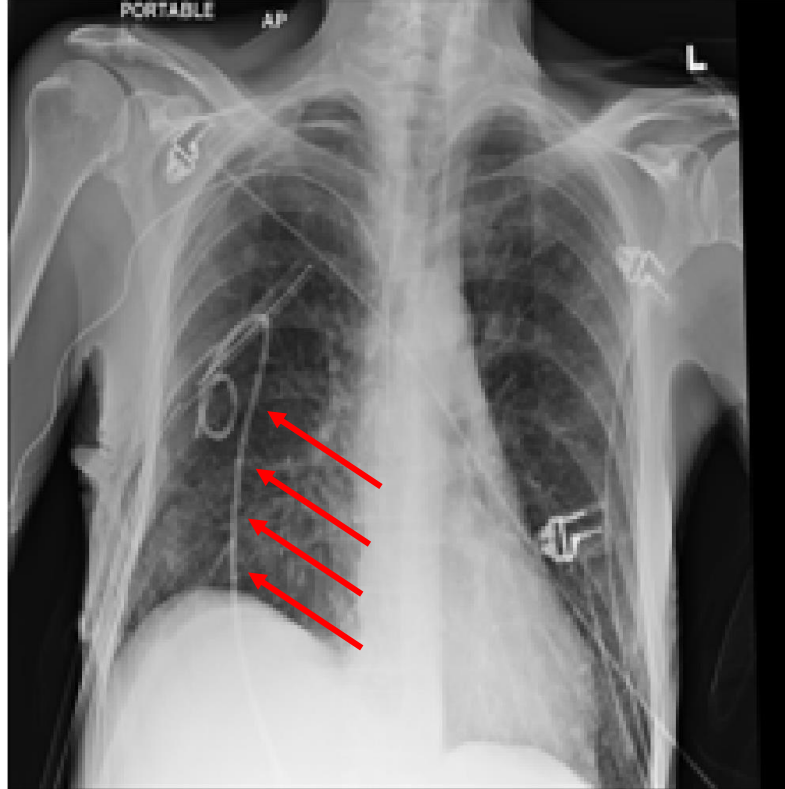
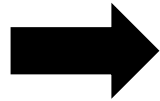
Class: **Pneumothorax**
Correlation: **Chest tube**

Dissecting Systematic Mistakes?

Mean Accuracy (Pneumothorax): ~**70%**



ResNet50

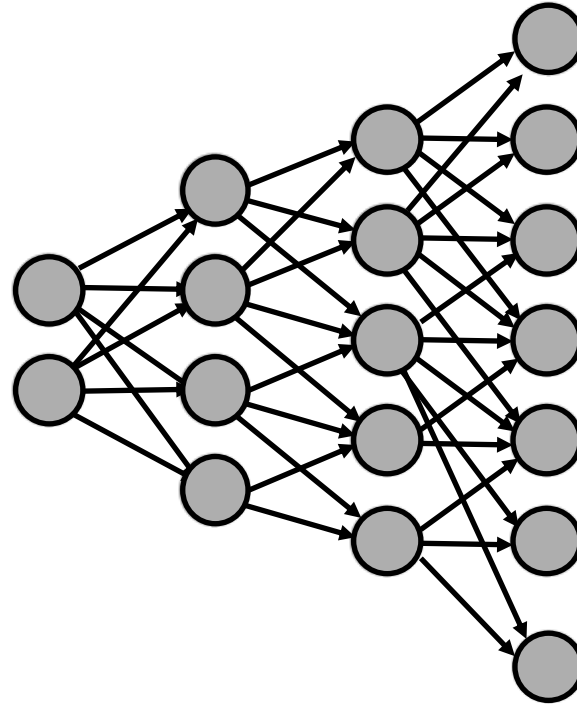


Class: **Pneumothorax**
Correlation: **Chest tube**
Accuracy: **90.4%**



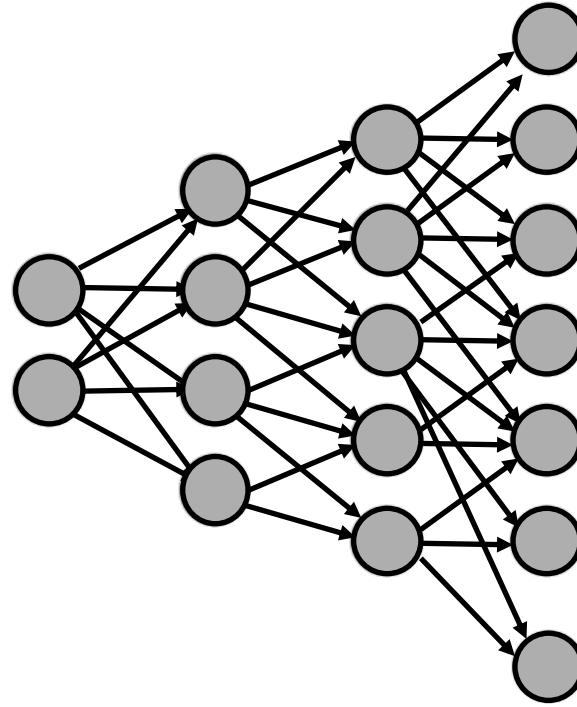
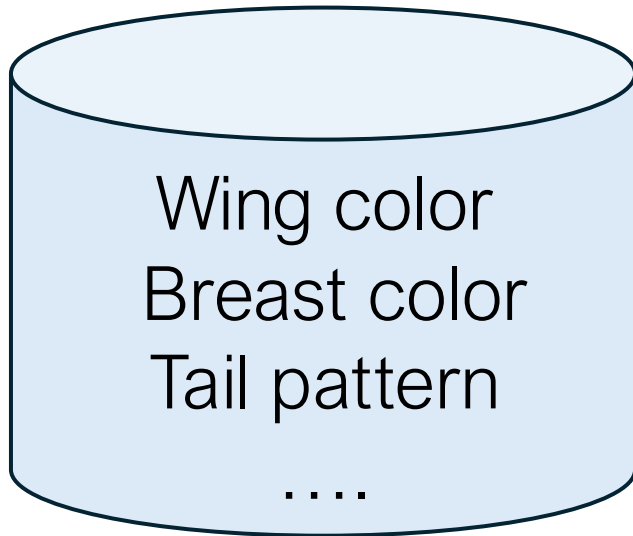
Class: **Pneumothorax**
Correlation : **Chest tube**
Accuracy: **60.2%**

How to detect? (Aim 1)



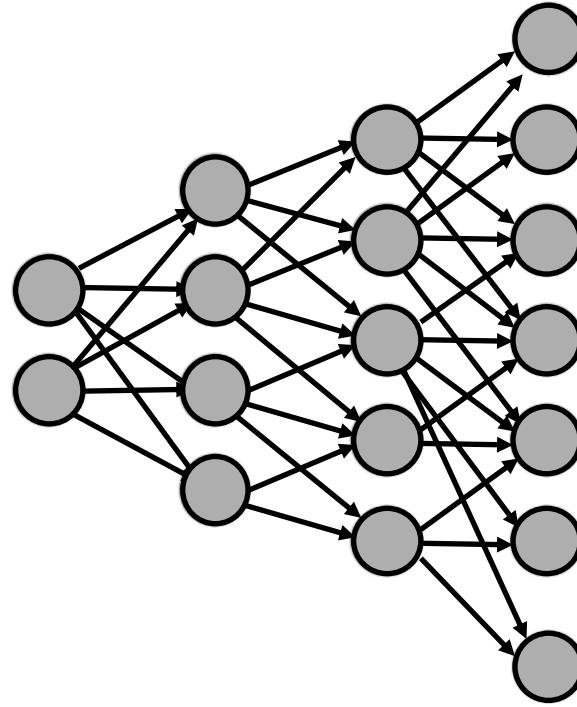
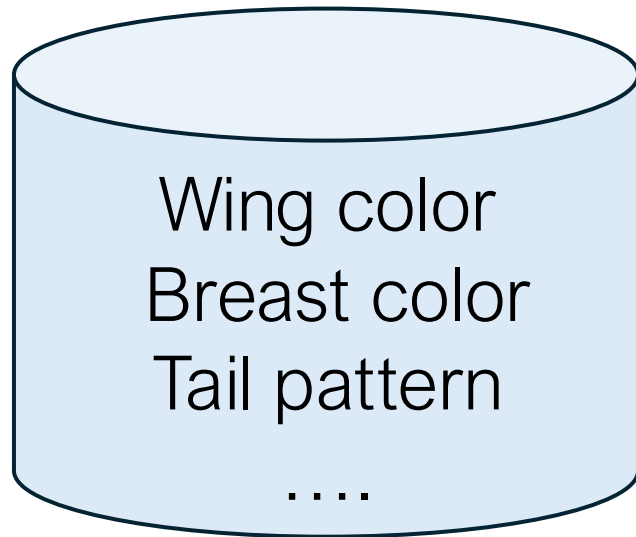
How to detect? (Aim 1)

Concept bank

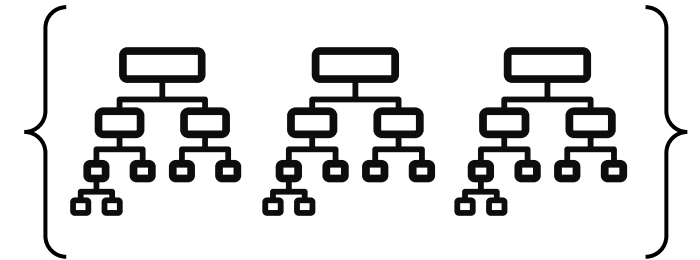


How to detect? (Aim 1)

Concept bank

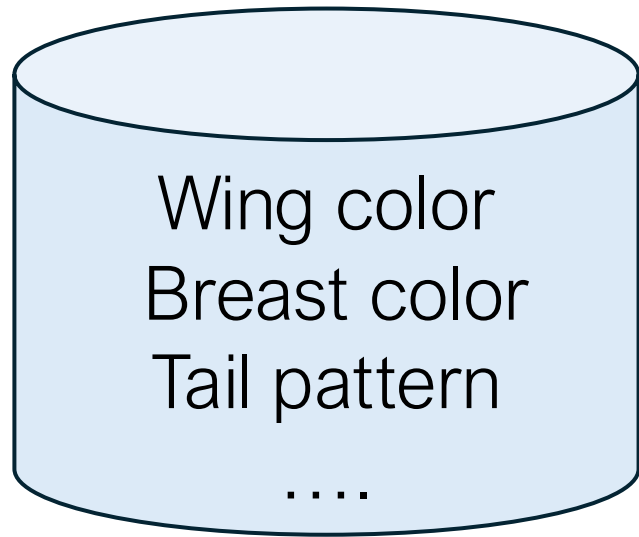


Olive sided Flycatcher \leftrightarrow breast_color_grey \wedge
tail_pattern_solid

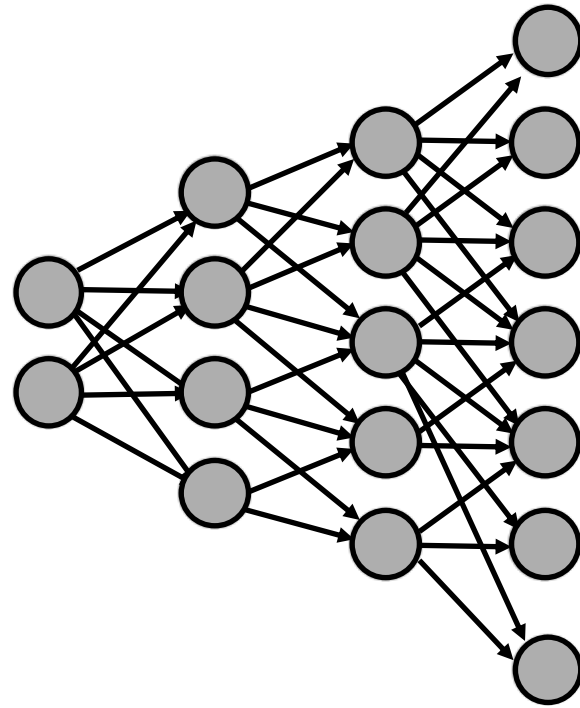


The goal: **Detect the Systematic Mistakes**

Concept bank

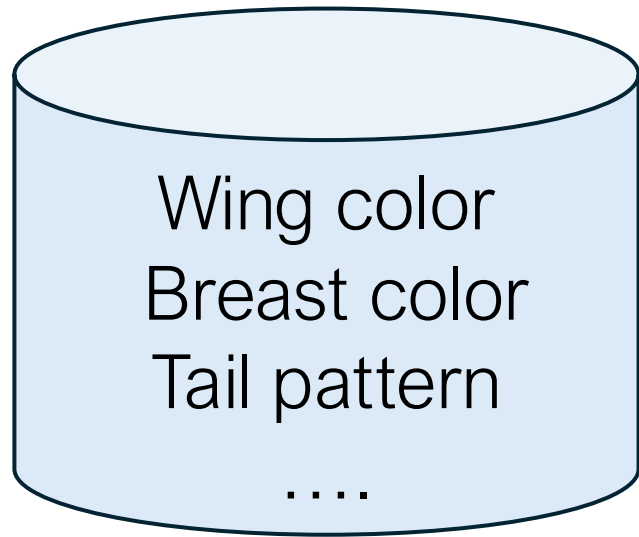


Expensive

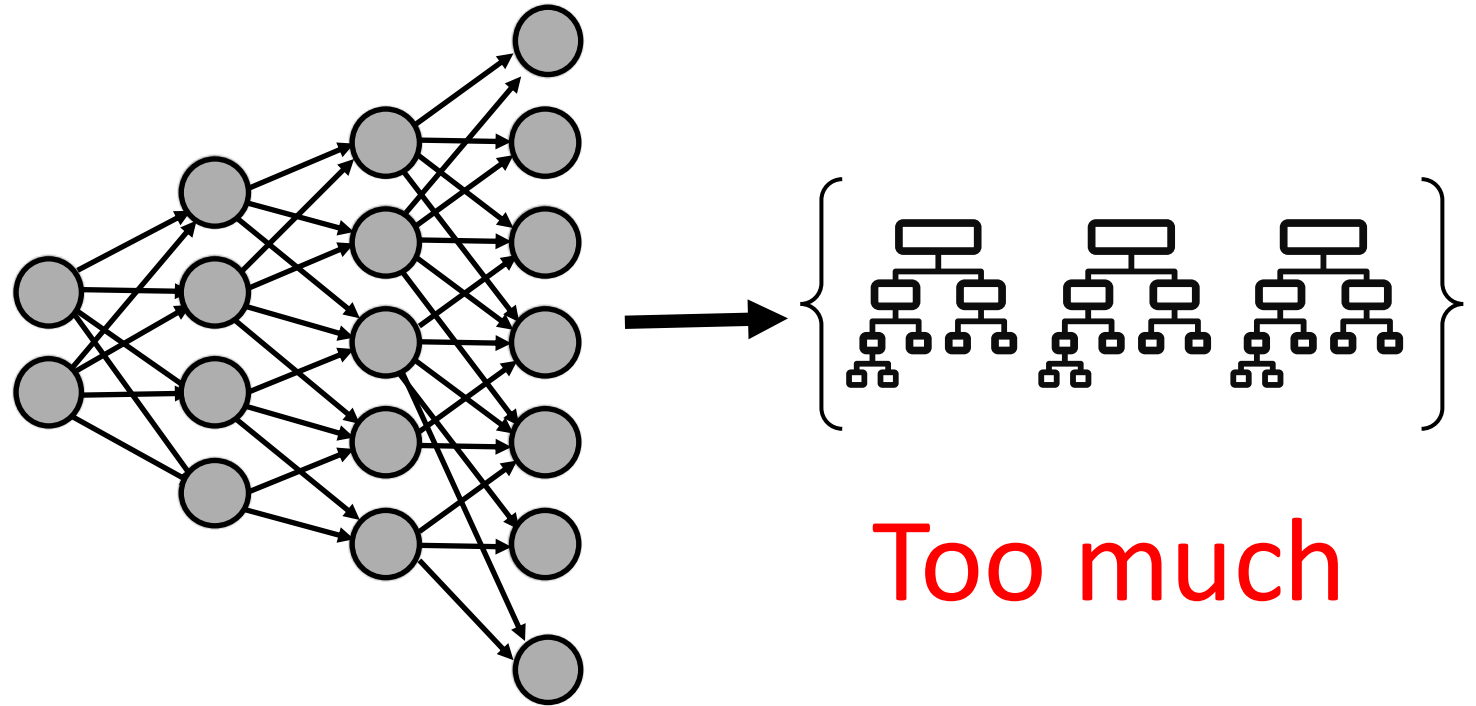


The goal: **Detect the Systematic Mistakes**

Concept bank



Expensive



Too much

Captions

1. A large seagull stands on a dock against a **backdrop of a harbor** with boats and a blue sky
2. A digitally altered image features a large bird, possibly an albatross, superimposed over a backdrop of industrial buildings **by a body of water**
3. A seagull stands on rocks **by the water** at sunset, with a lighthouse visible in the background

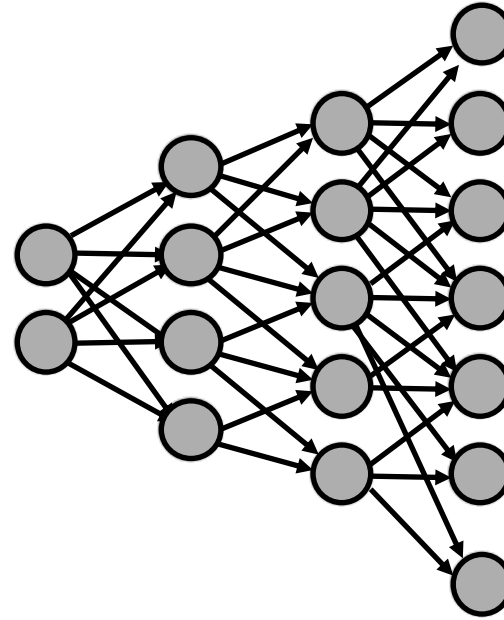
Report

1. perhaps mild increase in hydropneumothorax but with **chest tube**
2. other less likely possibility include expansion of known loculated hydropneumothorax (**chest tube** does not appear to be draining this region)
3. one of two right - **sided pleural tubes** has been removed in the interval
4. **3 chest tubes** remain in place and there is again an area of hydro pneumothorax

How to detect? (Aim 2 & 3)

Free-text

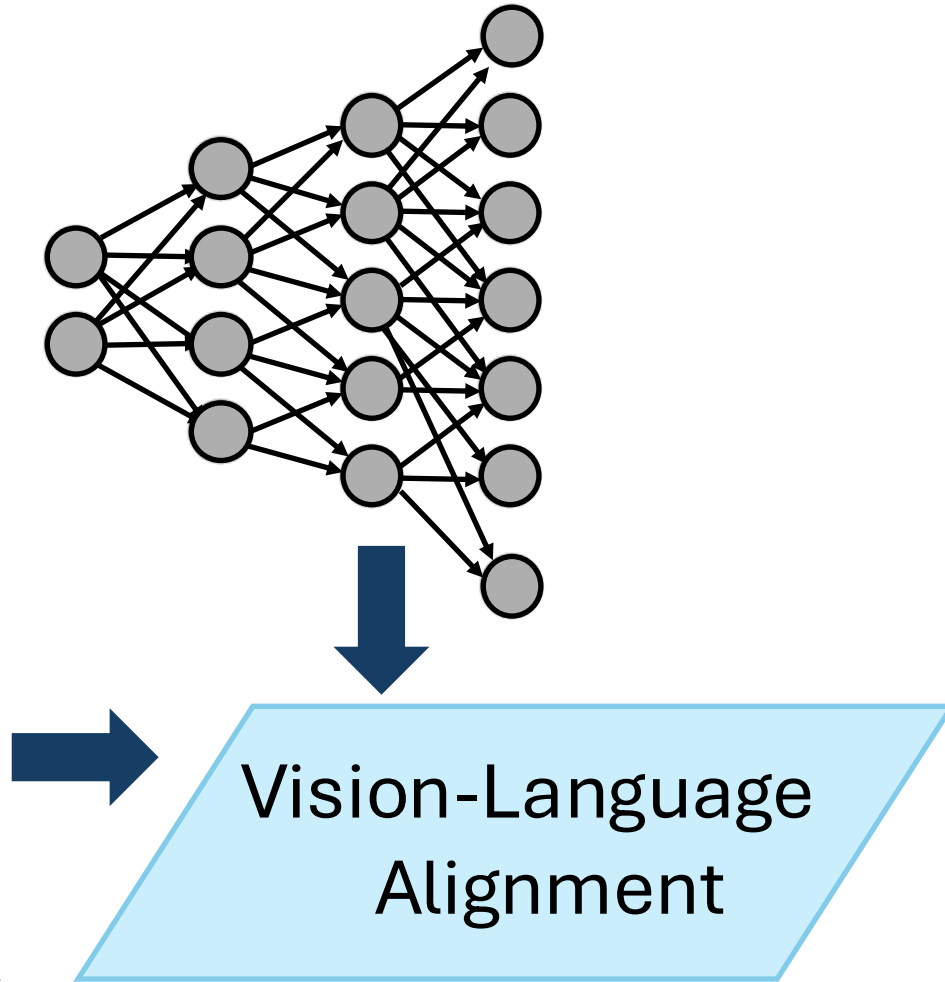
1. A large seagull stands on a dock against a **backdrop of a harbor** with boats and a blue sky
2. A digitally altered image features a large bird, possibly an albatross, superimposed over a backdrop of industrial buildings **by a body of water**
3. A seagull stands on rocks **by the water** at sunset, with a lighthouse visible in the background



How to detect? (Aim 2 & 3)

Free-text

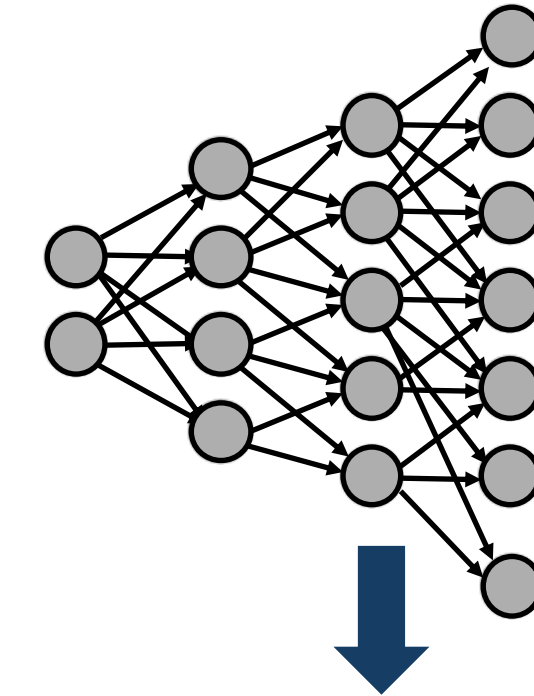
1. A large seagull stands on a dock against a **backdrop of a harbor** with boats and a blue sky
2. A digitally altered image features a large bird, possibly an albatross, superimposed over a backdrop of industrial buildings **by a body of water**
3. A seagull stands on rocks **by the water** at sunset, with a lighthouse visible in the background



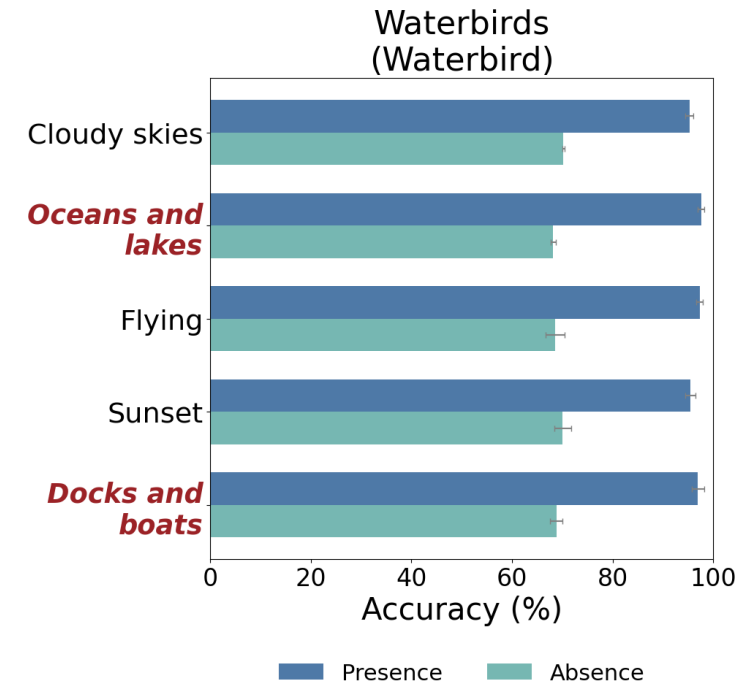
How to detect? (Aim 2 & 3)

Free-text

1. A large seagull stands on a dock against a **backdrop of a harbor** with boats and a blue sky
2. A digitally altered image features a large bird, possibly an albatross, superimposed over a backdrop of industrial buildings **by a body of water**
3. A seagull stands on rocks **by the water** at sunset, with a lighthouse visible in the background



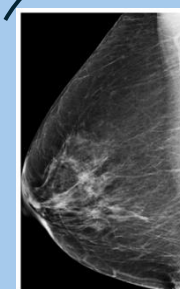
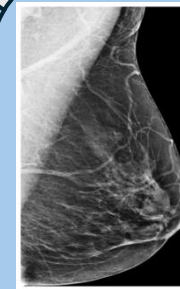
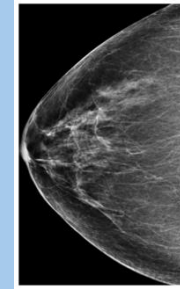
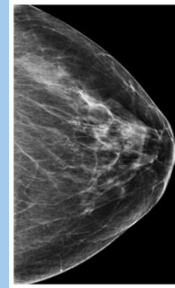
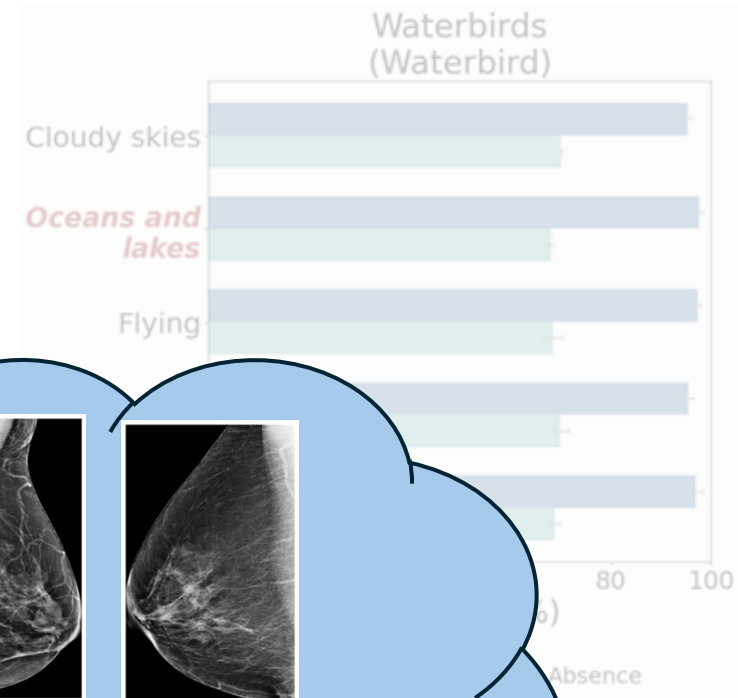
Vision-Language
Alignment



How to detect? (Aim 2)

Free-text

1. A large seagull stands on a dock against a backdrop of a harbor with boats and a blue sky
2. A digitally altered image features a large bird, possibly an albatross, superimposed over a backdrop of industrial buildings by a body of water
3. A seagull stands on rocks by water at sunset, with a lighthouse visible in the background



Vision-Language
Alignment

Full Length Article

Global burden of female breast cancer: new estimates in 2022, temporal trend and future projections up to 2050 based on the latest release from GLOBOCAN



Yunmeng Zhang^{1,†}, Yuting Ji^{1,†}, Siwen Liu¹, Jingjing Li¹, Jie Wu¹, Qianyun Jin¹, Xiaomin Liu¹, Hongyuan Duan¹, Zhuowei Feng¹, Ya Liu¹, Yacong Zhang², Zhangyan Lyu¹, Fangfang Song¹, Fengju Song¹, Lei Yang³, Hong Liu^{4,*}, Yubei Huang^{1,*}

Results: In 2022, an estimated 2.3 million new BC cases and 666,000 BC-related deaths occurred globally, accounting for 23.8 % and 15.4 % of all cancer cases and deaths in women, respectively. Regionally, Eastern Asia

Effect of mammographic screening from age 40 years on breast cancer mortality (UK Age trial): final results of a randomised, controlled trial



Stephen W Duffy*, Daniel Vulkan*, Howard Cuckle, Dharmishta Parmar, Shama Sheikh, Robert A Smith, Andrew Evans, Oleg Blyuss, Louise Johns, Ian O Ellis, Jonathan Myles, Peter D Sasieni*, Sue M Moss*



Summary

Background The appropriate age range for breast cancer screening remains a matter of debate. We aimed to estimate *Lancet Oncol* 2020; 21: 1165–72

Challenges of early screening

Annals of Internal Medicine

ORIGINAL RESEARCH

Estimation of Breast Cancer Overdiagnosis in a U.S. Breast Screening Cohort

Marc D. Ryser, PhD; Jane Lange, PhD; Lurdes Y.T. Inoue, PhD; Ellen S. O'Meara, PhD; Charlotte Gard, PhD; Diana L. Miglioretti, PhD; Jean-Luc Bulliard, PhD; Andrew F. Brouwer, PhD; E. Shelley Hwang, MD, MPH; and Ruth B. Etzioni, PhD

ORIGINAL ARTICLE

Open Access

Workload of diagnostic radiologists in the foreseeable future based on recent scientific advances: growth expectations and role of artificial intelligence

Thomas C. Kwee^{1*} and Robert M. Kwee²



A Systematic Review of Fatigue in Radiology: Is It a Problem?

Nadia Stec¹
Danielle Arje¹
Alan R. Moody¹
Elizabeth A. Krupinski²
Pascal N. Tyrrell^{1,3}

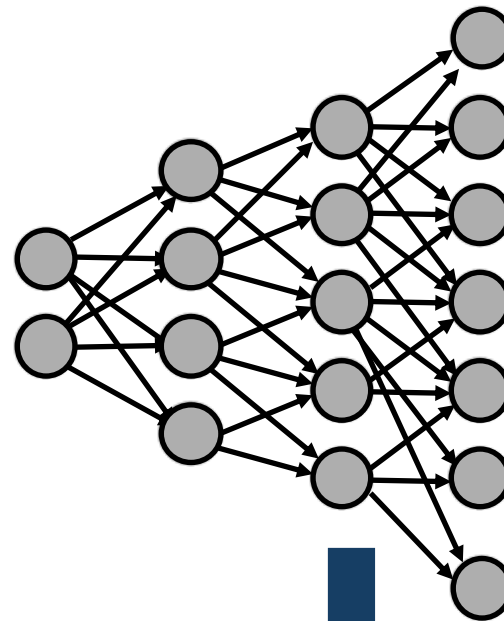
OBJECTIVE. The purpose of this study was to review current literature regarding radiologist fatigue.

MATERIALS AND METHODS. A literature search was performed using PubMed. Key words and Medical Subject Heading terms were used to generate refined queries with inclusion and exclusion criteria, focusing on fatigue and error. Results were selected according to these criteria: examined radiologist fatigue and radiologic error stemming from fatigue; experimental results measured as accuracy, error, or performance; and peer-reviewed publication. The risk of bias was addressed by including both quantitative and qualitative studies.

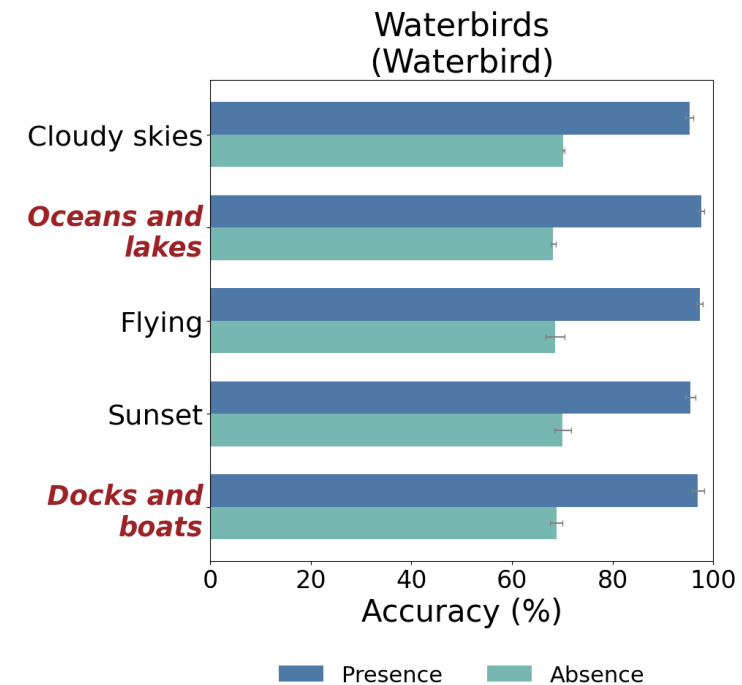
How to detect? (Aim 3)

Free-text

1. A large seagull stands on a dock against a **backdrop of a harbor** with boats and a blue sky
2. A digitally altered image features a large bird, possibly an albatross, superimposed over a backdrop of industrial buildings **by a body of water**
3. A seagull stands on rocks **by the water** at sunset, with a lighthouse visible in the background



Vision-Language
Alignment



Aim 1

The goal: Extract mixture of **Interpretable** models from the Blackbox **Post-hoc** using **FOL**

Why Post-hoc

✓ Does not alter the
Black box

✗ No intervention

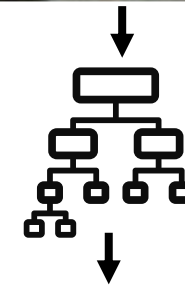
Why Post-hoc

✓ Does not alter the
Black box

✗ No intervention

Why Interpretable

✓ Supports interventions



Prediction: Brewer Blackbird ✗

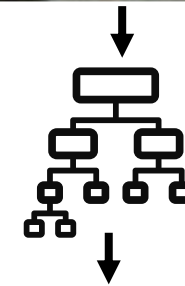
Why Post-hoc

✓ Does not alter the
Black box

✗ No intervention

Why Interpretable

✓ Supports interventions



Prediction: Brewer Blackbird ✗

Concepts	Concept values
bill_length_shorter_than_head	0.89
bill_shape_allpurpose	0.42
wing_shape_roundedwings	0.40
⋮	⋮

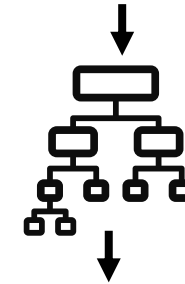
Why Post-hoc

✓ Does not alter the
Black box

✗ No intervention

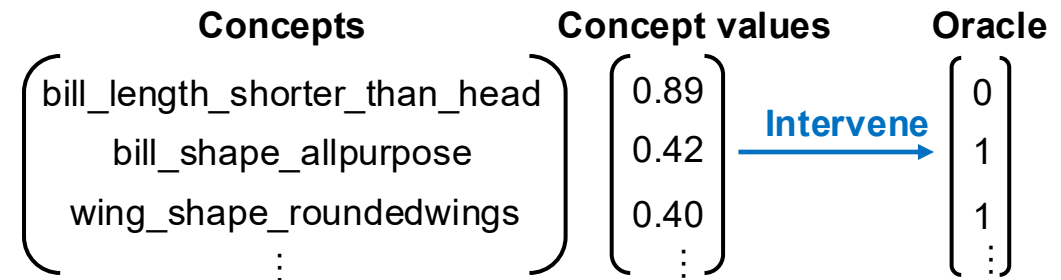
Why Interpretable

✓ Supports interventions



Prediction: Brewer Blackbird ✗

Fish Crow ✓



Why Post-hoc

✓ Does not alter the
Black box

✗ No intervention

Why Interpretable

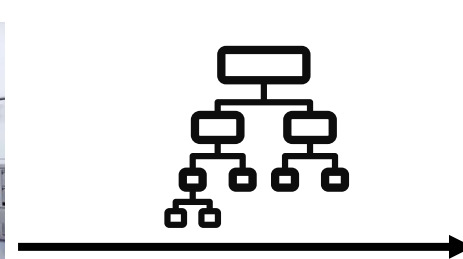
✓ Supports interventions

Why FOL?

Cardiomegaly \leftrightarrow heart_size \wedge enlarge



Site A



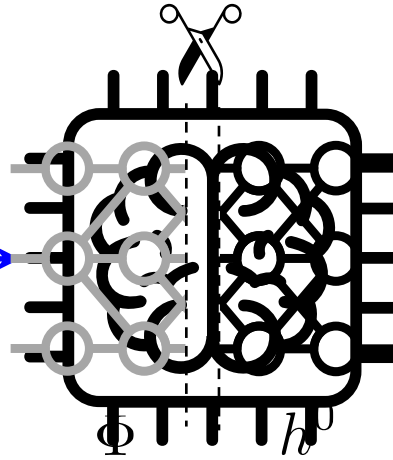
Site B

Problem Set Up

 \mathcal{X}  \mathcal{C}

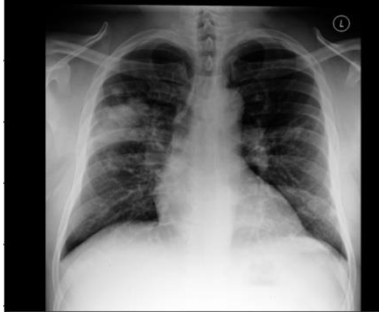
Wing color grey
Breast color white
Tail pattern

....

 \mathcal{Y}

Olive
Sided
Flycatcher

Problem Set Up

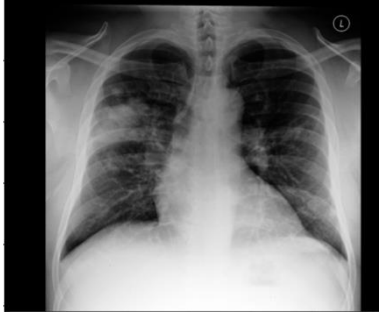


+

Report:

Right upper lobe consolidation with adjacent. While this may be infectious in nature, a CT scan is recommended for further clarification.

Problem Set Up



+

Report:

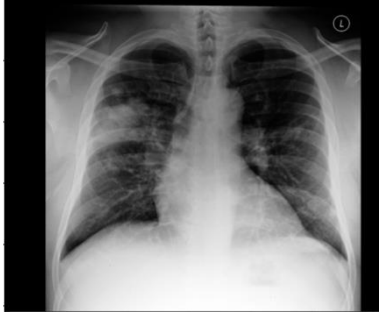
Right upper lobe **consolidation** with adjacent.
While this **may** be **infectious** in nature, a CT
scan is recommended for further clarification.

↓
parse the reports to get the concepts

\mathcal{C}

right upper lobe
left lower lobe
heart size
....

Problem Set Up

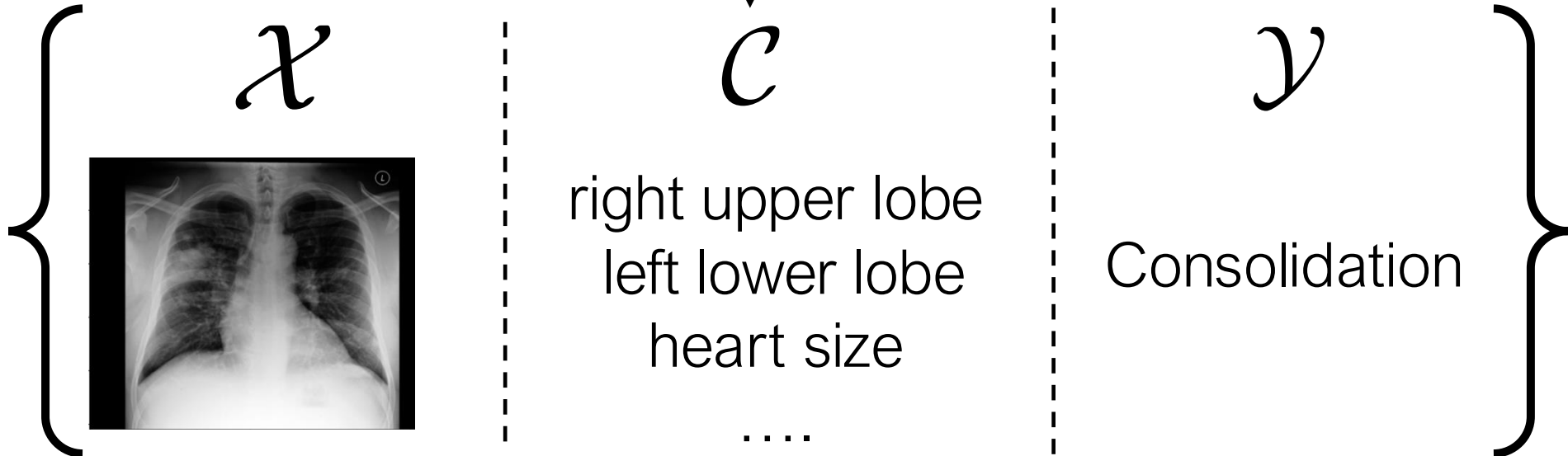


+

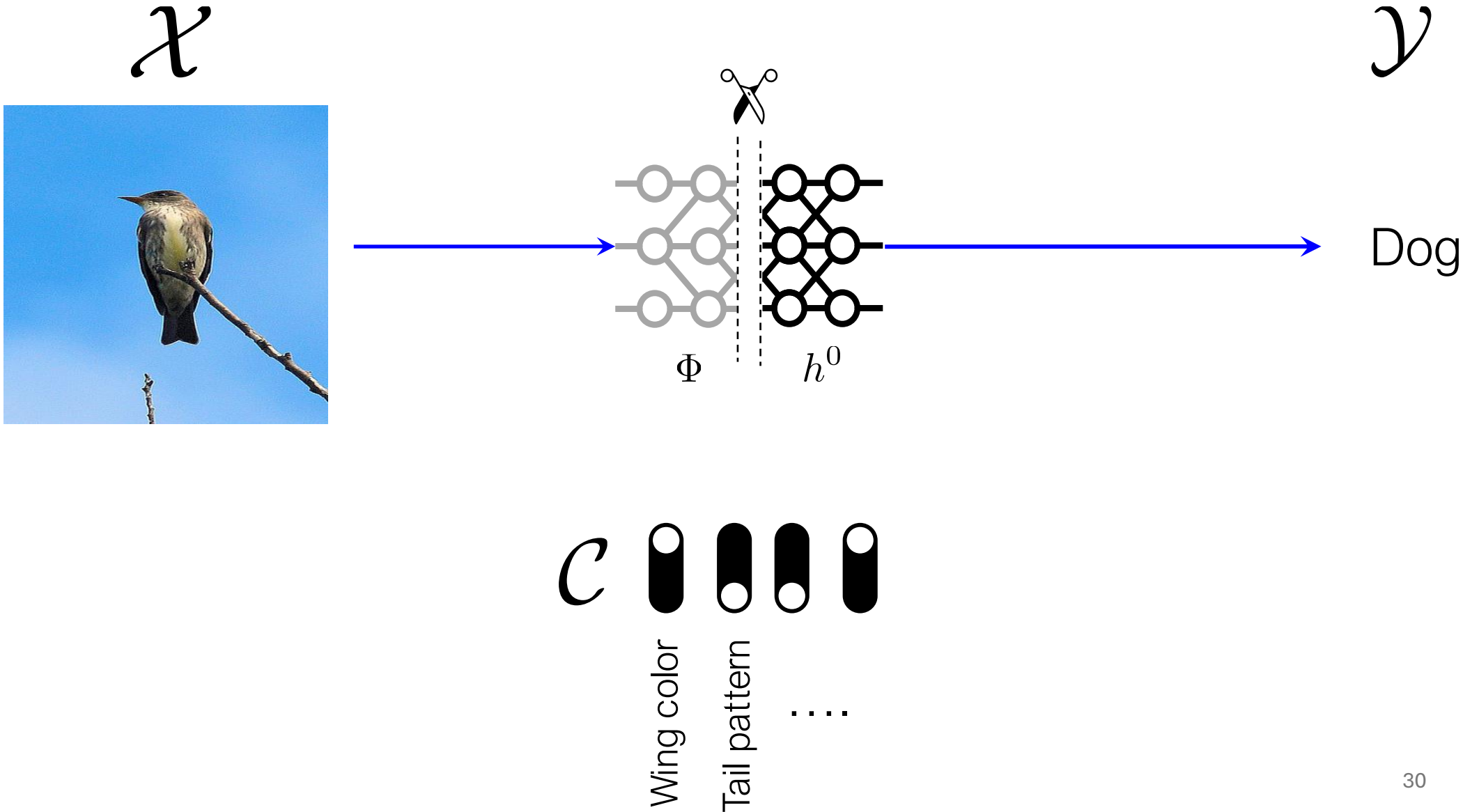
Report:

Right upper lobe consolidation with adjacent.
While this may be infectious in nature, a CT
scan is recommended for further clarification.

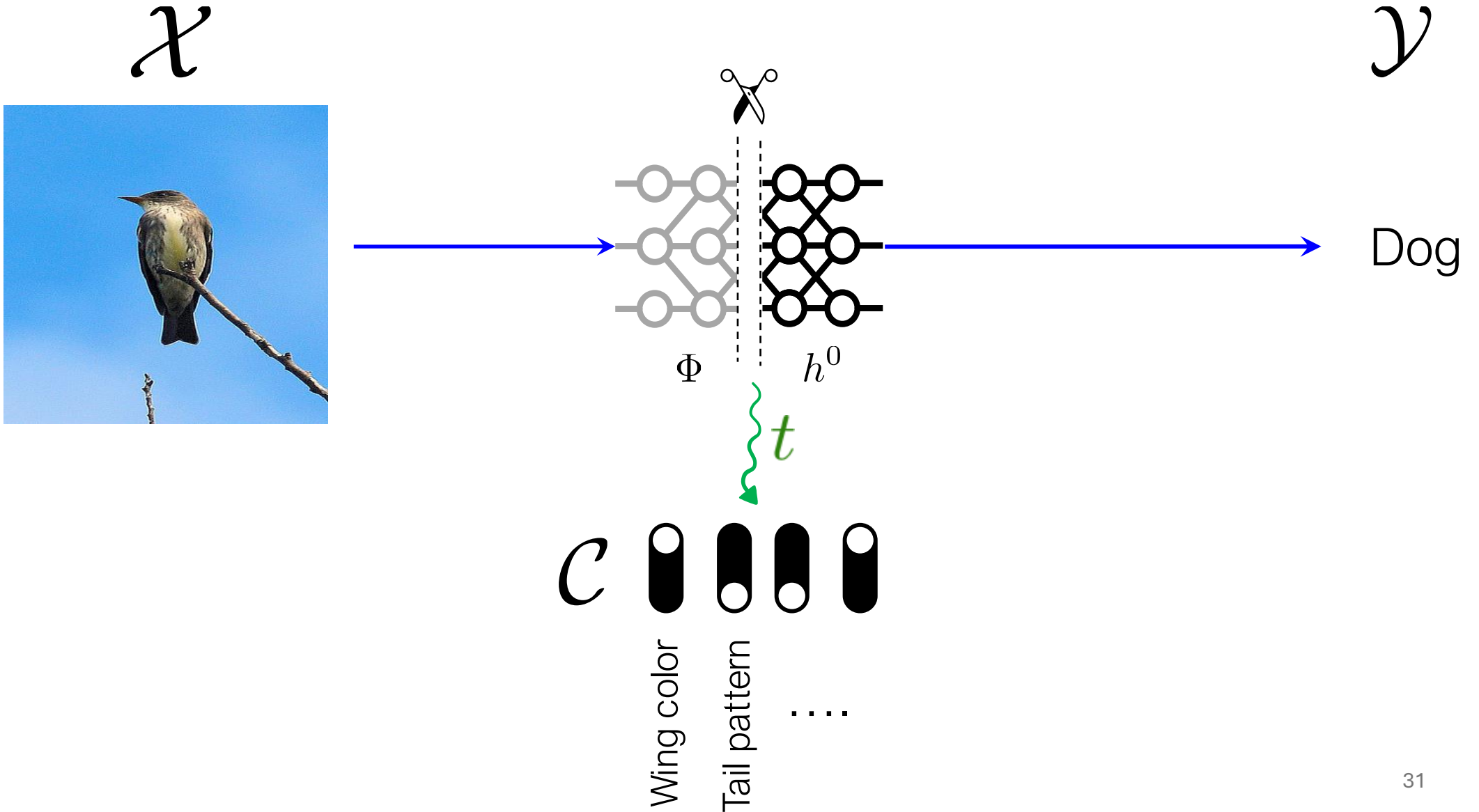
parse the reports to get the concepts



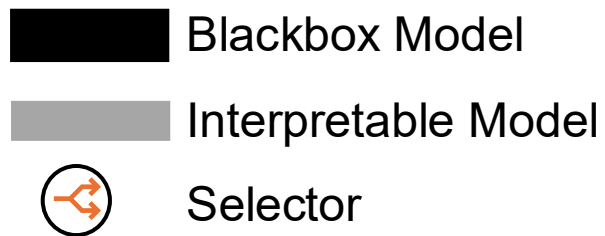
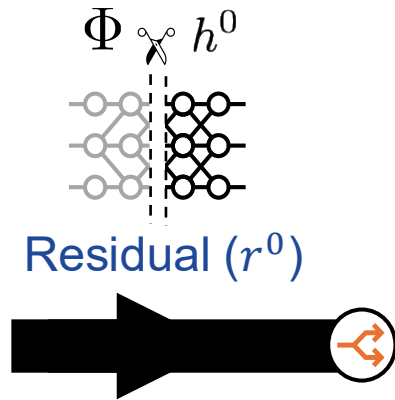
Discovering Hidden Concepts



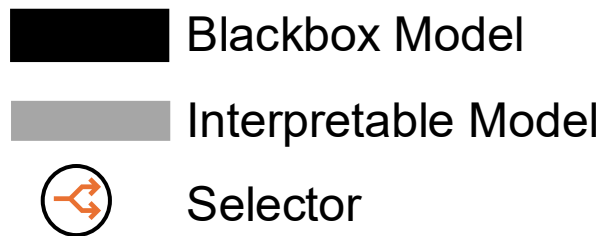
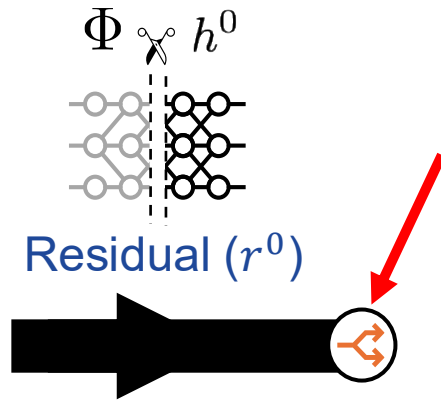
Discovering Hidden Concepts



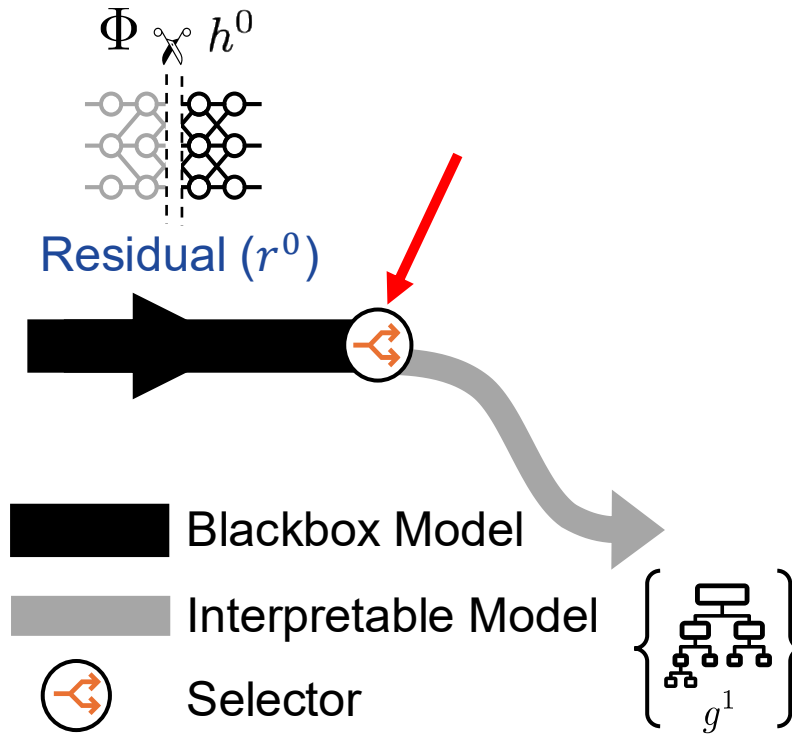
Carving out Interpretable Models



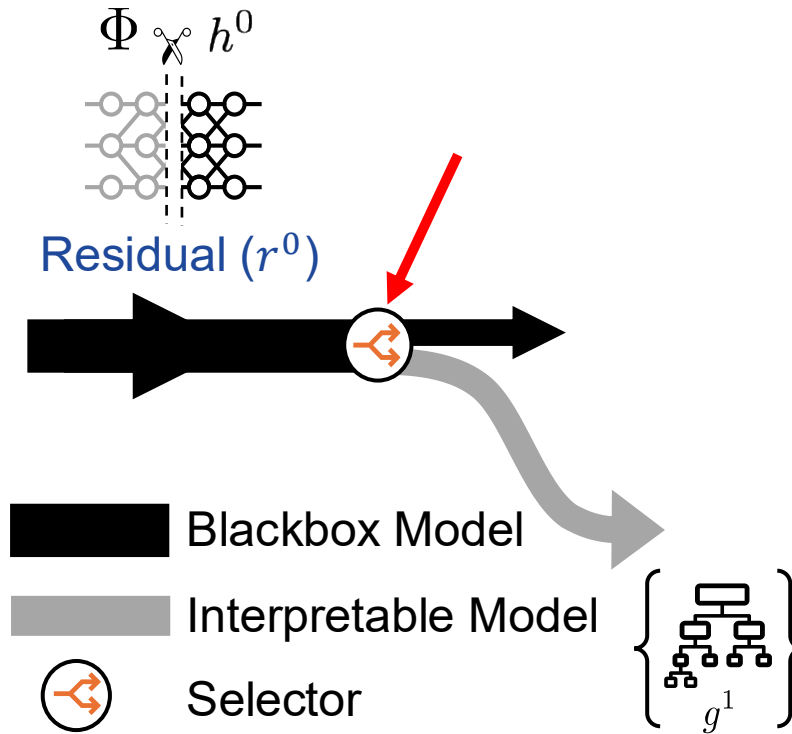
Carving out Interpretable Models



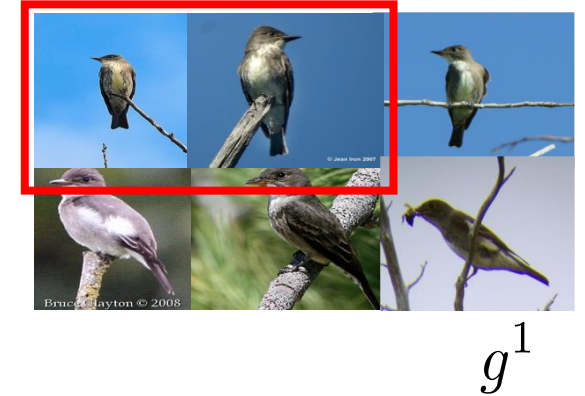
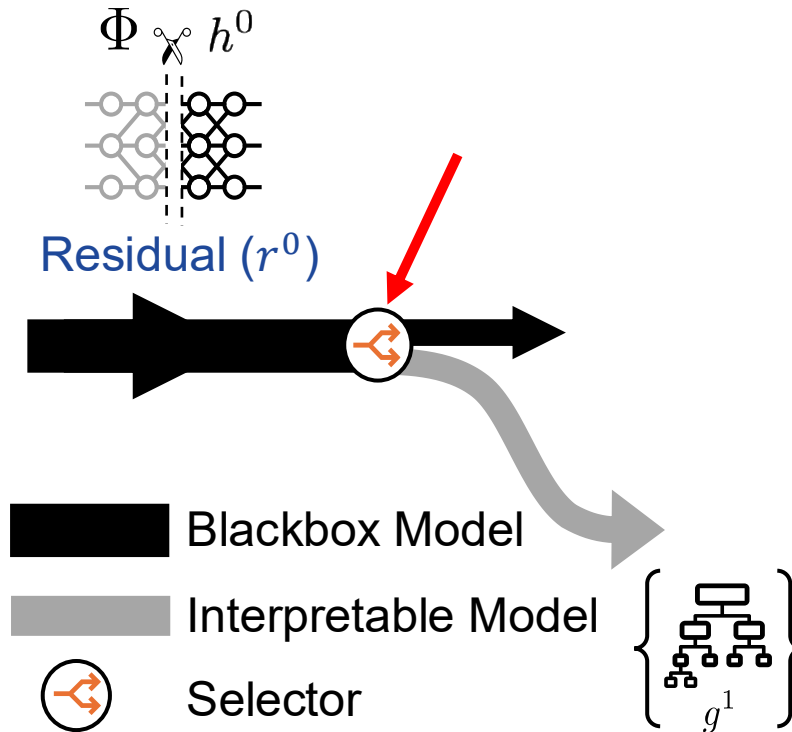
Carving out Interpretable Models



Carving out Interpretable Models

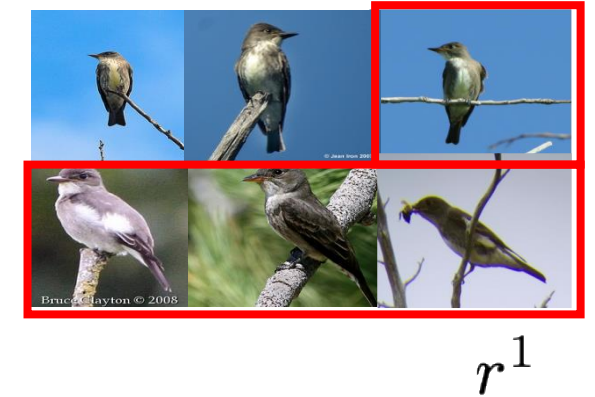
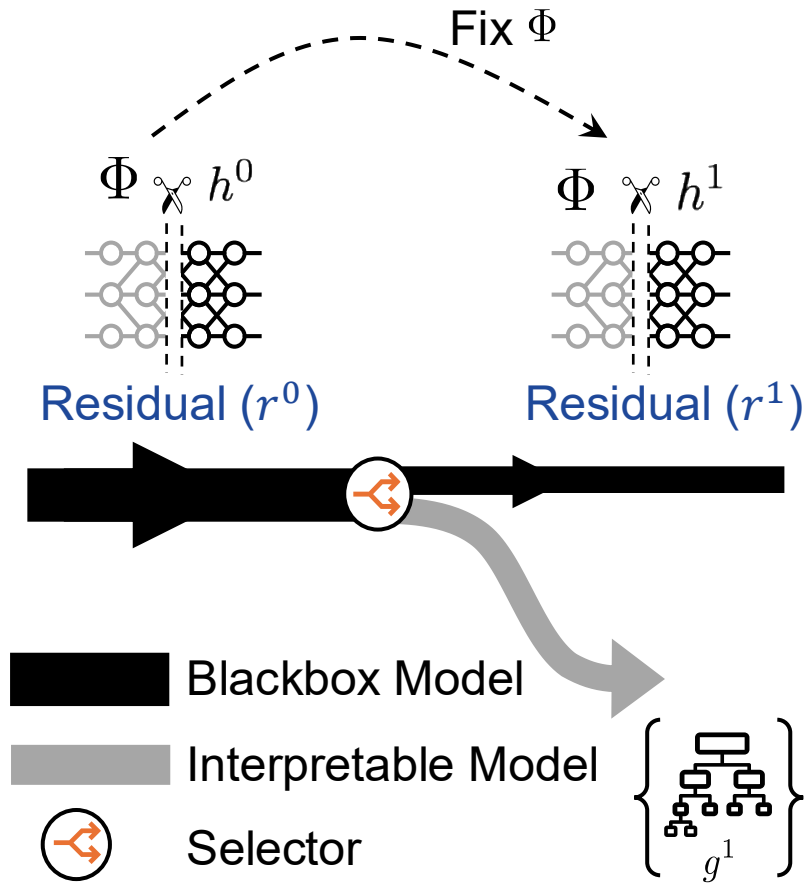


Carving out Interpretable Models

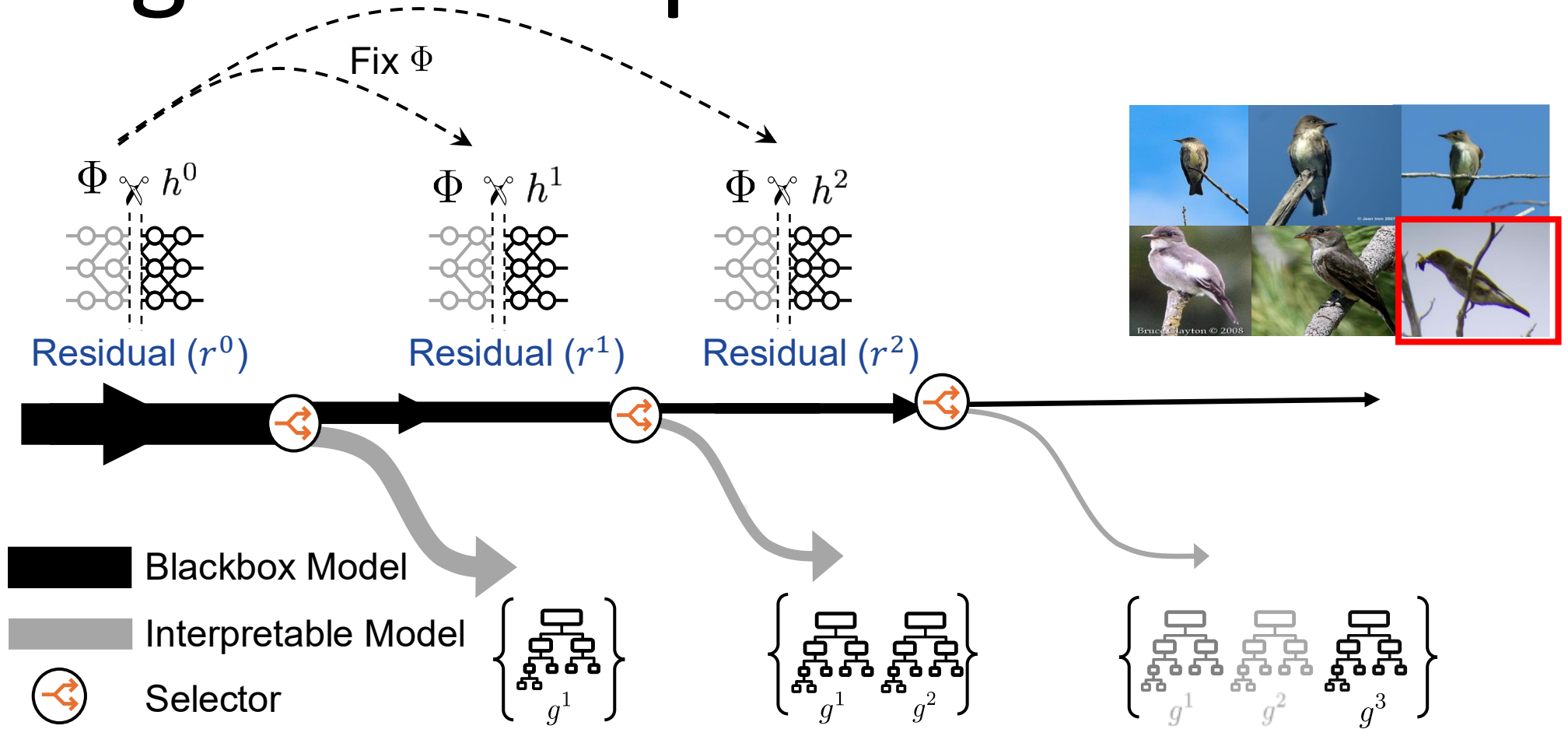


Olive sided Flycatcher \leftrightarrow `breast_color_grey` \wedge
`tail_pattern_solid`

Carving out Interpretable Models

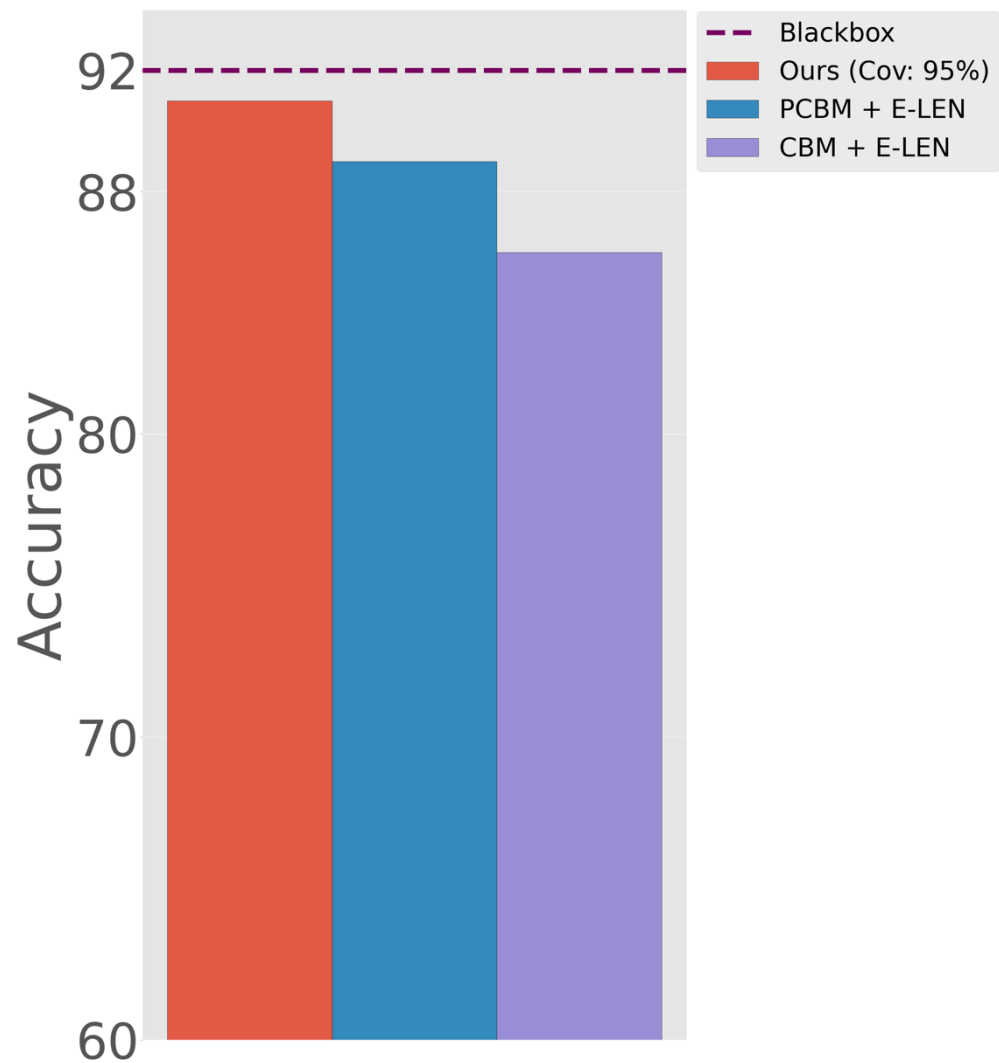


Carving out Interpretable Models



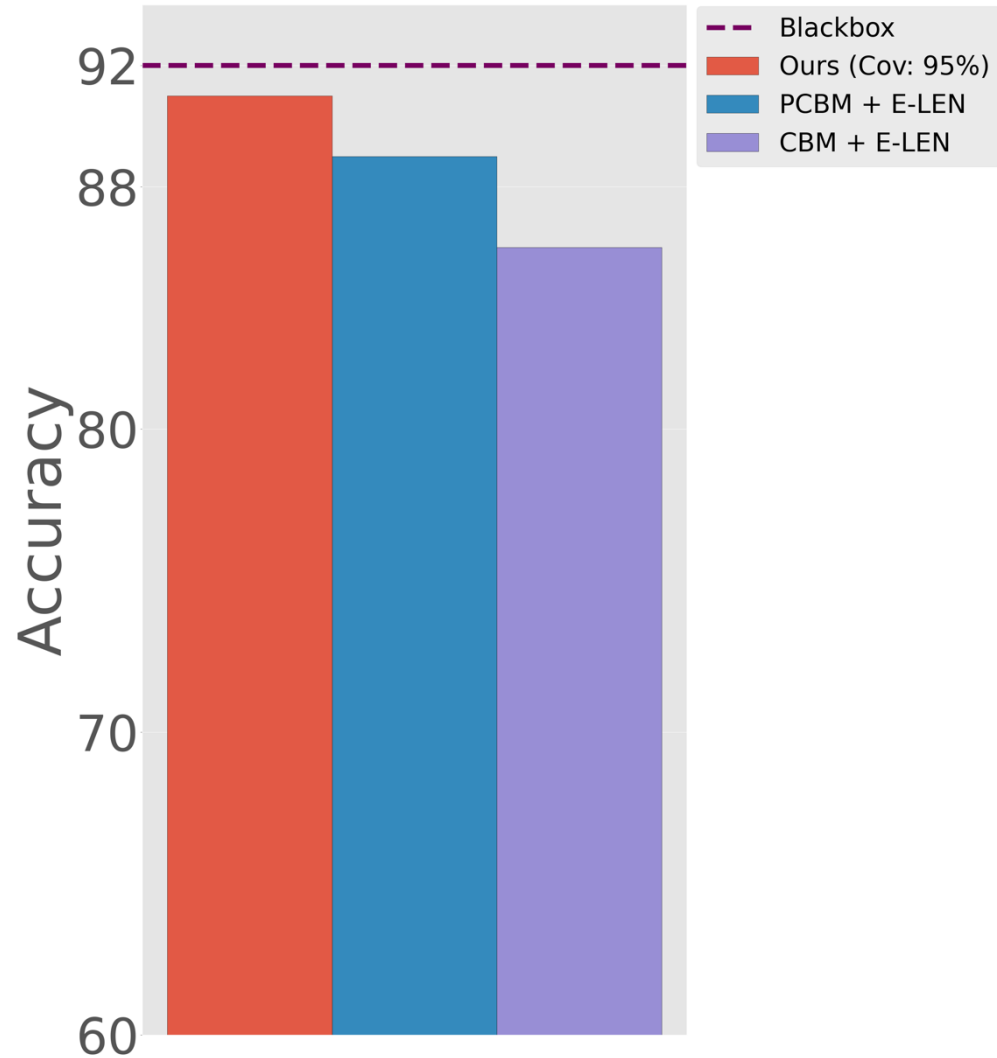
Comparing Performance

CUB-200 with ViT

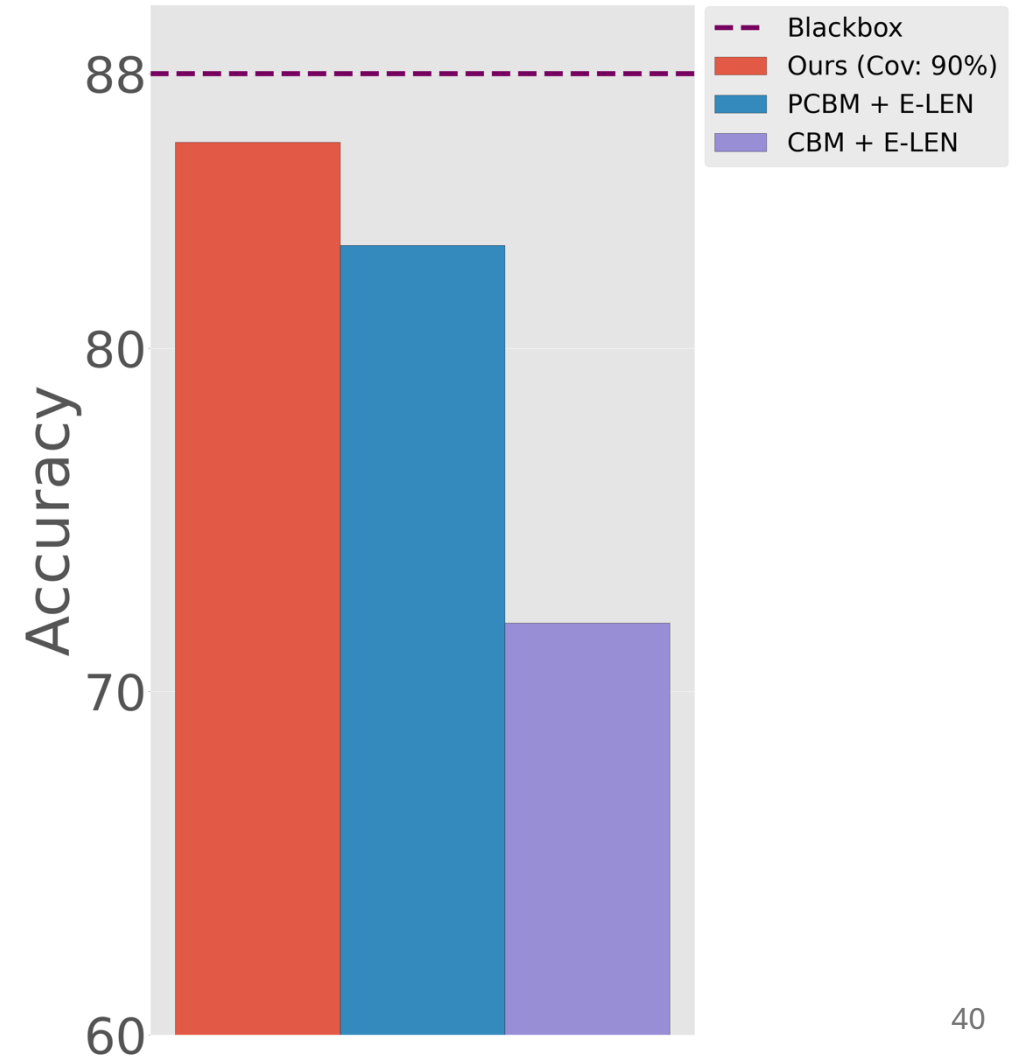


Comparing Performance

CUB-200 with ViT

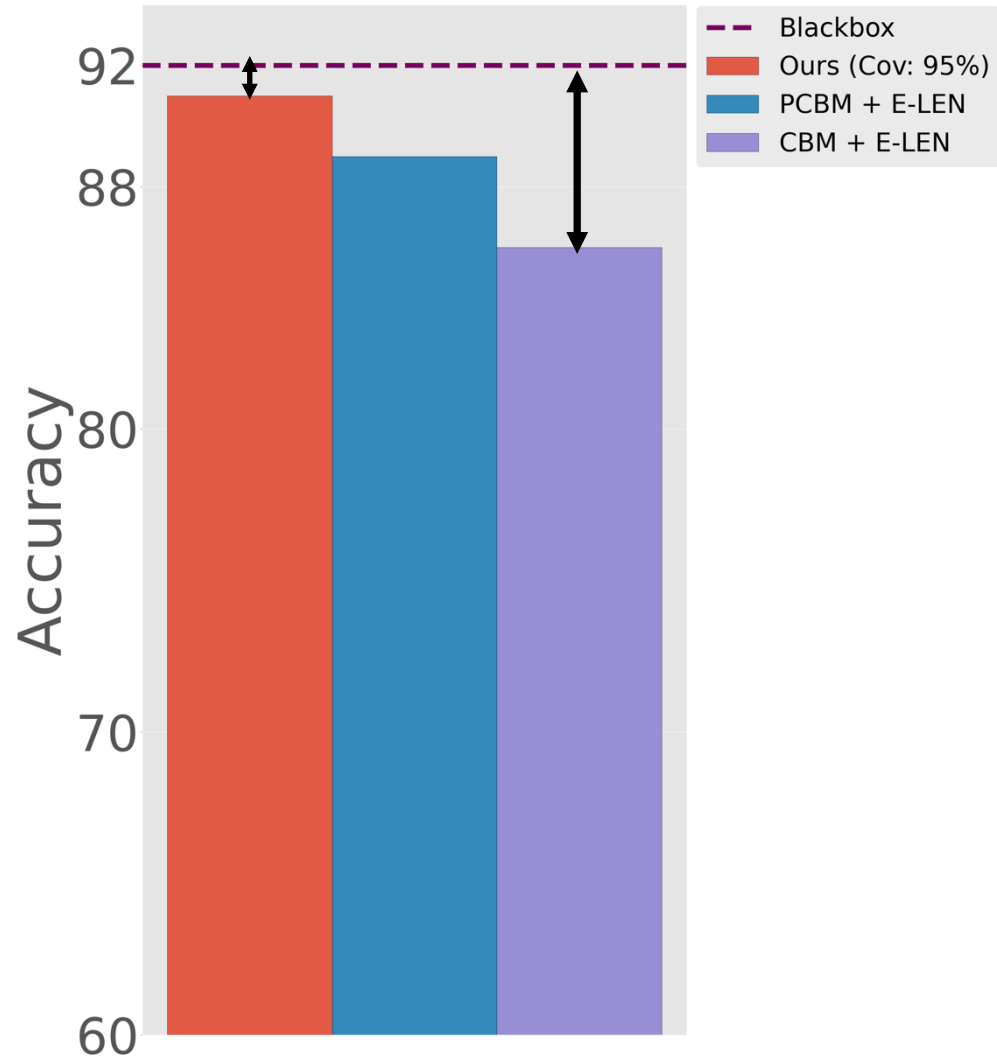


CUB-200 with ResNet101

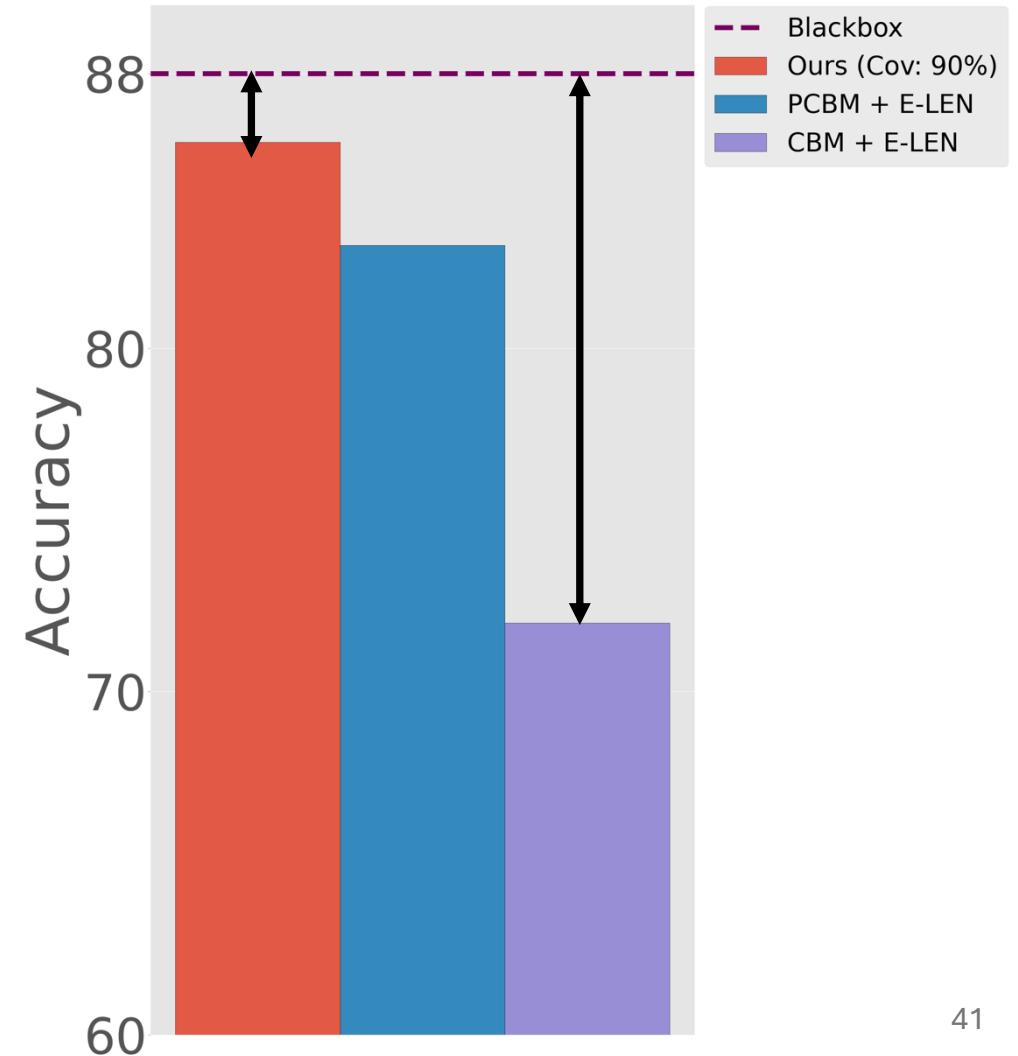


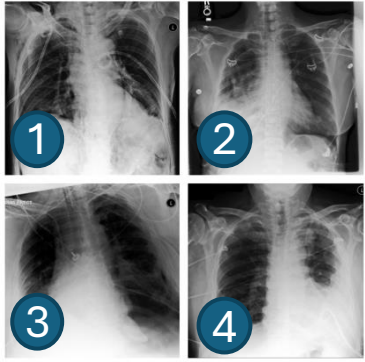
Comparing Performance

CUB-200 with ViT

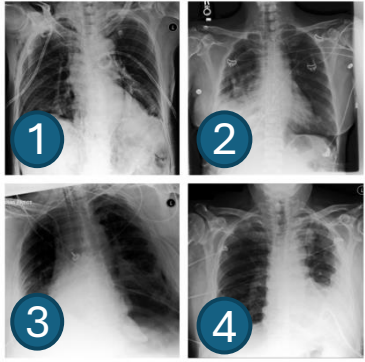


CUB-200 with ResNet101

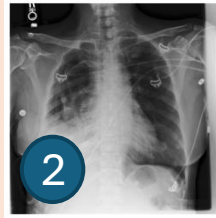
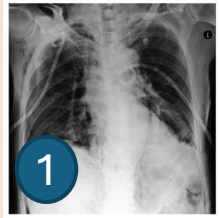




Examples on Chest X-ray



Examples on Chest X-ray



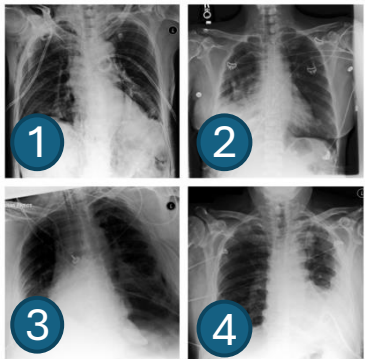
Expert 1

Pneumothorax

\leftrightarrow right_apical_lung

\wedge right_lung_unspec

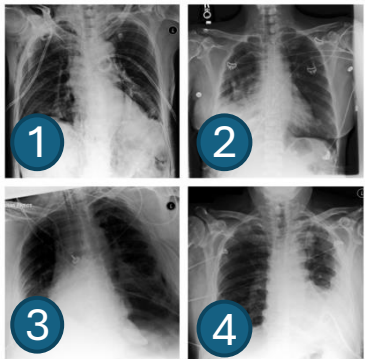
* Right lung unspec refers a malignant neoplasm or cancer in an unspecified part of the right bronchus or lung.



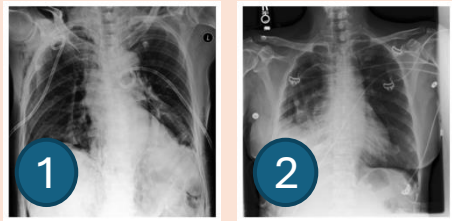
Examples on Chest X-ray

Expert 1
Pneumothorax
 \leftrightarrow right_apical_lung
 \wedge right_lung_unspec


Expert 2
Pneumothorax
 \leftrightarrow emphysema
 \wedge left_lung_unspec



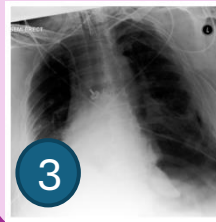
Examples on Chest X-ray



Expert 1
Pneumothorax
 \leftrightarrow right_apical_lung
 \wedge right_lung_unspec



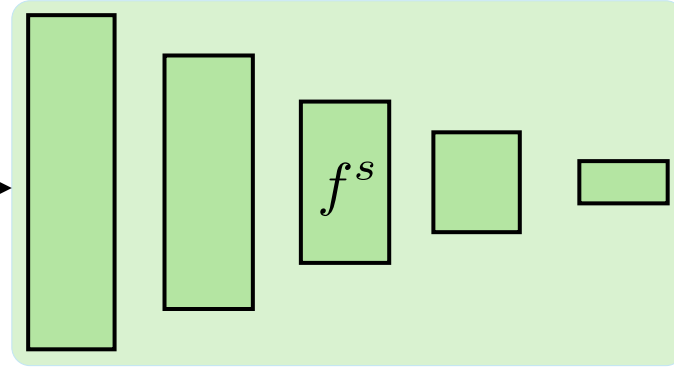
Expert 2
Pneumothorax
 \leftrightarrow emphysema
 \wedge left_lung_unspec



Expert 3
Pneumothorax
 \leftrightarrow left_apical_lung

Application: Data-Efficient Fine-tuning

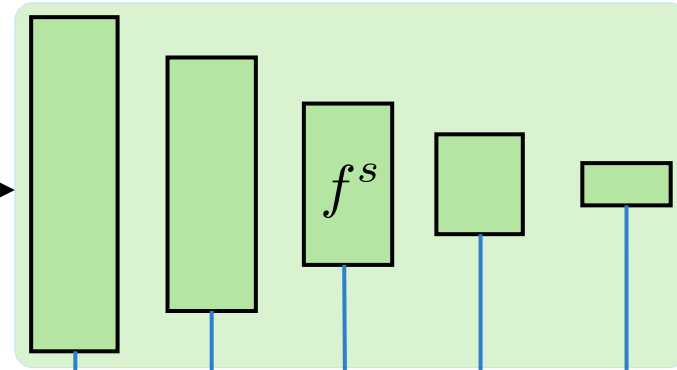
MIMIC-CXR



Pneumothorax

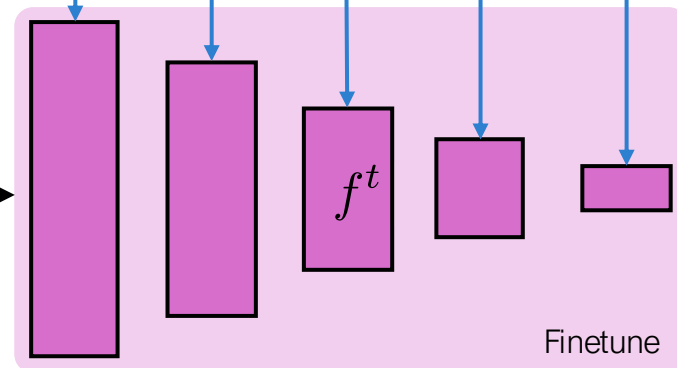
Application: Data-Efficient Fine-tuning

MIMIC-CXR



Pneumothorax

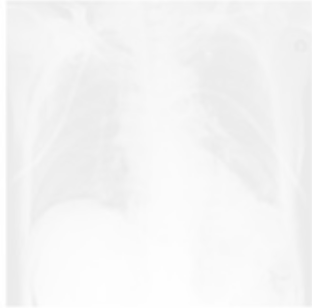
Stanford-CXR



Pneumothorax

Application: Data-Efficient Fine-tuning

MIMIC-CXR



Data and Computationally inefficient

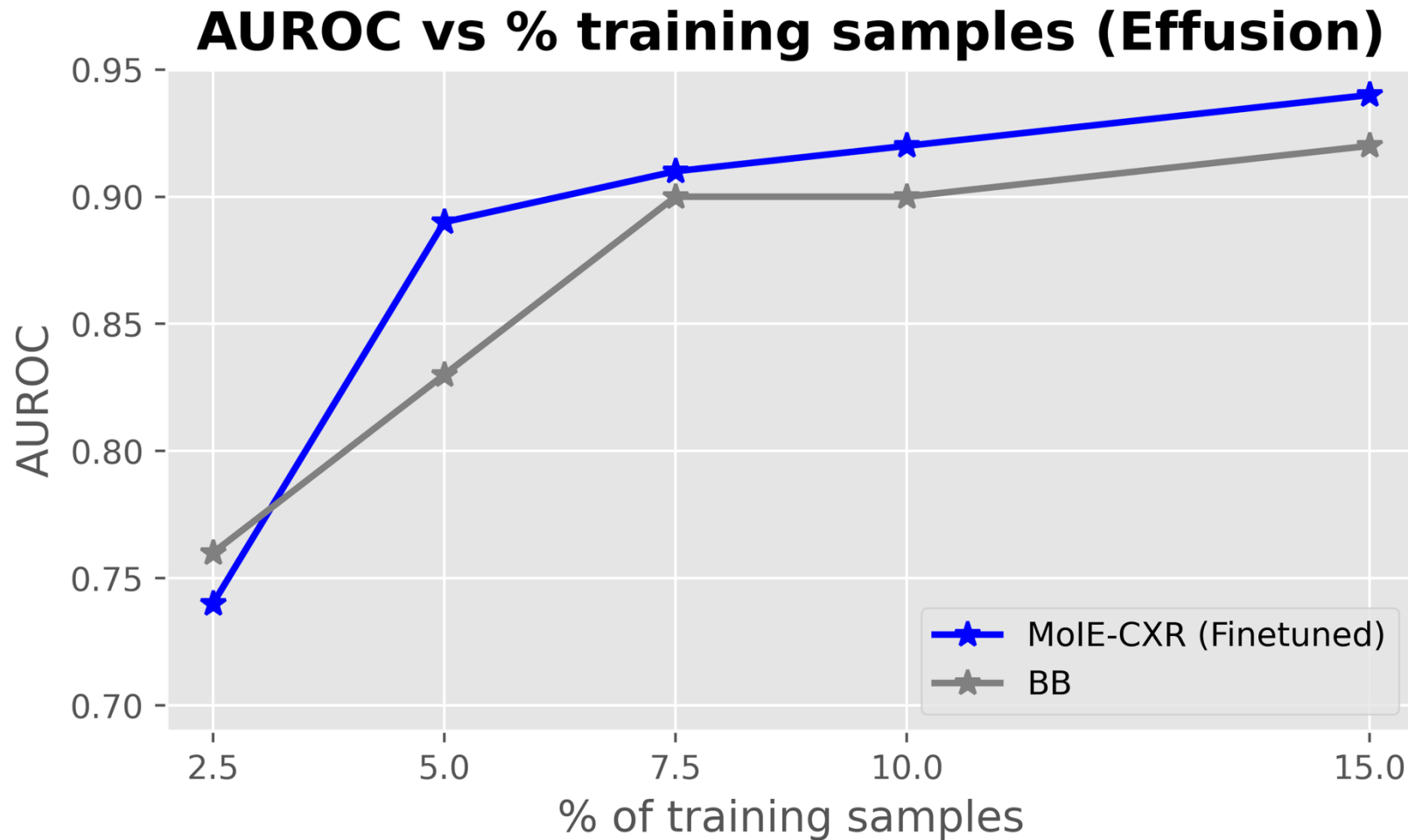
Stanford-CXR



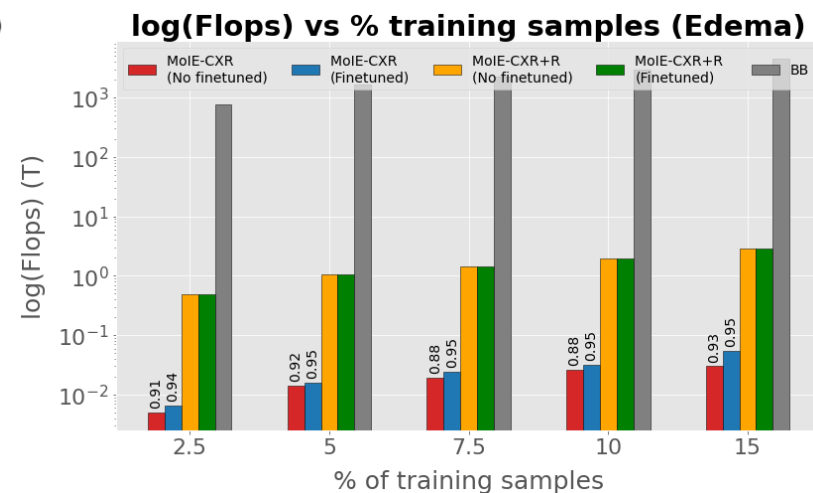
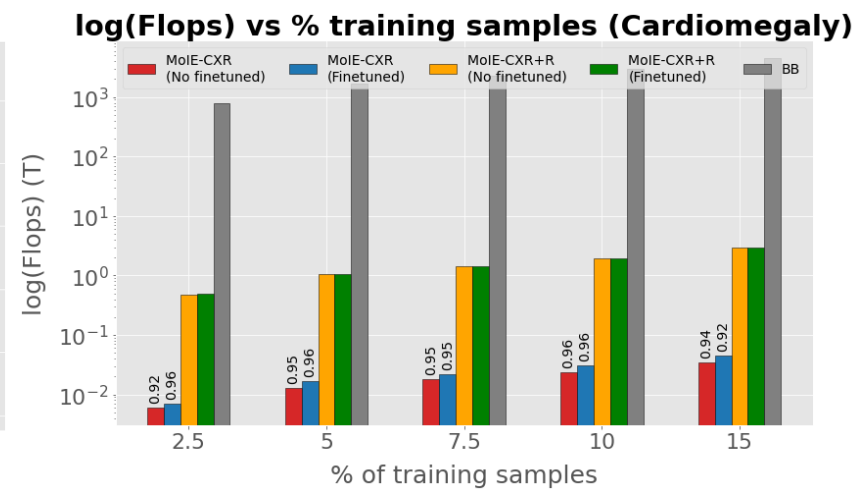
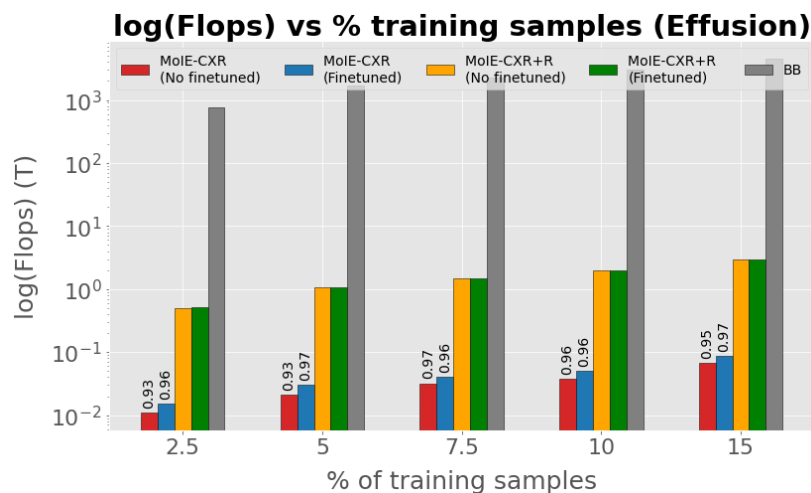
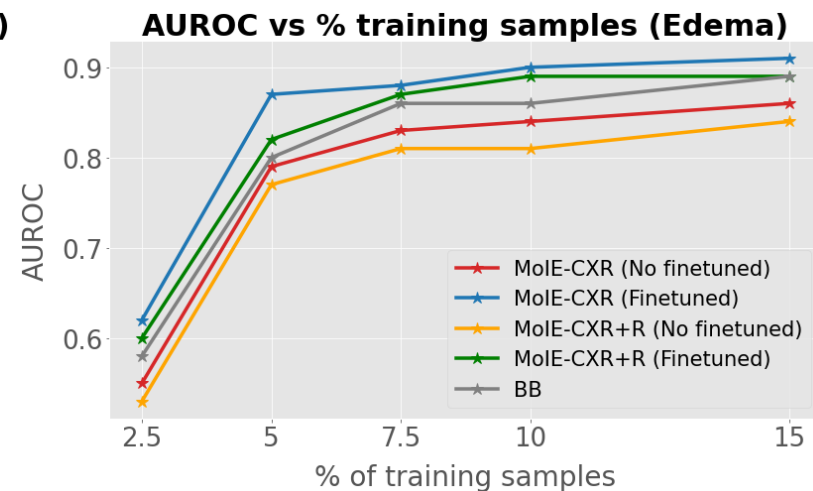
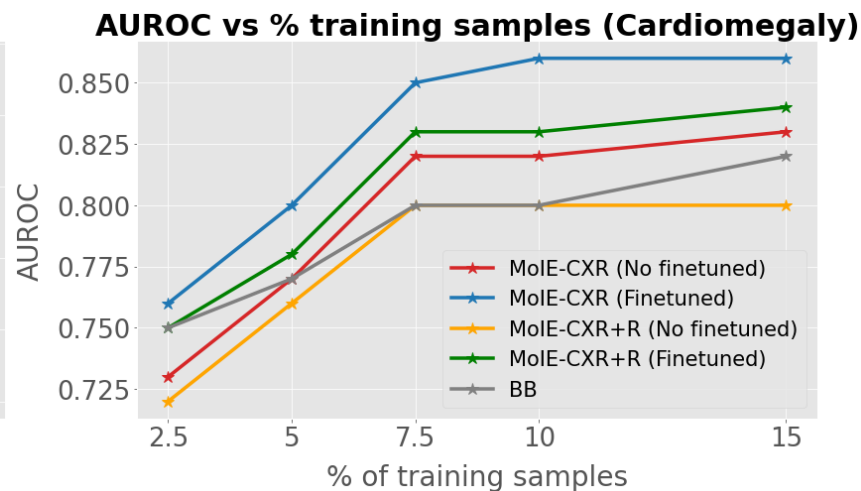
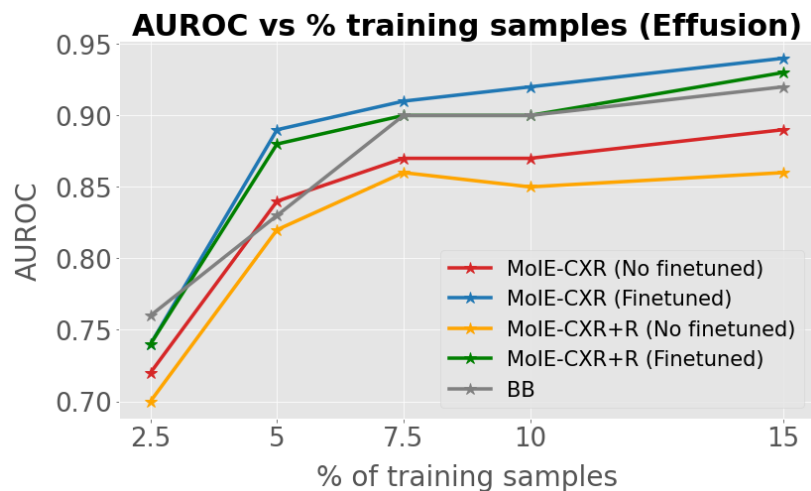
The clinical rules are “invariant”

Finetune

Transferring to Stanford-CXR



Transferring to Stanford-CXR



Conclusion from Aim 1

1. Domain invariant rules learned: A mixture of interpretable models are *carved out* of a Blackbox model offering best of both worlds. [ICML 2023]

They effectively learn domain invariant rules.

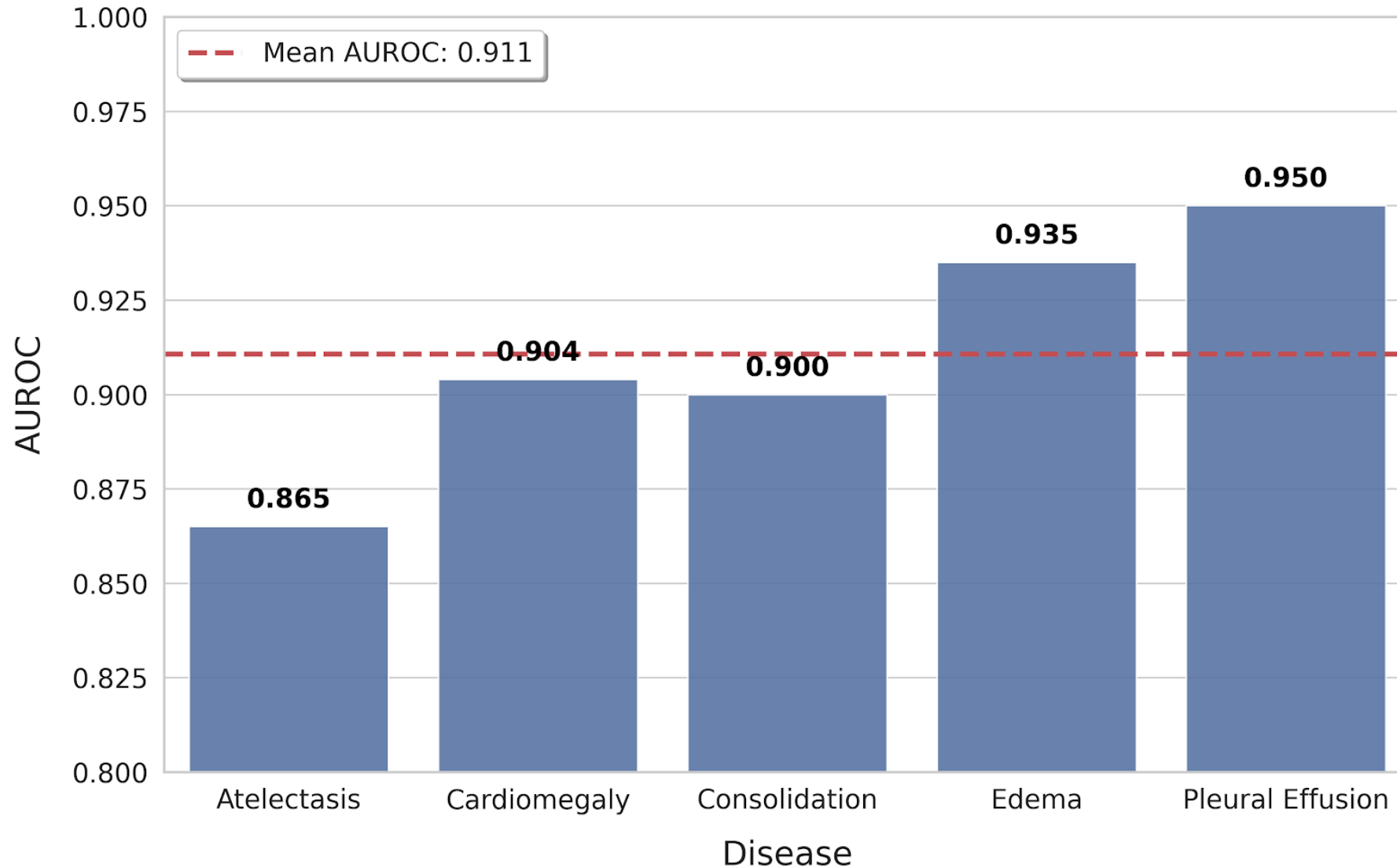
2. Efficient transfer learning: Transfer Learning is more efficient *using limited training data* with the new interpretable model. [MICCAI 2023. Top 14%]

Aim 2

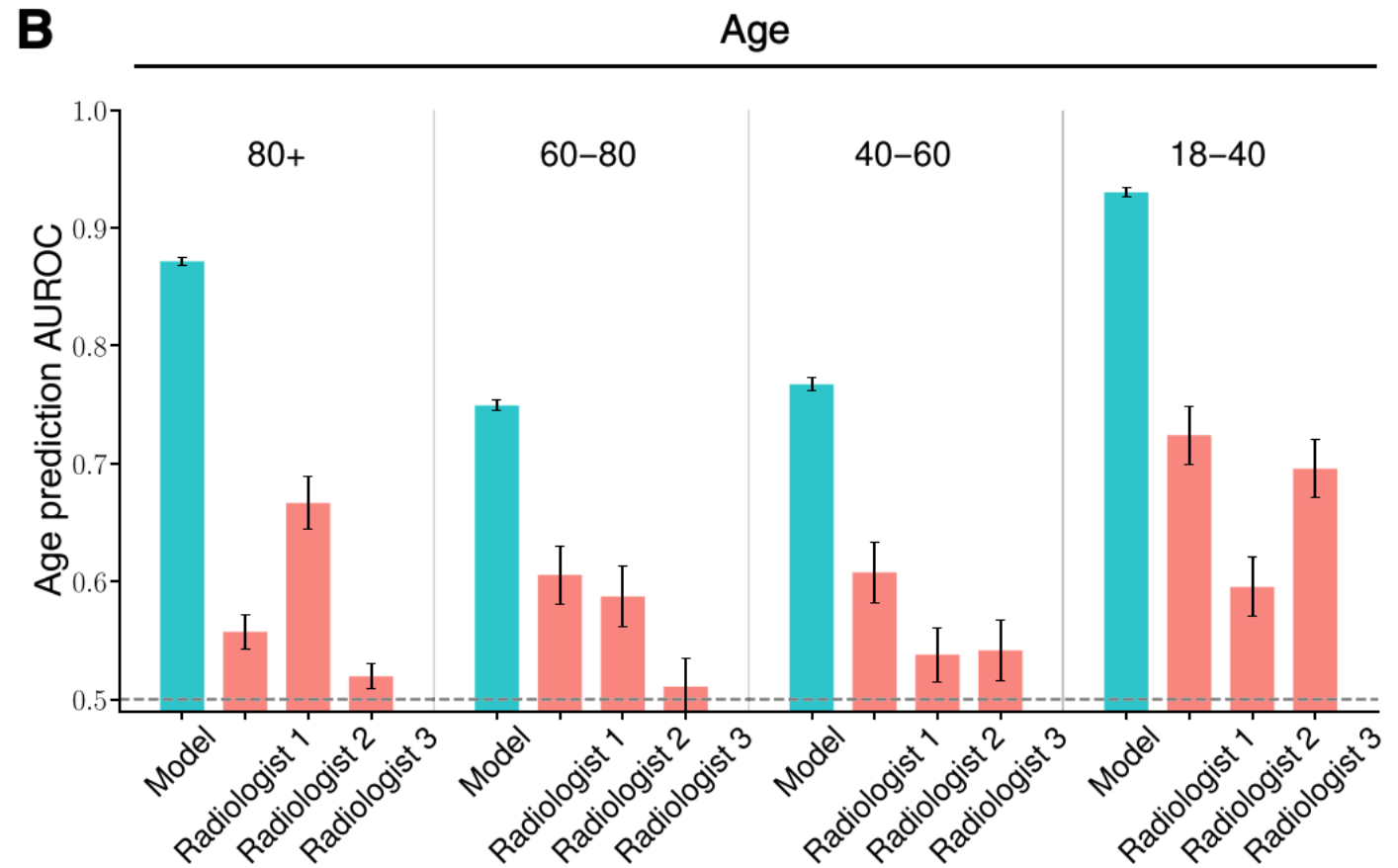
The goal: Develop a large **VLM** for
Mammography

Why SSL based VLMs

BioMedCLIP Zero-Shot Performance on MIMIC-CXR

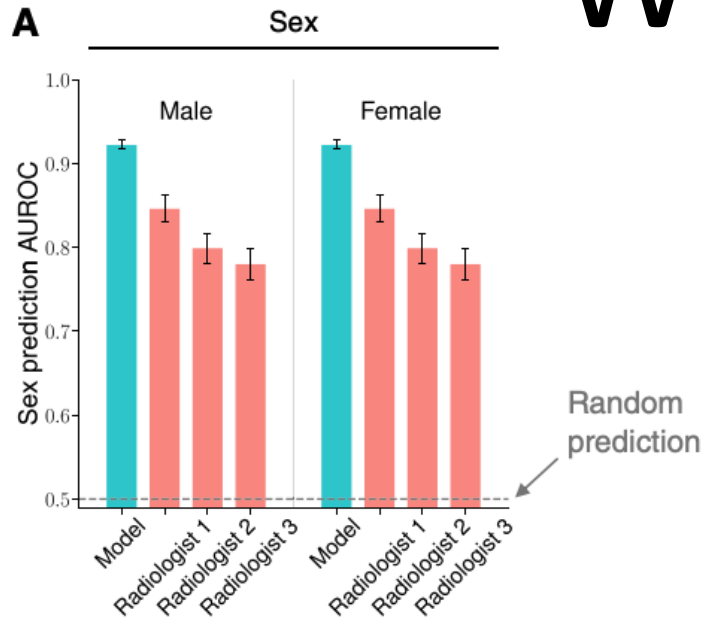


Why SSL based VLMs

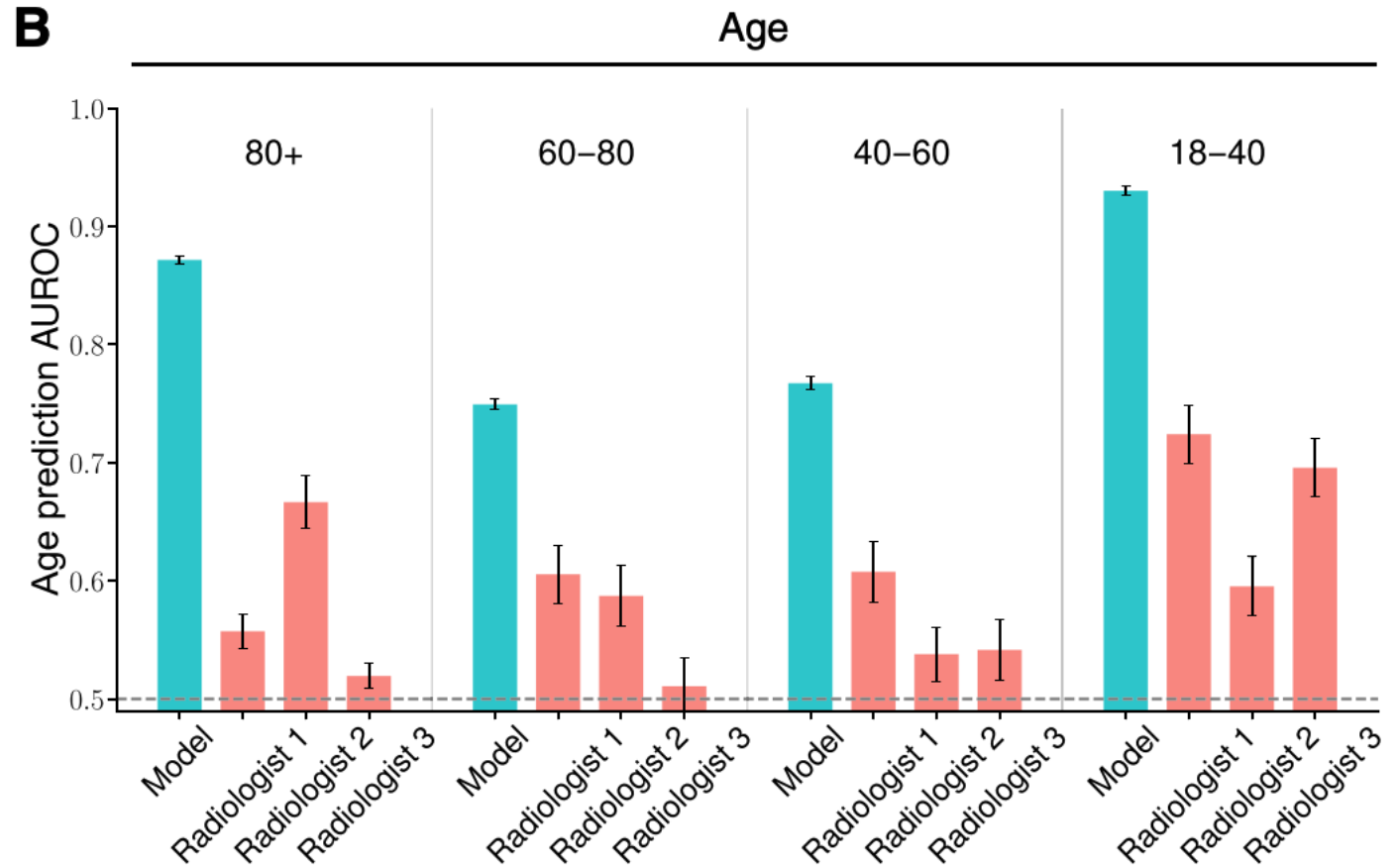


Why SSL based VLMs

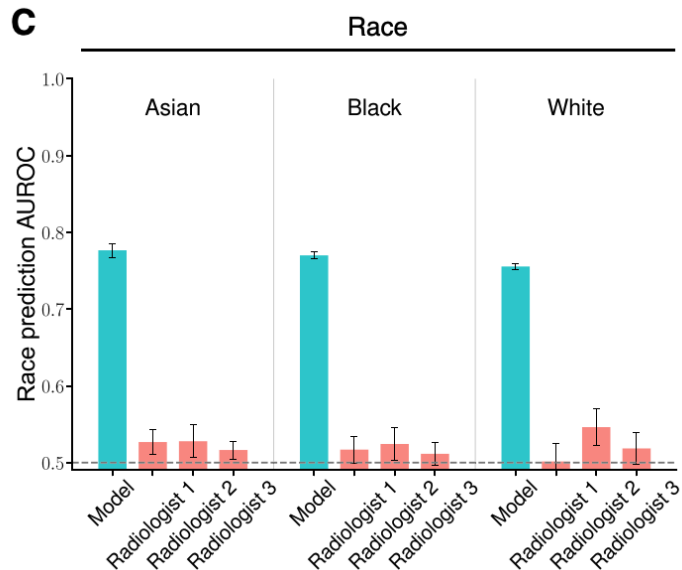
A



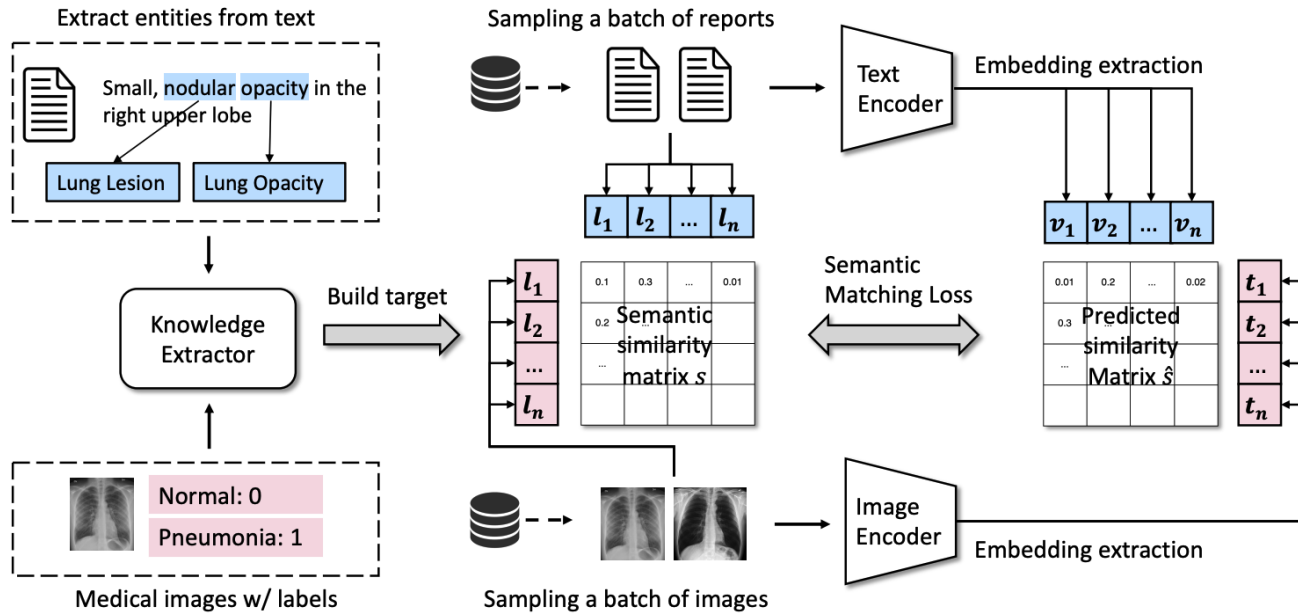
B



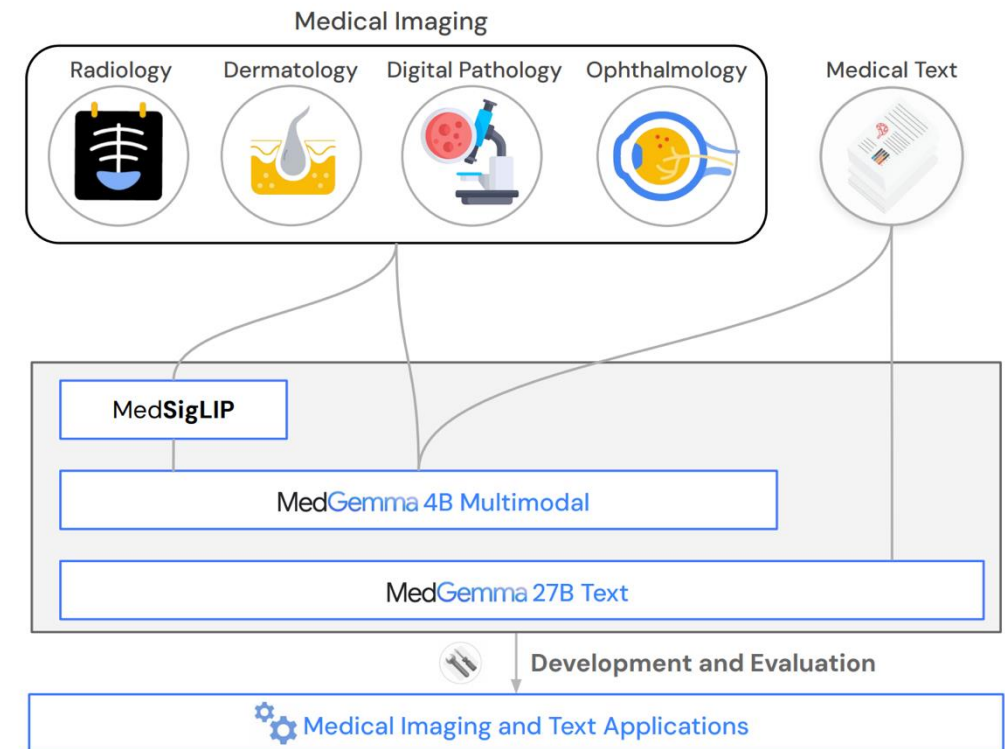
C



VLM in Mammography: Strategies

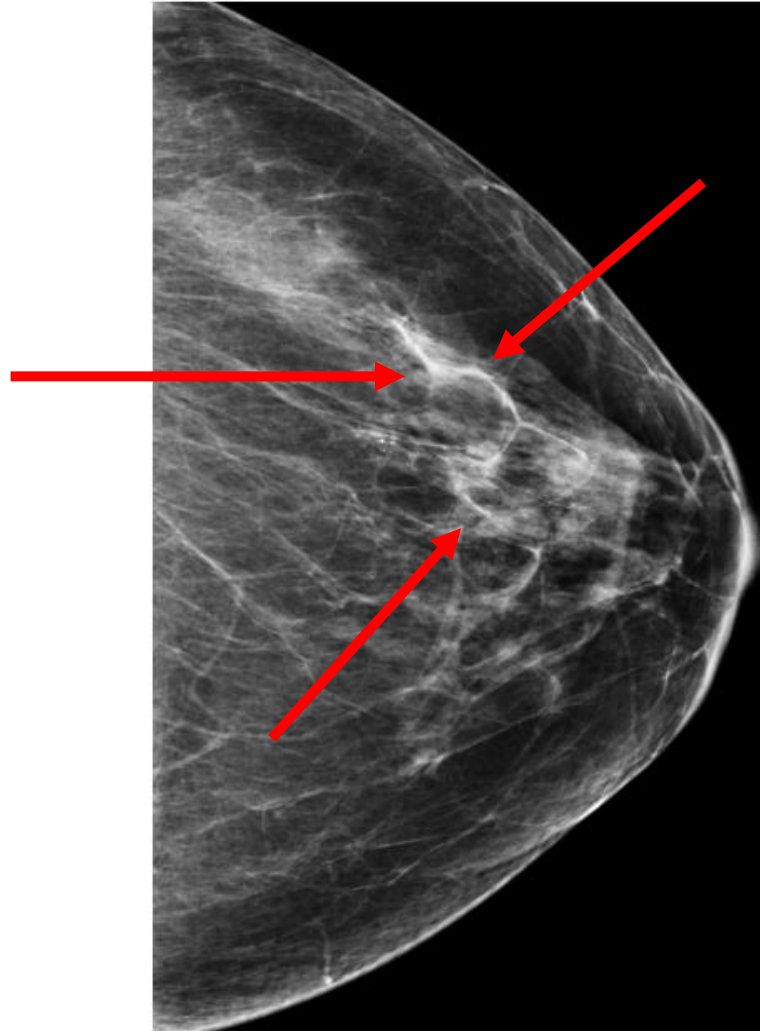


MedCLIP. EMNLP 2022



We need domain-specific VLMs

Mammograms are large and have subtle cues



VLMs can help

Sample screening mammogram report:

Ground truth: *Findings:* The breast tissue is heterogeneously dense, which could obscure detection of small masses. There are calcifications in the left upper, slightly outer breast. There is also an asymmetry in the left breast. Otherwise, no suspicious masses, clustered microcalcifications or areas of architectural distortion are seen. *Impression:* BI-RADS: 0. Calcifications and asymmetry in the left breast, which needs additional imaging.

Mammo-FM: pretraining

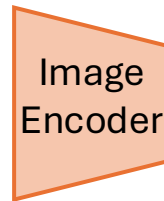
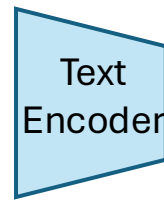
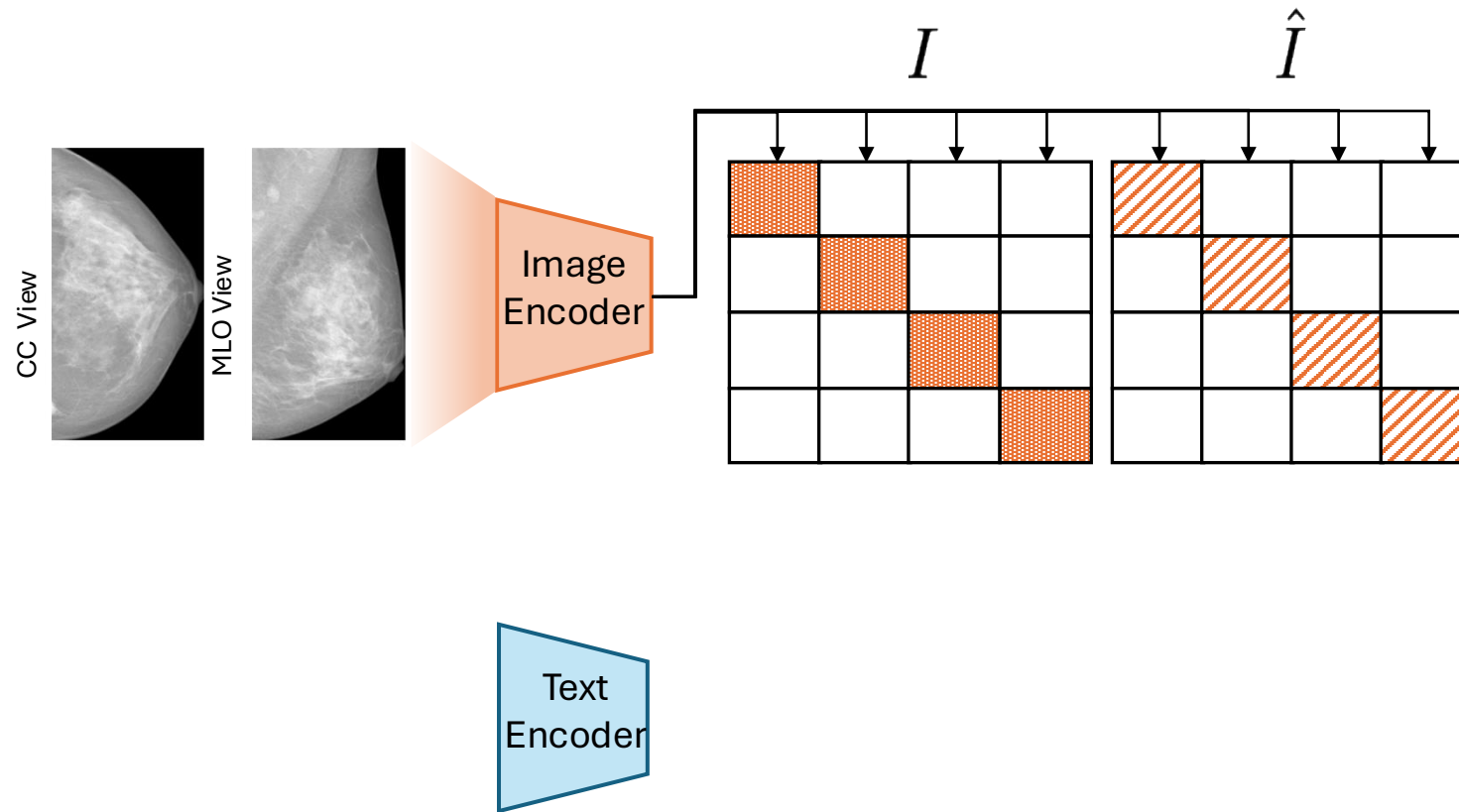


Image
Encoder

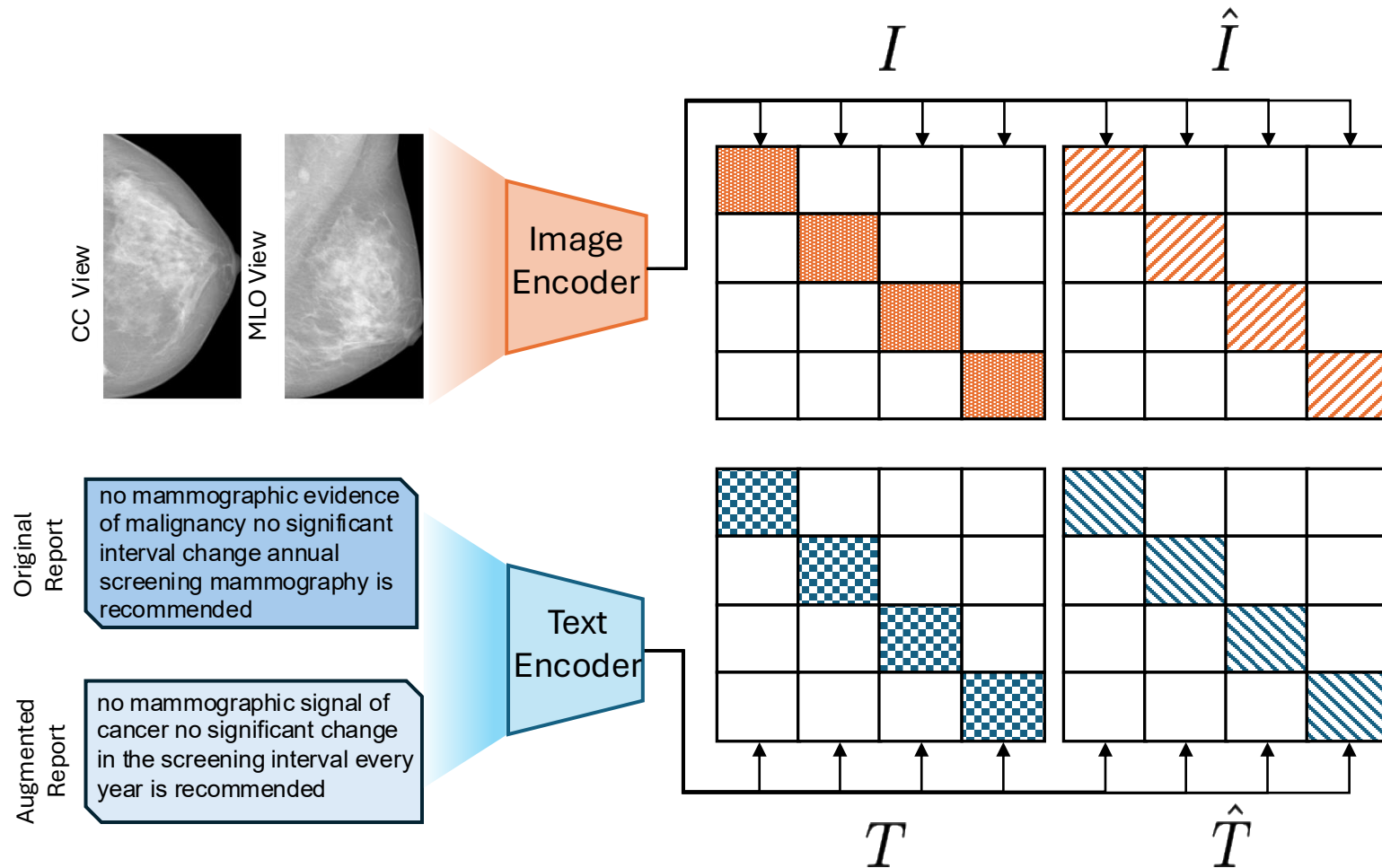


Text
Encoder

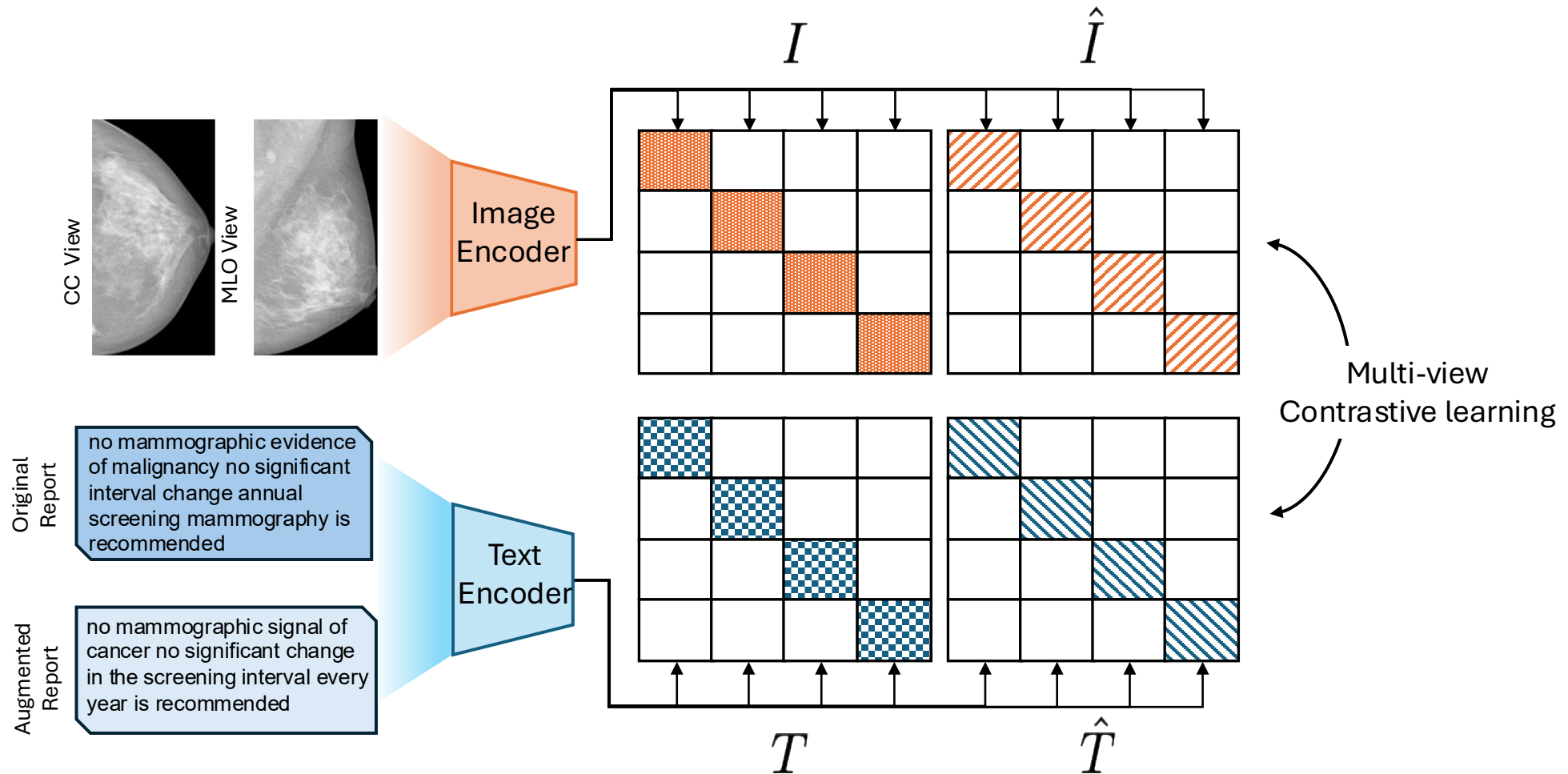
Mammo-FM: pretraining



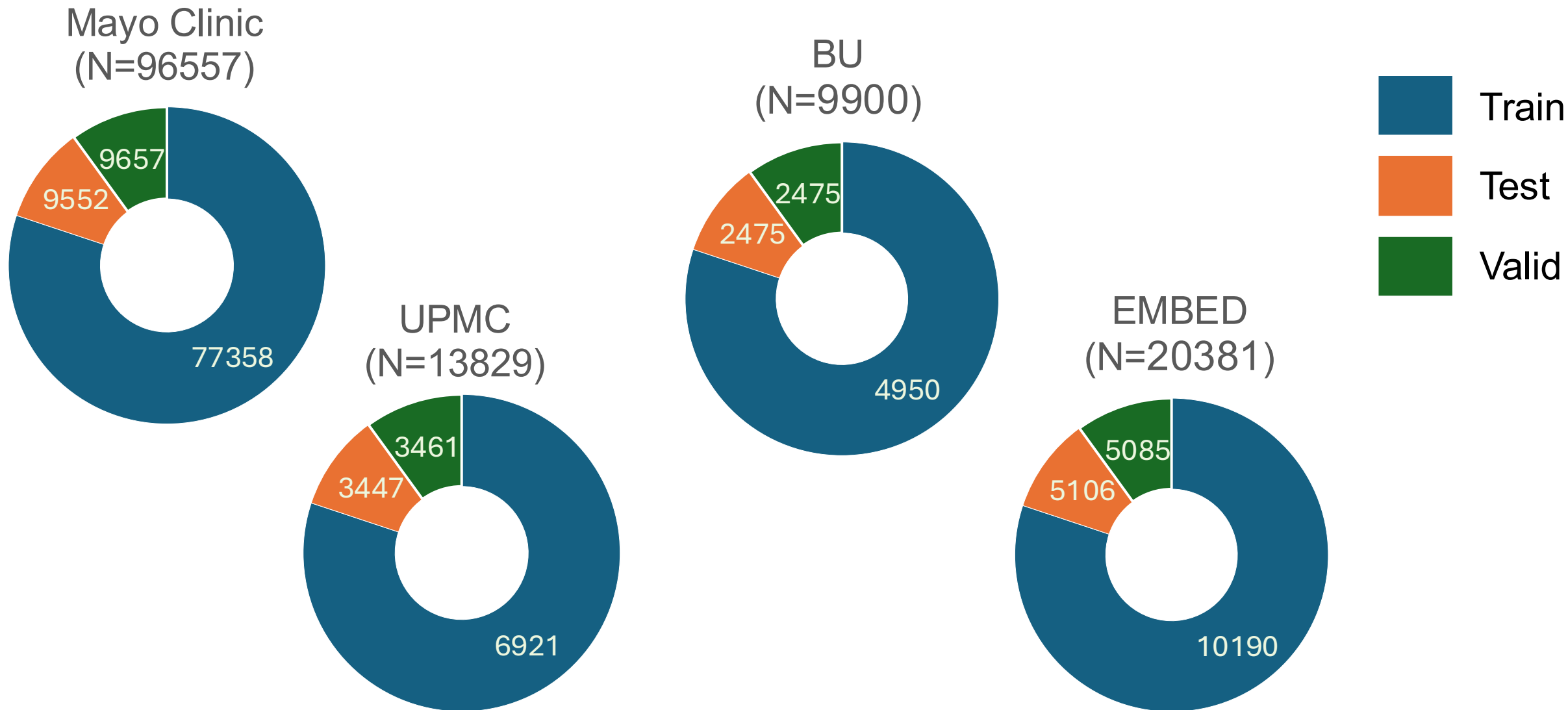
Mammo-FM: pretraining



Mammo-FM: pretraining



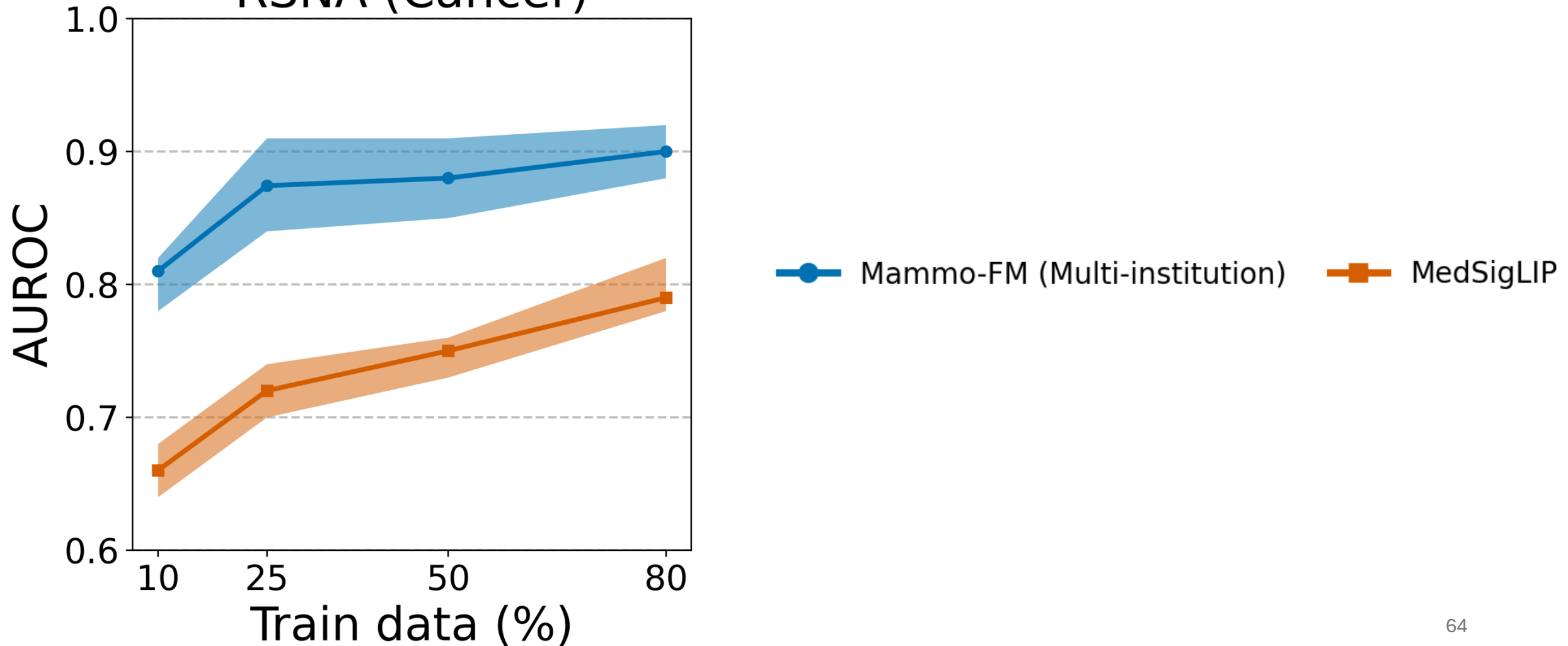
Mammo-FM: pretraining



Mammo-FM: data efficiency

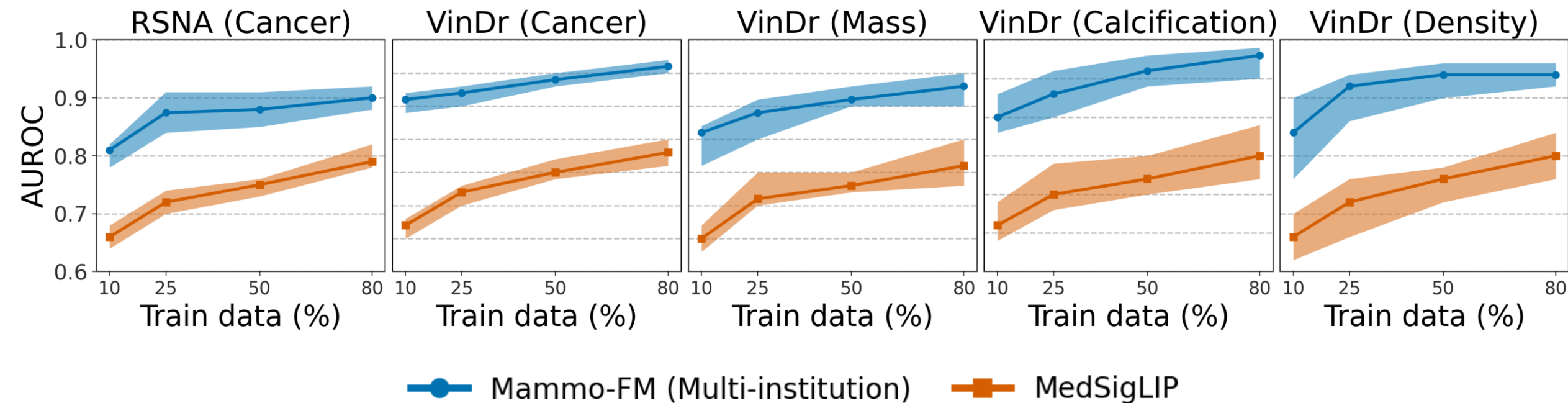
Out of distribution evaluation

RSNA (Cancer)



Mammo-FM: data efficiency

Out of distribution evaluation



Conclusion from Aim 2

1. Robust mammography features learned: We learn generalized features for mammography by pre-training on the **largest** and **most diverse** mammography datasets. **[MICCAI 2024. Top 11%]**

It's diagnostic performance is better than the SOTA **generalist industrial models**.

2. Two Applications: It helps to **interpret** the risks from any SOTA **risk predictors**.

We develop the **1st report generator** for mammography using Mammo-FM. **[ArXiv 2025]**

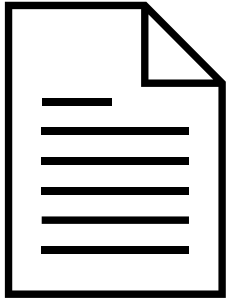
Aim 3

The goal: **Detect** the **Systematic Mistakes**
using **Language**

Aim 3

The goal: **Detect the Systematic Mistakes**
using **Language**

Captions



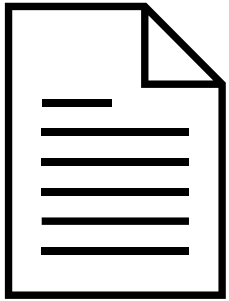
Reports



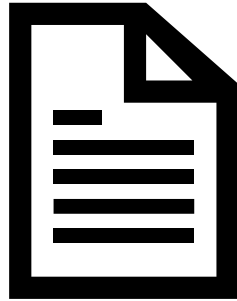
Aim 3

The goal: **Detect the Systematic Mistakes**
using **Language**

Captions



Reports



Vision-Language
Alignment

Natural image

CLIP (ICML 2020)

CXRs

1. **GLORIA**
(ICCV2021)

2. **MedCLIP**
(EMNLP 2022)

3. **CXR-CLIP**
(MICCAI 2023)

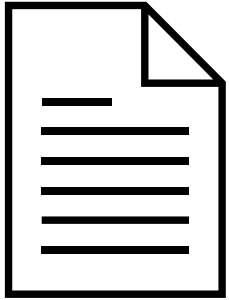
Mammo



Aim 3

The goal: **Detect the Systematic Mistakes**
using **Language**

Captions



Reports



Vision-Language
Alignment

Natural image

CLIP (ICML 2020)

CXRs

1. ***GLORIA***
(ICCV2021)

2. ***MedCLIP***
(EMNLP 2022)

3. ***CXR-CLIP***
(MICCAI 2023)

Mammo

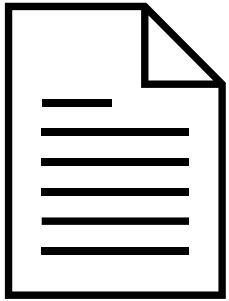
Mammo-FM
(MICCAI 2024,
ArXiv 2025)

Aim 3

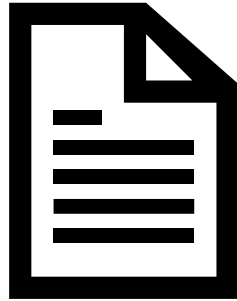
(Visual + non-Visual)

The goal: **Detect the Systematic Mistakes**
using **Language**

Captions



Reports



Vision-Language
Alignment

Natural image

CLIP (ICML 2020)

CXRs

1. **GLORIA**
(ICCV2021)

2. **MedCLIP**
(EMNLP 2022)

3. **CXR-CLIP**
(MICCAI 2023)

Mammo

Mammo-FM
(MICCAI 2024,
ArXiv 2025)

Tracing non-visual mistakes

Population



Age: [32-88]
Race: 80% Non-Hispanic
White, 20% Asian
....

Individual



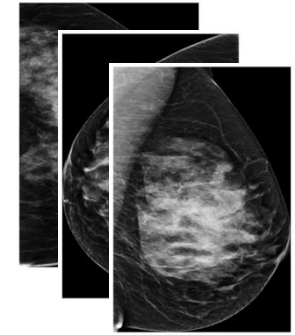
Reason for Visit: [....]
Blood Pressure: [...]
Lab Test: [....]
...

Site



Manufacturer: [....]
X-ray Dosage: [...]
Aperture Setting: [....]
...

Preprocessing



Photometric Interpretation: [
Monochrome 1 vs Monochrome2]
Crop ratio: [...]
...

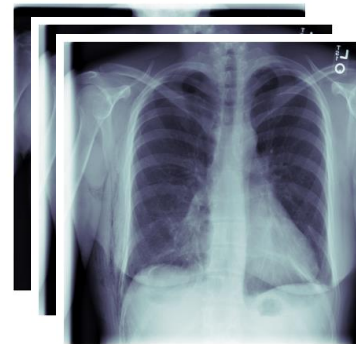
How a human would do?



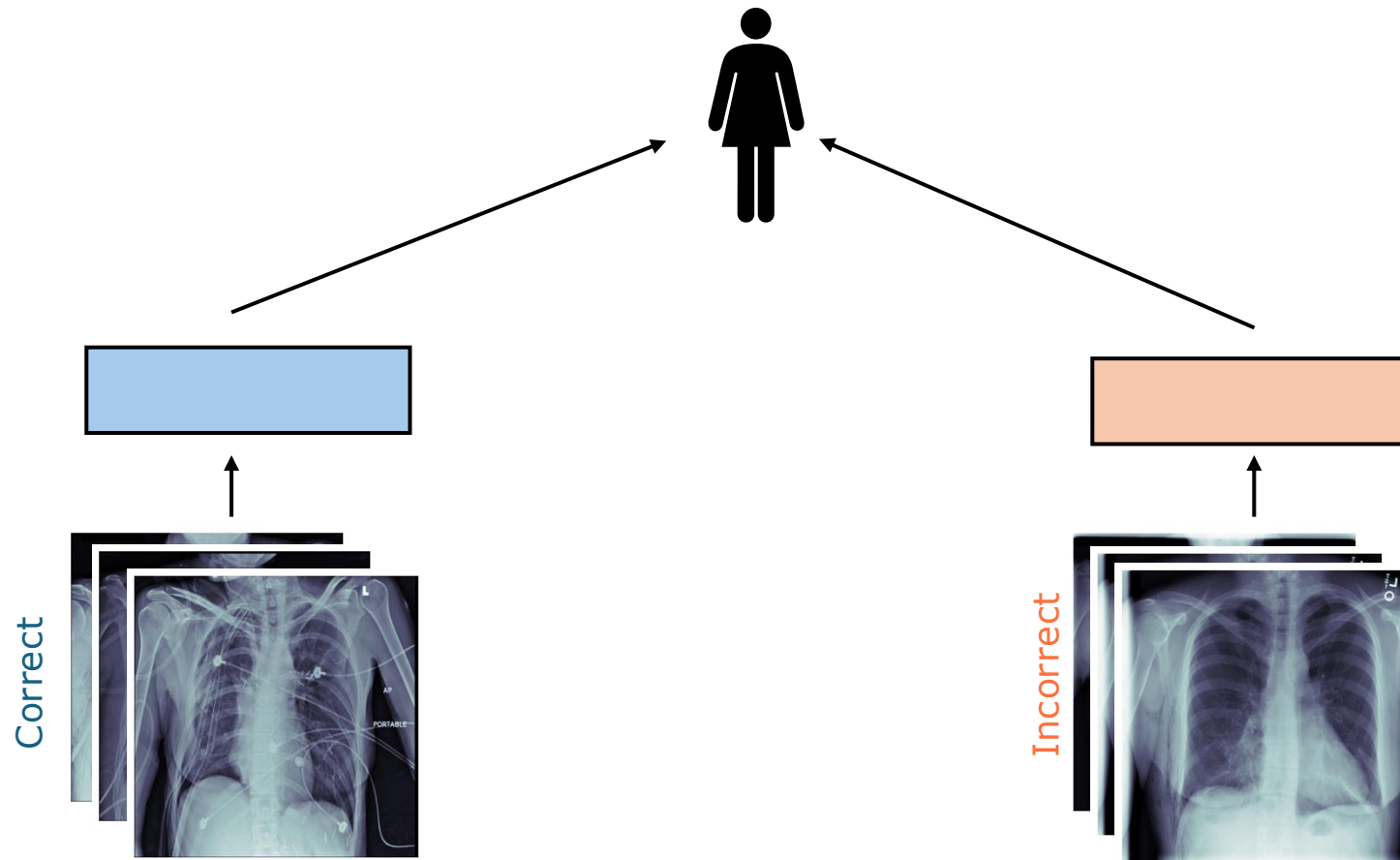
Correct



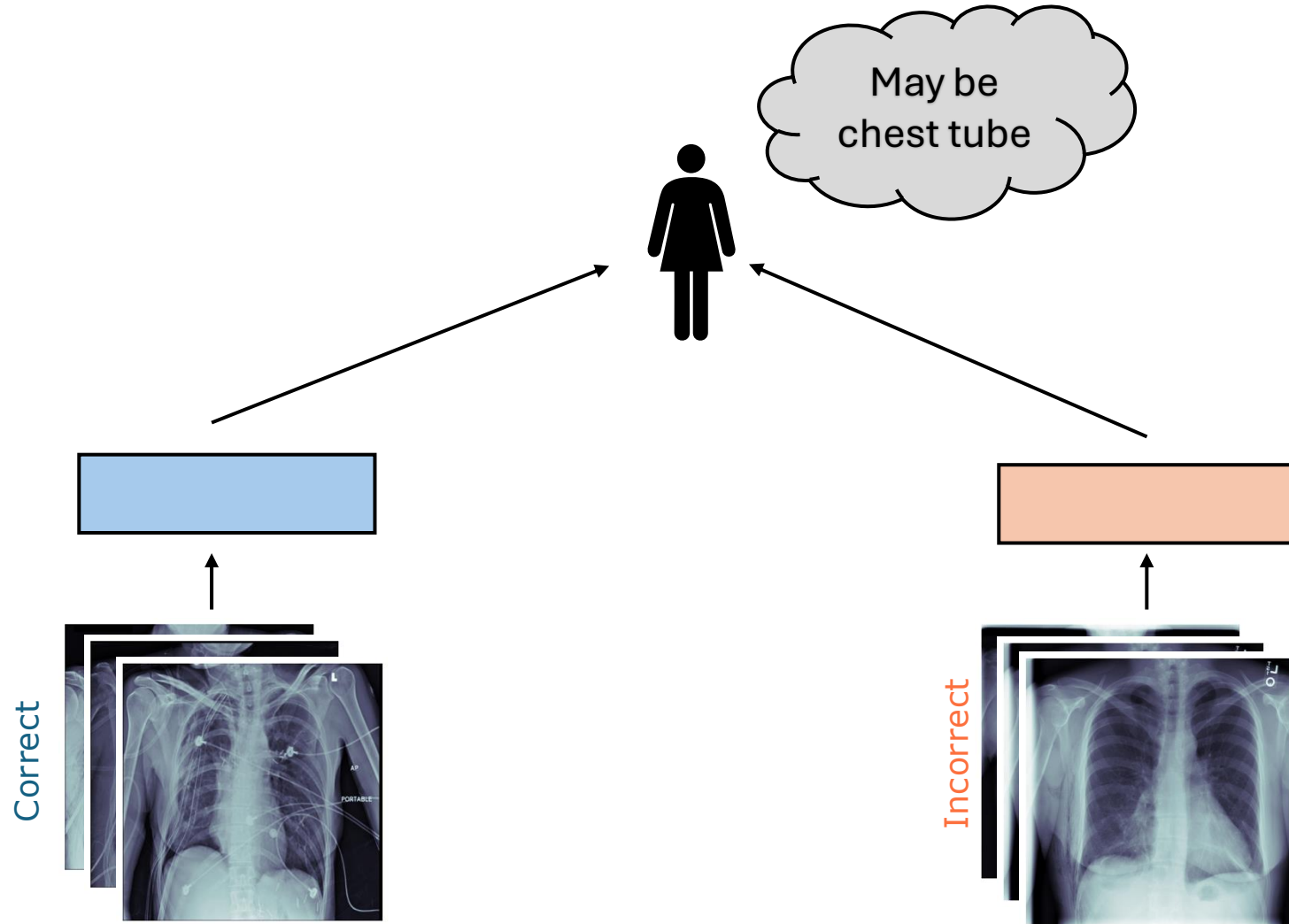
Incorrect



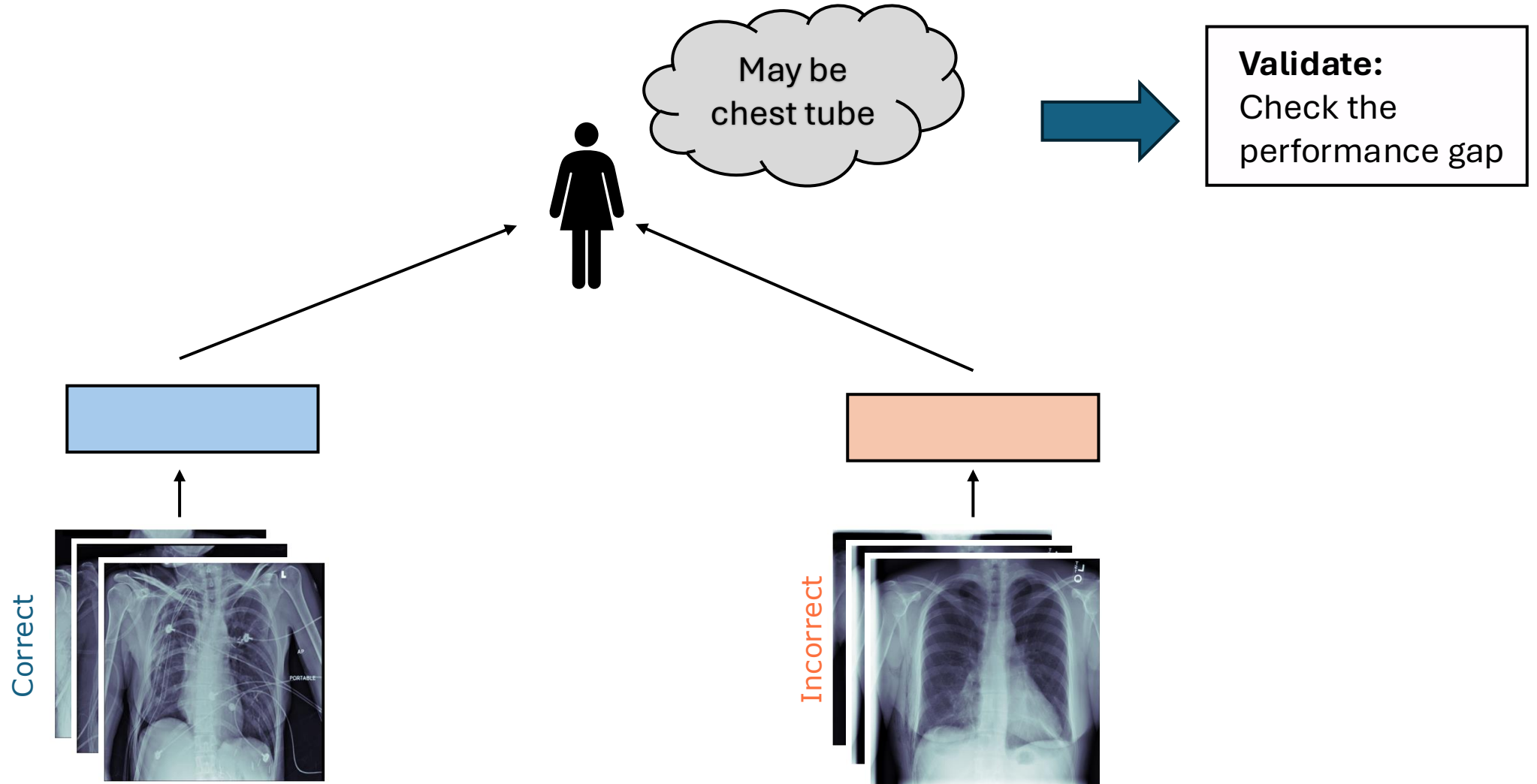
How a human would do?



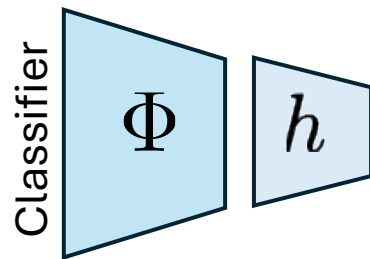
How a human would do?



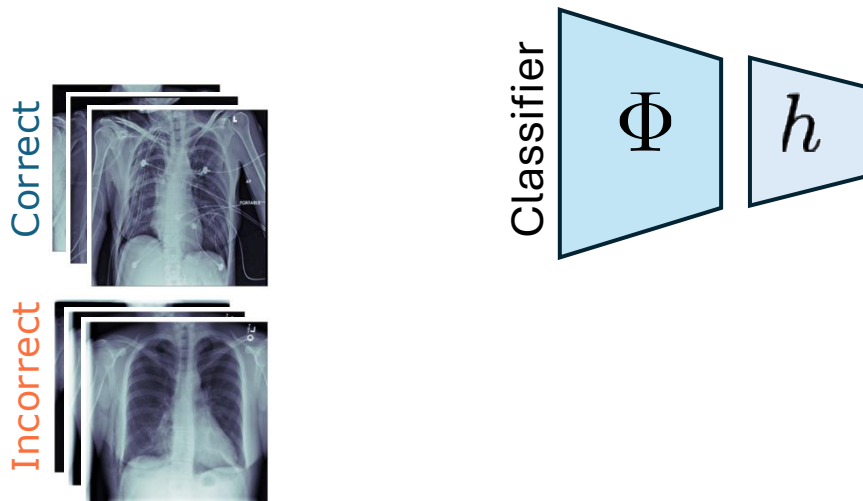
How a human would do?



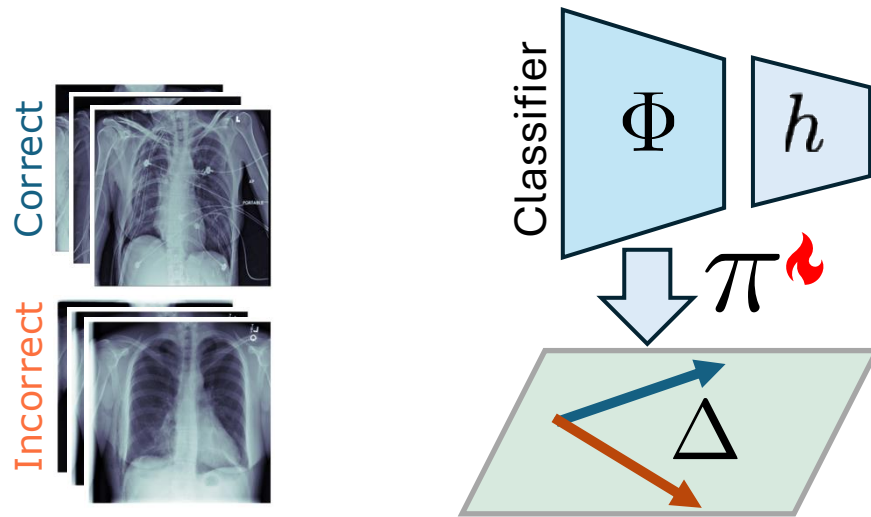
Ladder



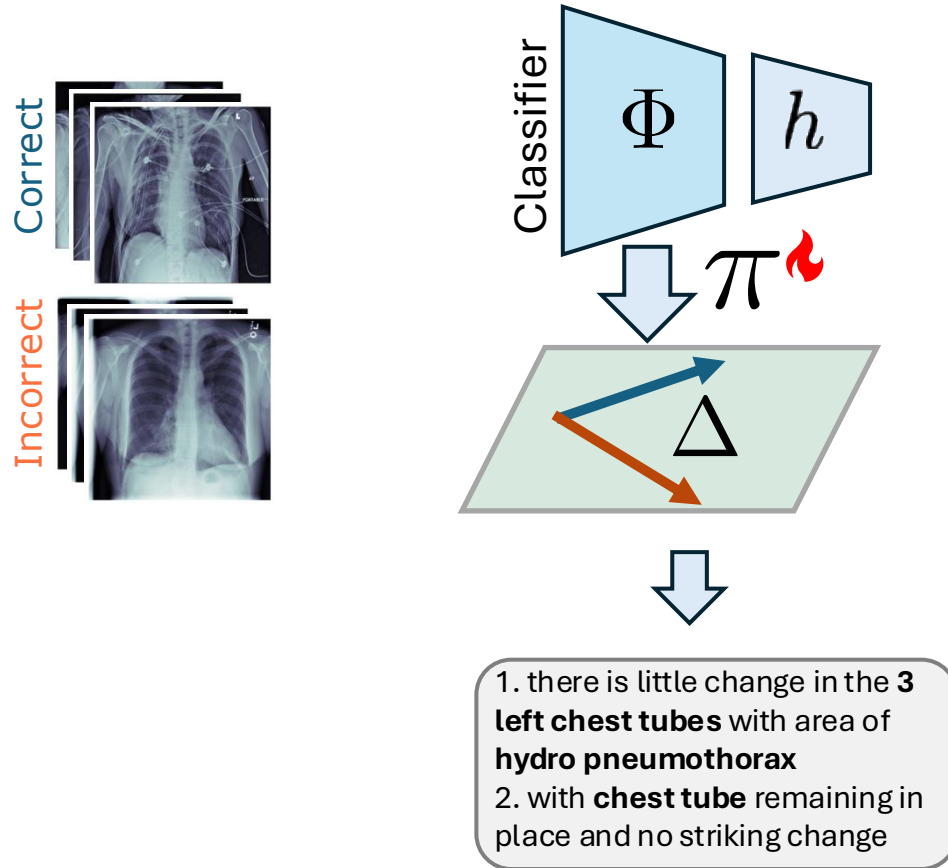
Ladder



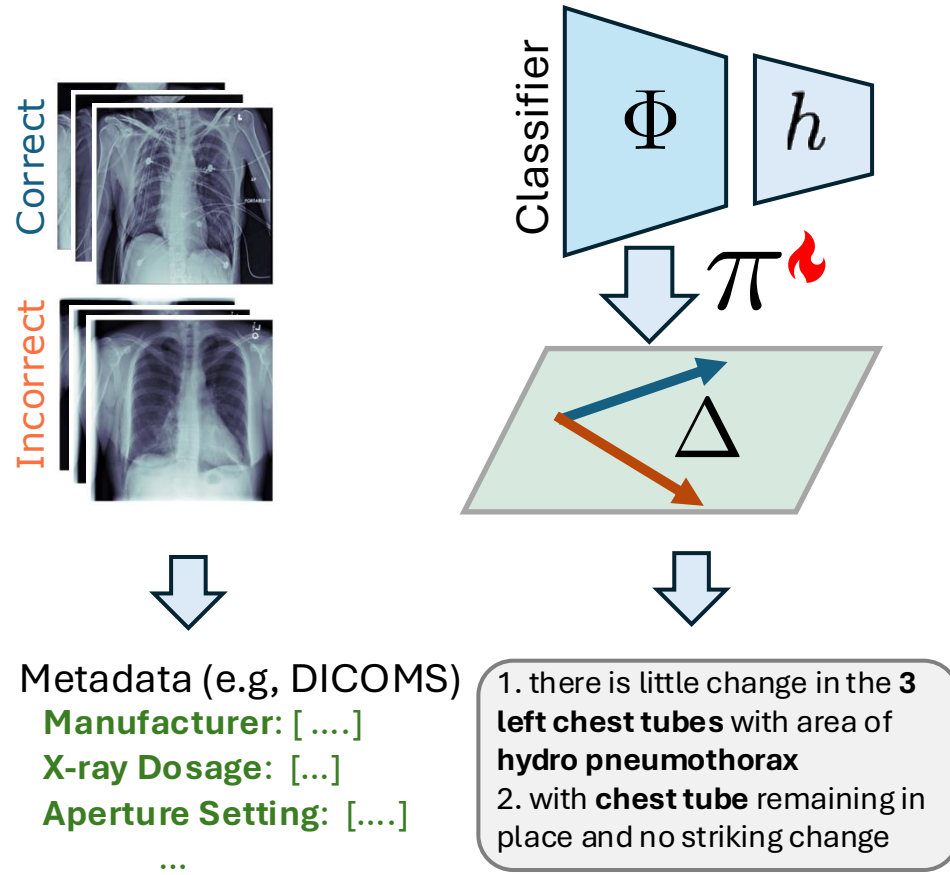
Ladder



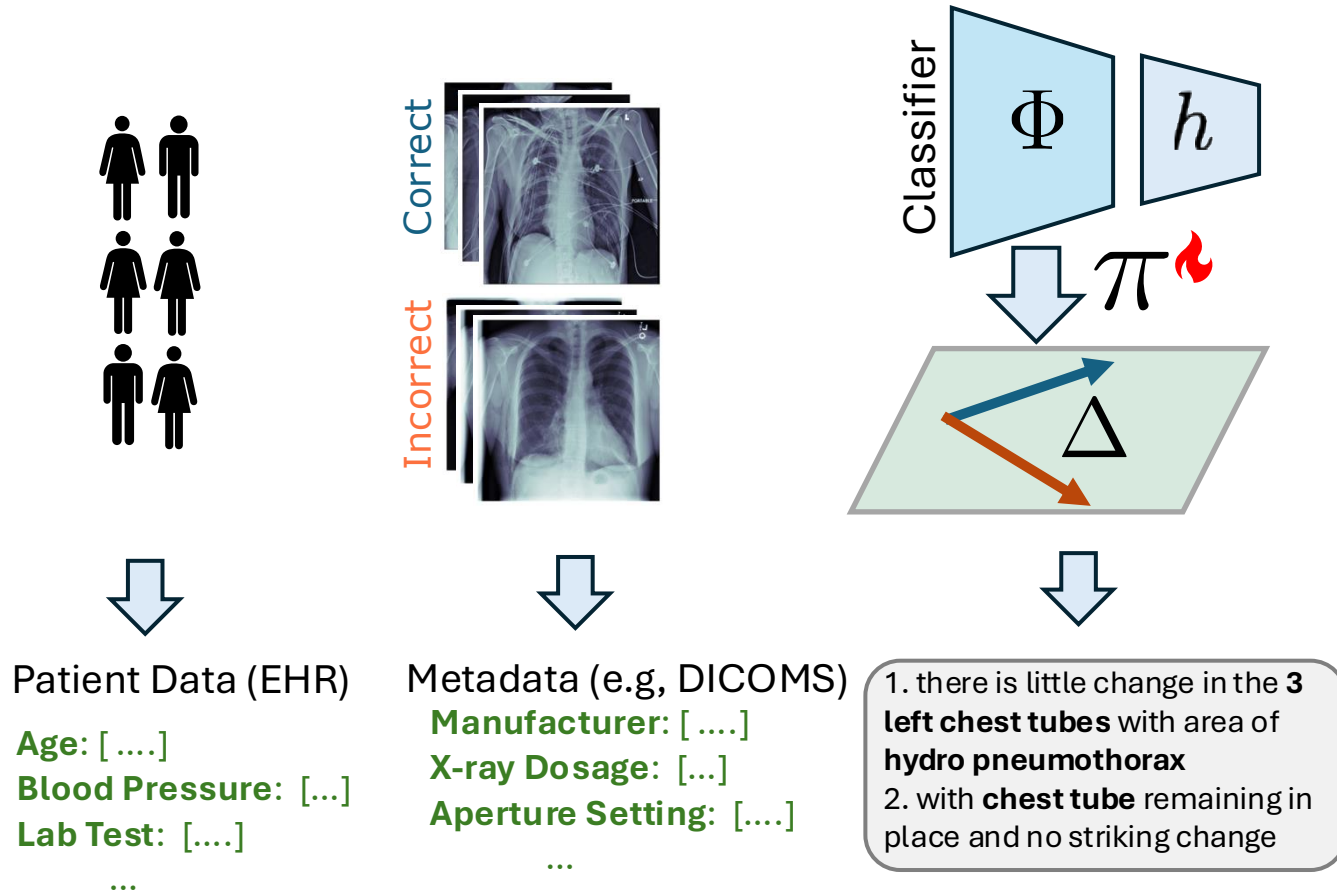
Ladder



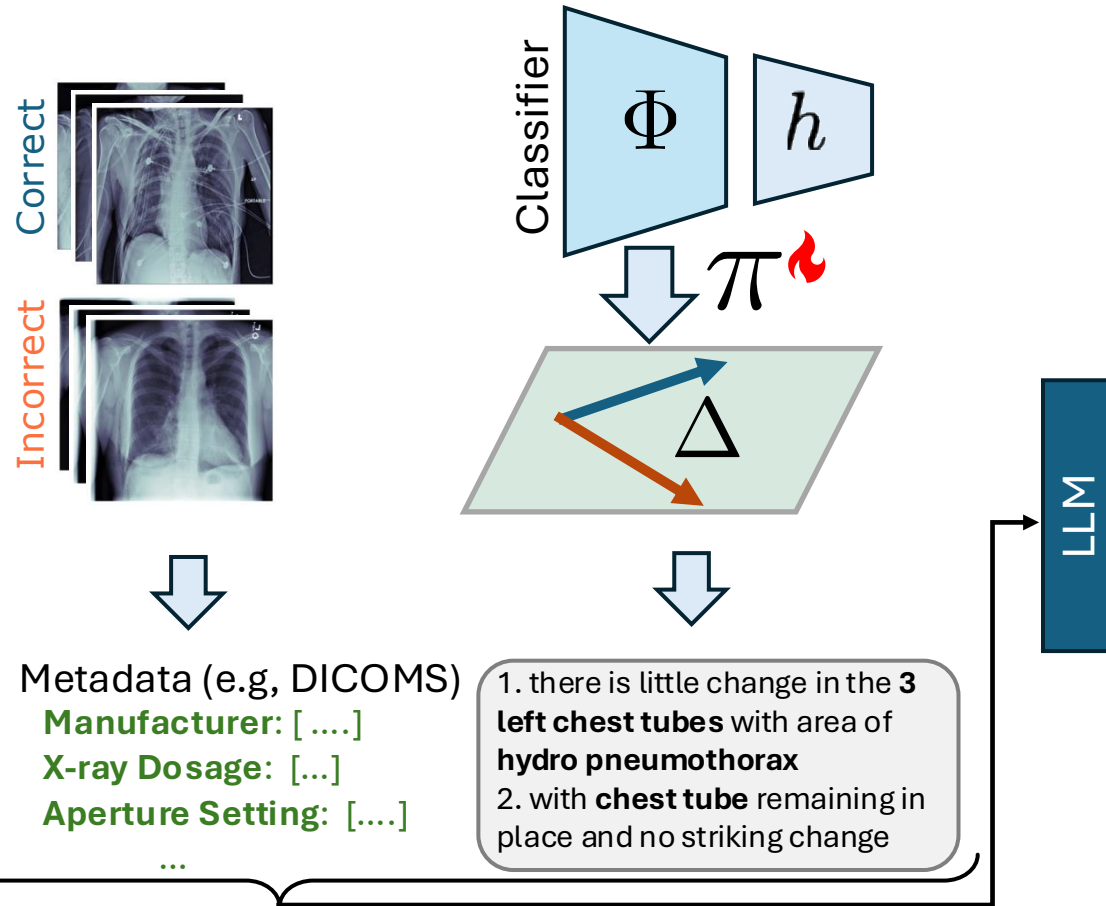
Ladder



Ladder

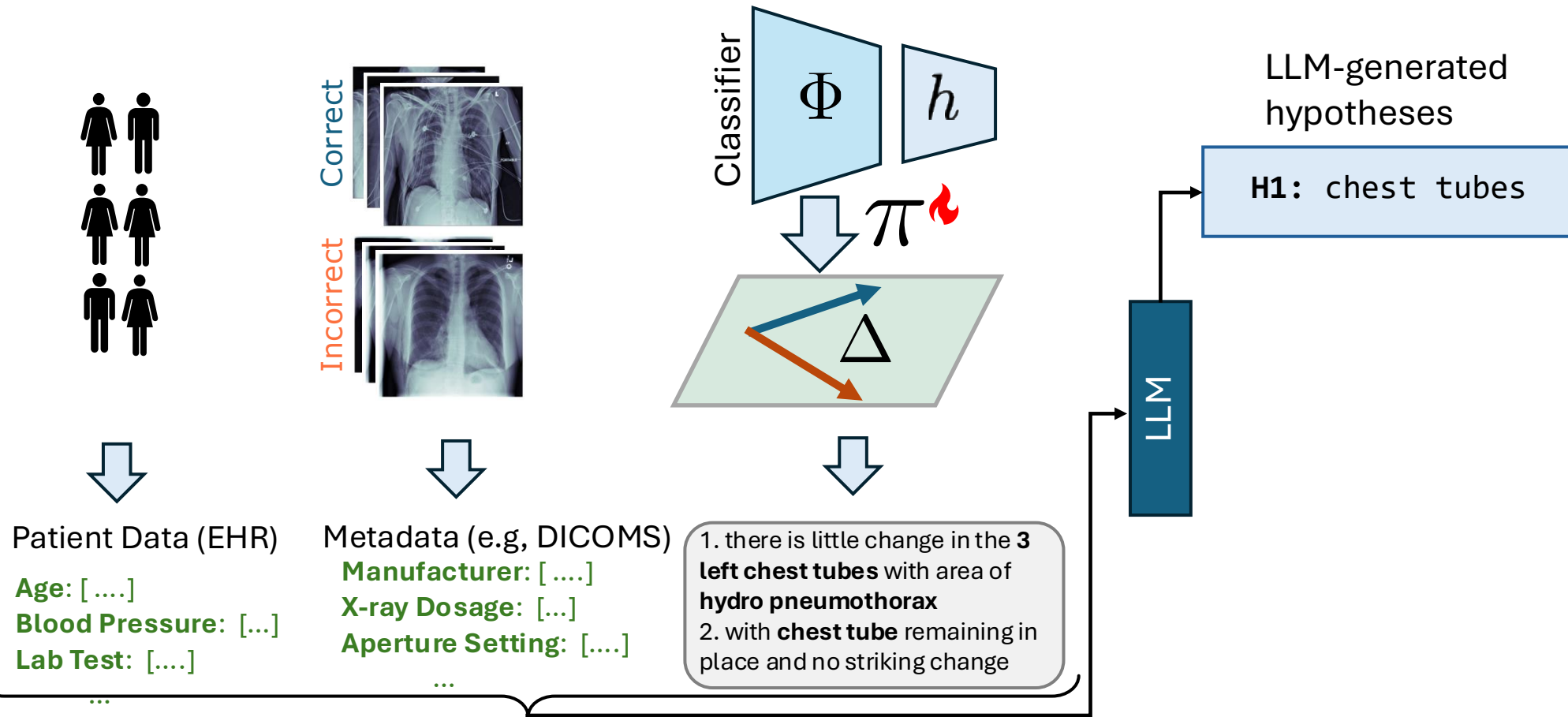


Ladder



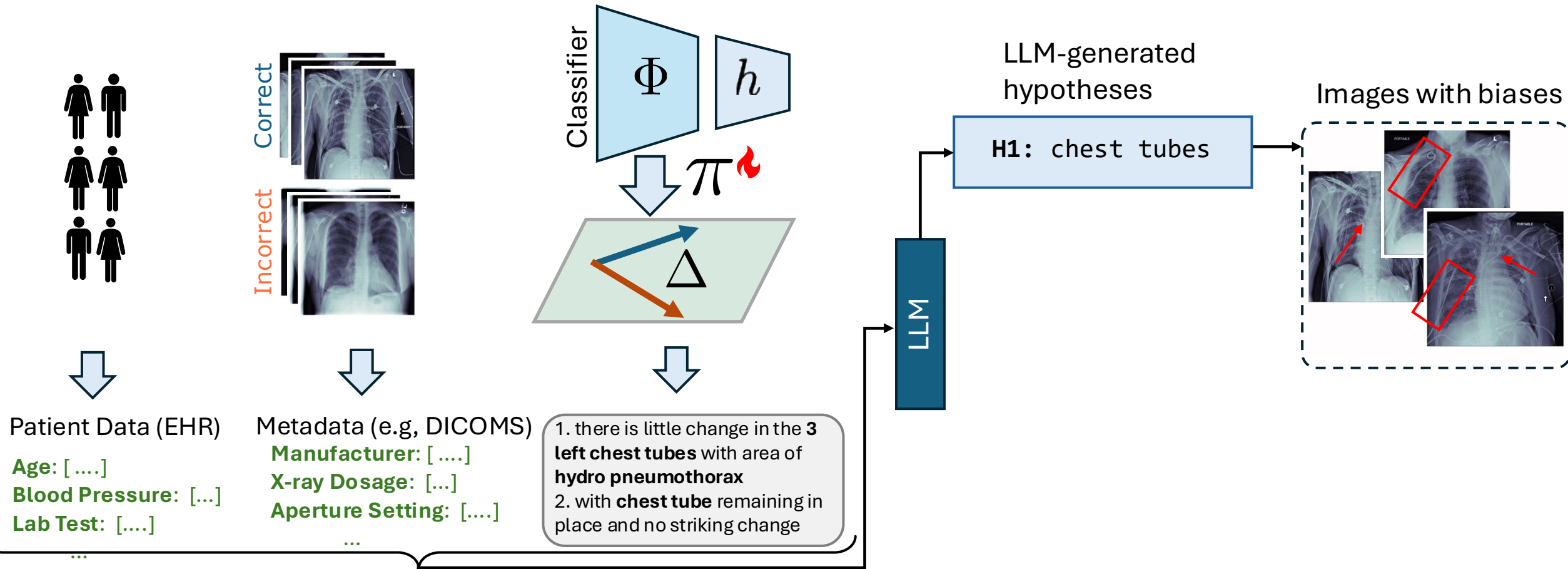
The DICOMS and EHR is an illustrative example here. NIH-CXR does not have that. We use RSNA (Cancer) dataset to perform this experiment

Ladder



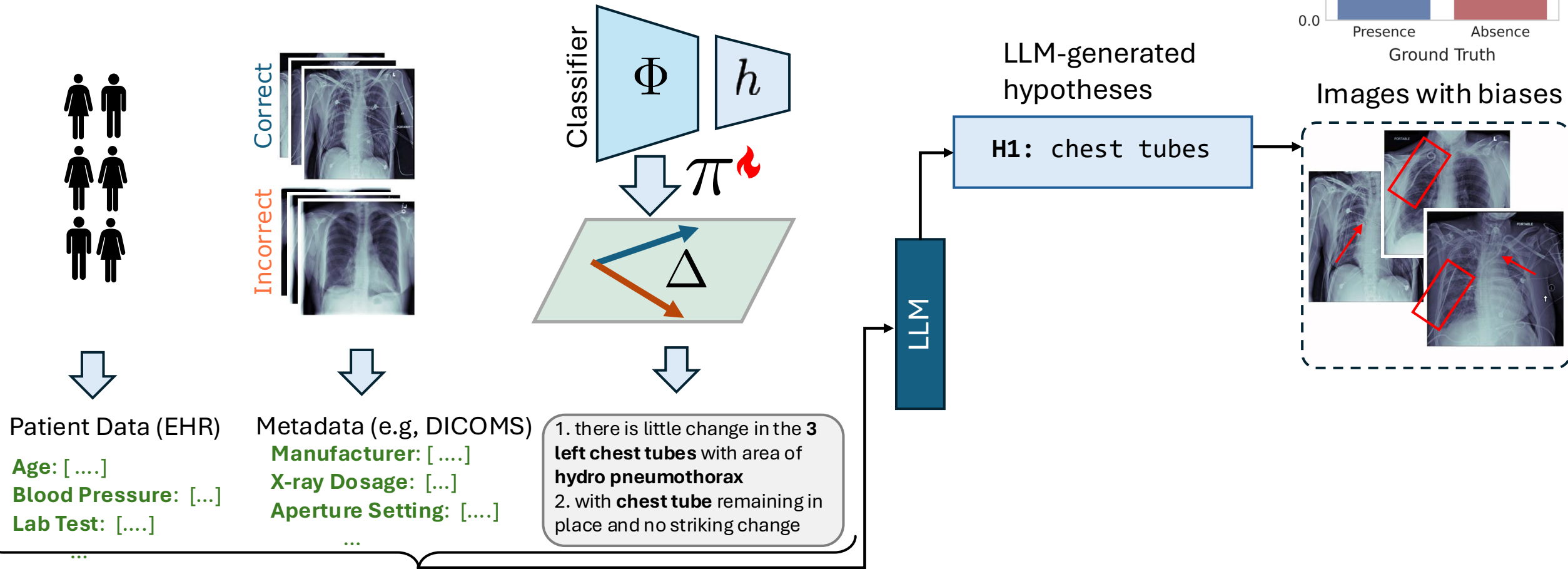
The DICOMS and EHR is an illustrative example here. NIH-CXR does not have that. We use RSNA (Cancer) dataset to perform this experiment

Ladder



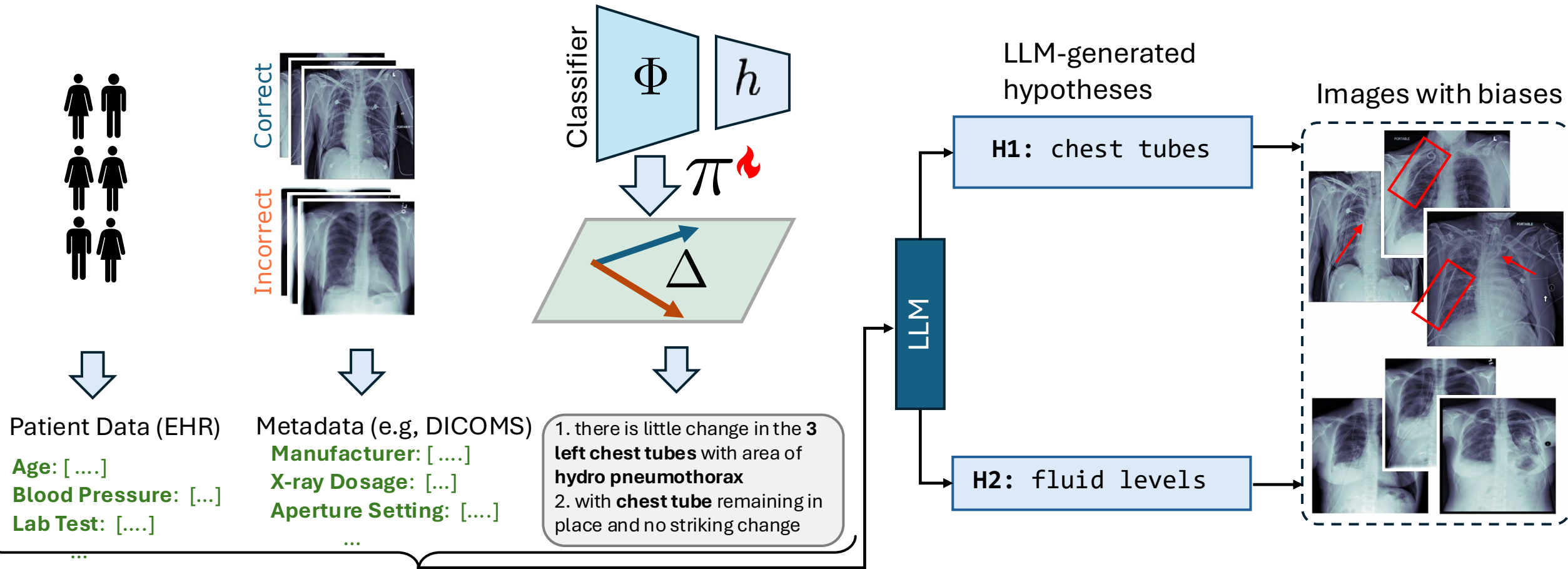
The DICOMS and EHR is an illustrative example here. NIH-CXR does not have that. We use RSNA (Cancer) dataset to perform this experiment

Ladder



The DICOMS and EHR is an illustrative example here. NIH-CXR does not have that. We use RSNA (Cancer) dataset to perform this experiment

Ladder



The DICOMS and EHR is an illustrative example here. NIH-CXR does not have that. We use RSNA (Cancer) dataset to perform this experiment

Qualitative Results

Dataset: **Waterbirds**

Target: **Waterbird**

Bias: Specific background elements like docks and boats



Qualitative Results

Dataset: **Waterbirds**

Target: **Waterbird**

Bias: Specific background elements like docks and boats



Presence: **97.2 %**

Classifier Performance

Absence: **68.8%**

Qualitative Results

Dataset: **NIH**

Target: **Pneumothorax**

Bias: Chest tubes

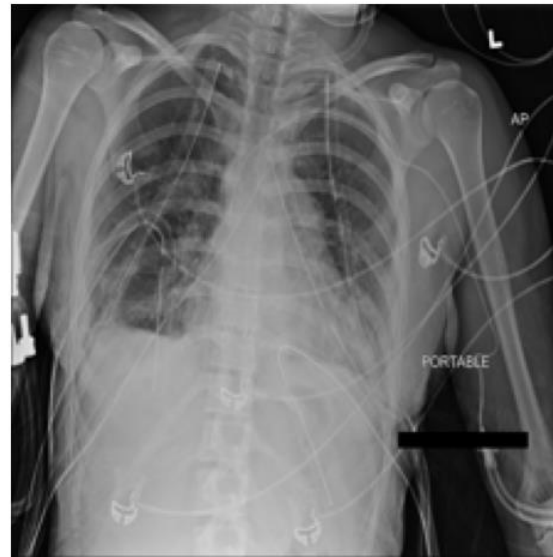


Qualitative Results

Dataset: **NIH**

Target: **Pneumothorax**

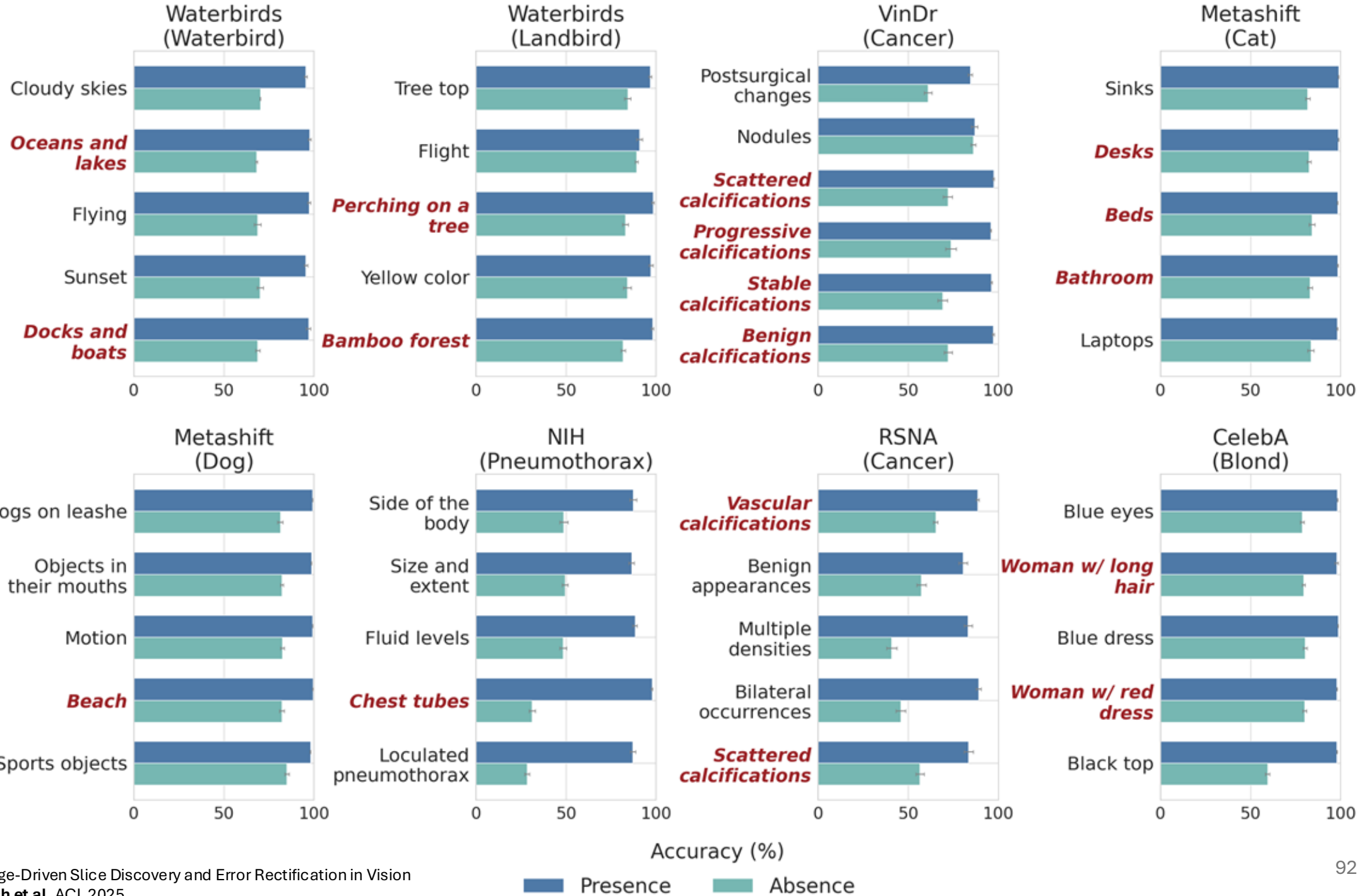
Bias: Chest tubes



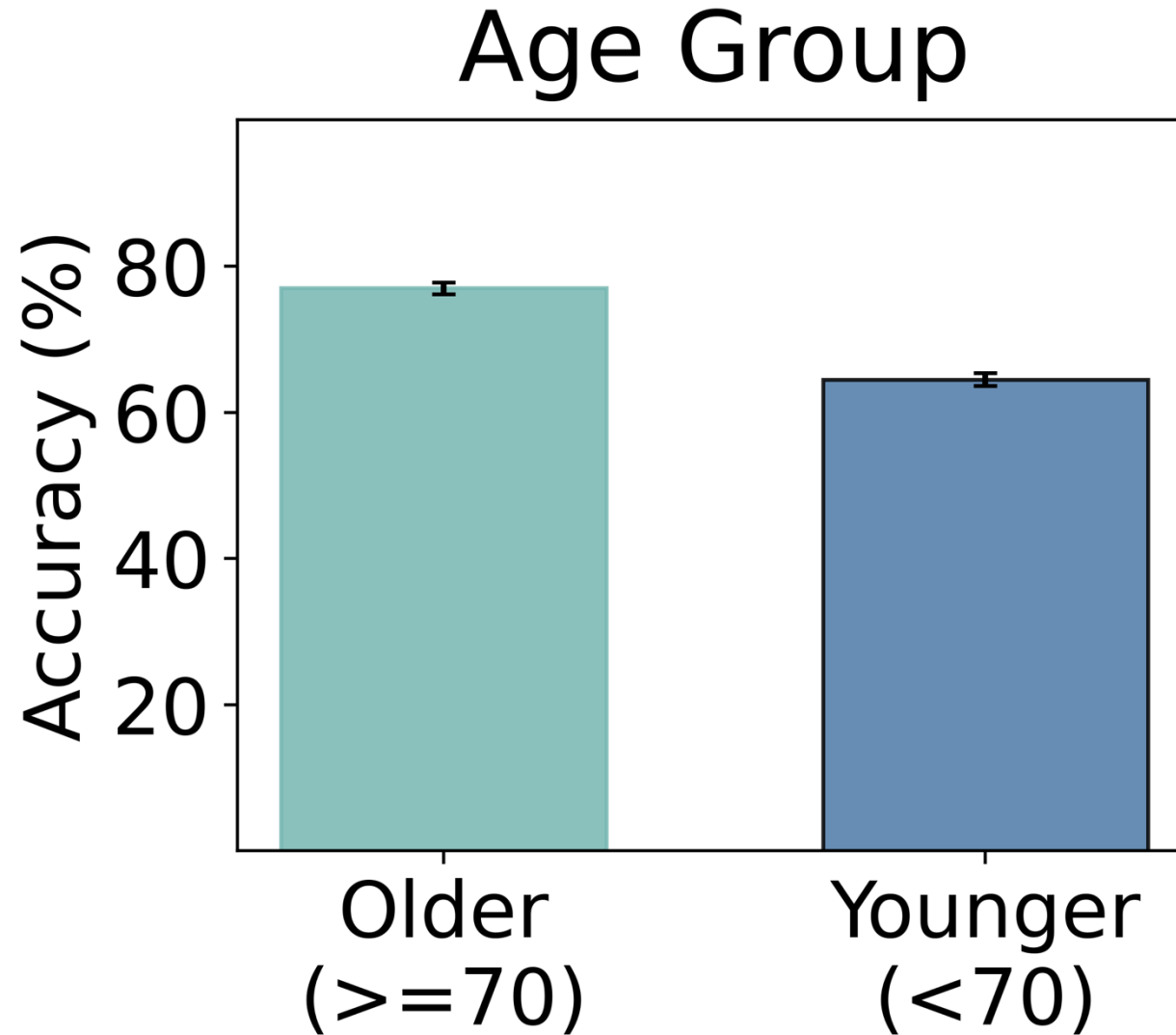
Presence: **76.2 %**

Classifier Performance

Absence: **34.8%**



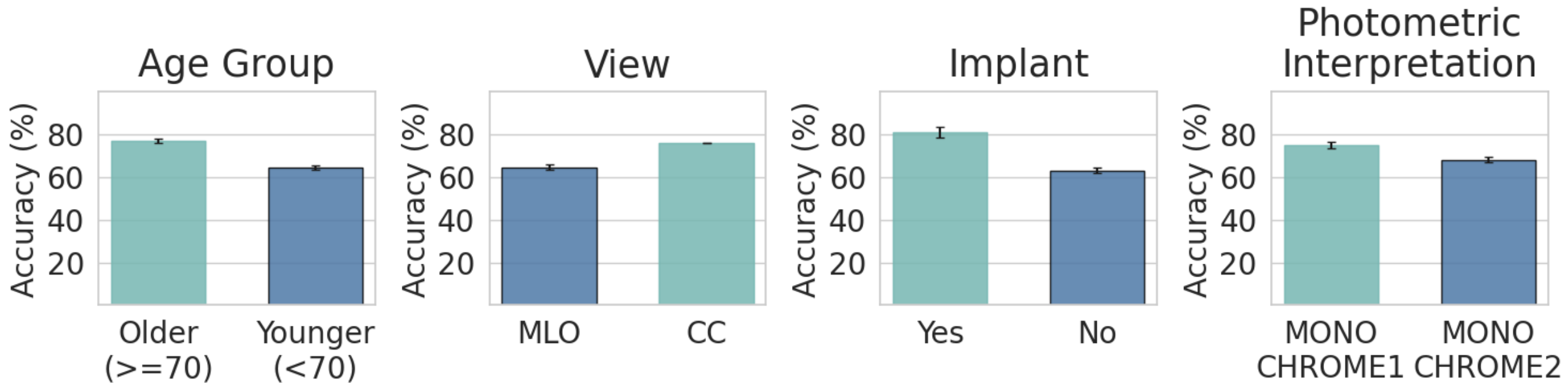
Detecting non-visual mistakes

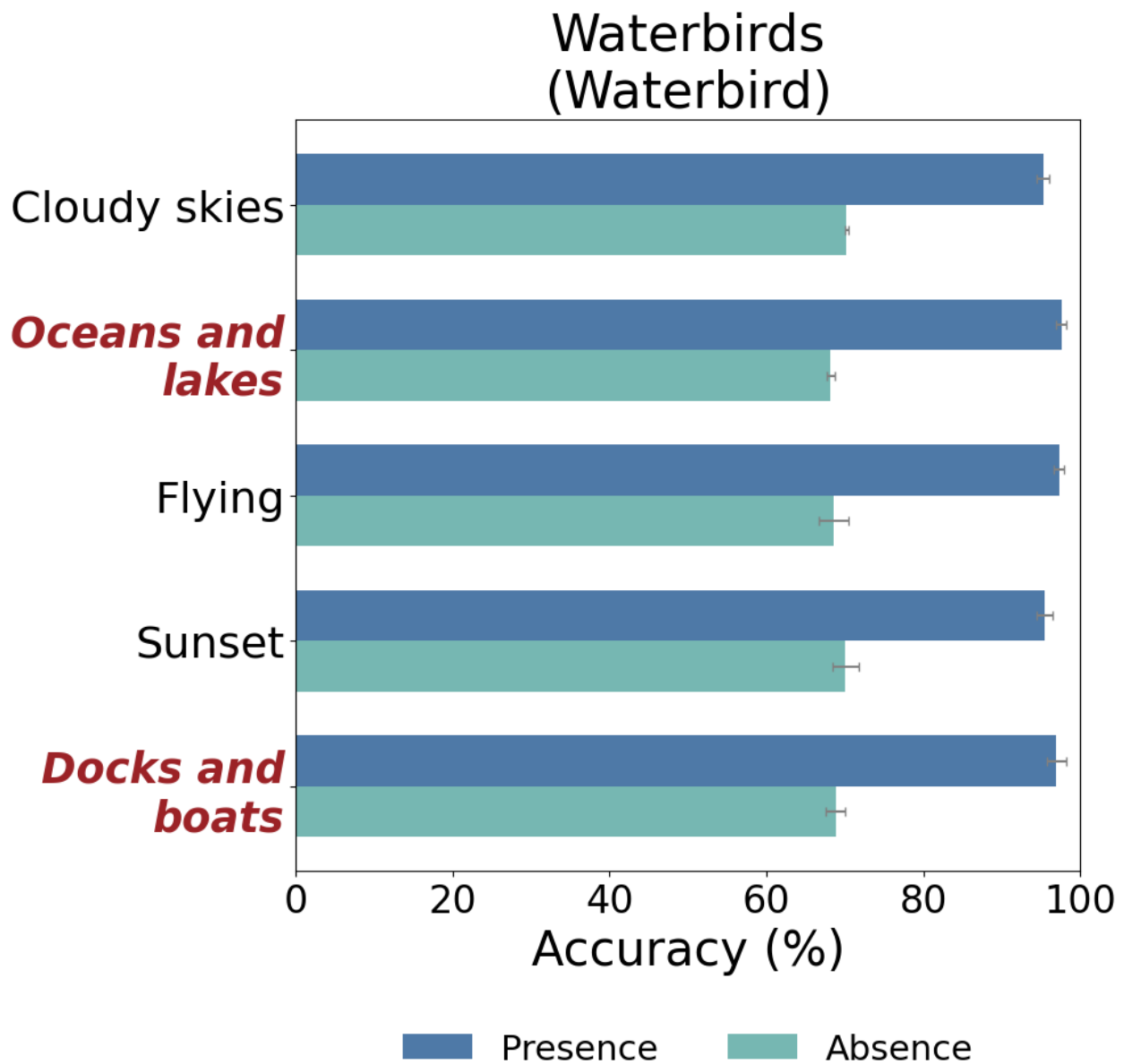


Detecting non-visual mistakes

Dataset: **RSNA**

Target: **Breast cancer**





$$\underbrace{\{A_1, A_2, \dots, A_k\}}_{\substack{\text{Explained by data} \\ \text{distribution shift} \\ \textbf{(Reject the hypotheses)}}}, \underbrace{\{A_{k+1} \dots A_n\}}_{\text{Completely unexplained}}$$

Lets a take step back

$$G_a = \mathbb{E}[m(R, Y) \mid A = a] - \mathbb{E}[m(R, Y)]$$

Lets a take step back

Performance of the subgroup of interest,
e.g. Age or younger patient

$$G_a = \overbrace{\mathbb{E}[m(R, Y) \mid A = a]} - \mathbb{E}[m(R, Y)]$$

Lets a take step back

Performance of the subgroup of interest,
e.g. Age or younger patient

Avg Performance of the model

$$G_a = \overbrace{\mathbb{E}[m(R, Y) \mid A = a]} - \overbrace{\mathbb{E}[m(R, Y)]}$$

Lets a take step back

Performance of the subgroup of interest,
e.g. Age or younger patient

Avg Performance of the model

$$G_a = \mathbb{E}[m(R, Y) \mid A = a] - \mathbb{E}[m(R, Y)]$$

Metric (accuracy, recall)

Prediction from model

Ground-truth label

Lets a take step back

Performance of the subgroup of interest,
e.g. Age or younger patient

Avg Performance of the model

$$G_a = \mathbb{E}[m(R, Y) \mid A = a] - \mathbb{E}[m(R, Y)]$$

Metric (accuracy, recall)

Prediction from model

Ground-truth label

m

X: Covariates or images

Y: Label

A: Subgroups

Lets a take step back

Performance of the subgroup of interest,
e.g. Age or younger patient

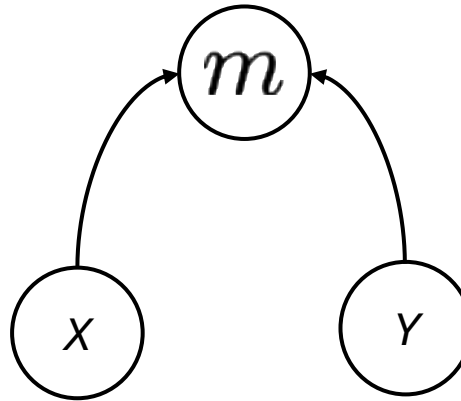
Avg Performance of the model

$$G_a = \mathbb{E}[m(R, Y) \mid A = a] - \mathbb{E}[m(R, Y)]$$

Metric (accuracy, recall)

Prediction from model

Ground-truth label



X: Covariates or images

Y: Label

A: Subgroups

Lets a take step back

Performance of the subgroup of interest,
e.g. Age or younger patient

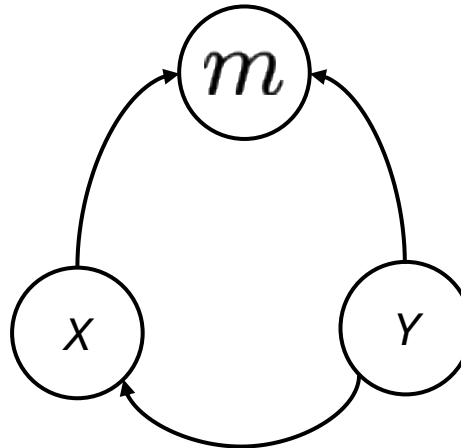
Avg Performance of the model

$$G_a = \mathbb{E}[m(R, Y) \mid A = a] - \mathbb{E}[m(R, Y)]$$

Metric (accuracy, recall)

Prediction from model

Ground-truth label



X: Covariates or images

Y: Label

A: Subgroups

Lets a take step back

Performance of the subgroup of interest,
e.g. Age or younger patient

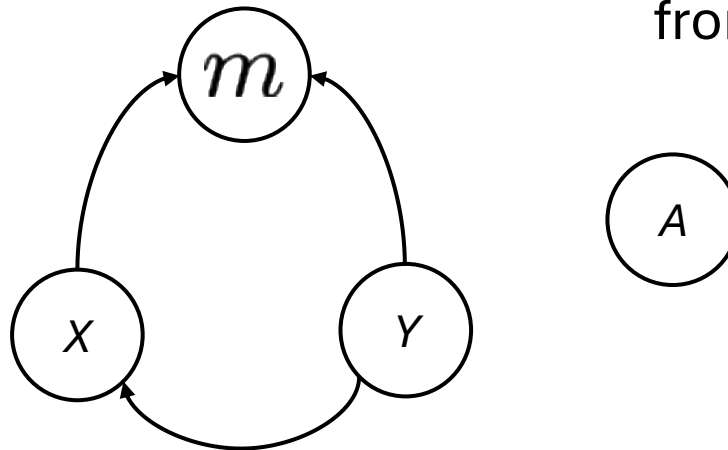
Avg Performance of the model

$$G_a = \mathbb{E}[m(R, Y) \mid A = a] - \mathbb{E}[m(R, Y)]$$

Metric (accuracy, recall)

Prediction from model

Ground-truth label



X : Covariates or images

Y : Label

A : Subgroups

Lets a take step back

Performance of the subgroup of interest,
e.g. Age or younger patient

Avg Performance of the model

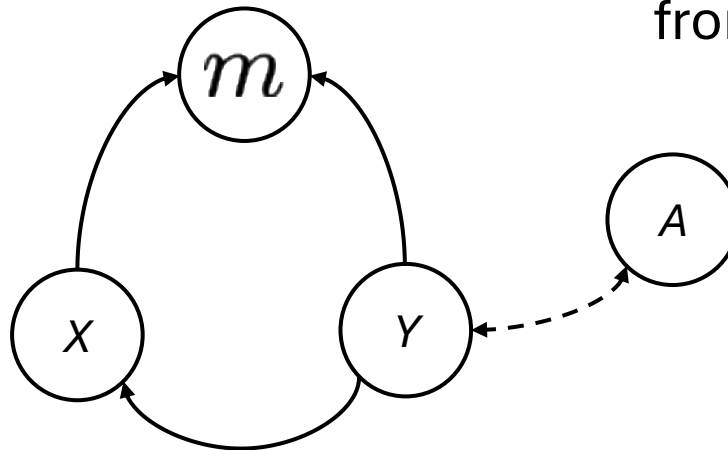
$$G_a = \mathbb{E}[m(R, Y) \mid A = a] - \mathbb{E}[m(R, Y)]$$

Metric (accuracy, recall)

Prediction from model

Ground-truth label

Label Shift



X : Covariates or images

Y : Label

A : Subgroups

←-----→ Unobserved confounding

————→ Causal direction

Lets a take step back

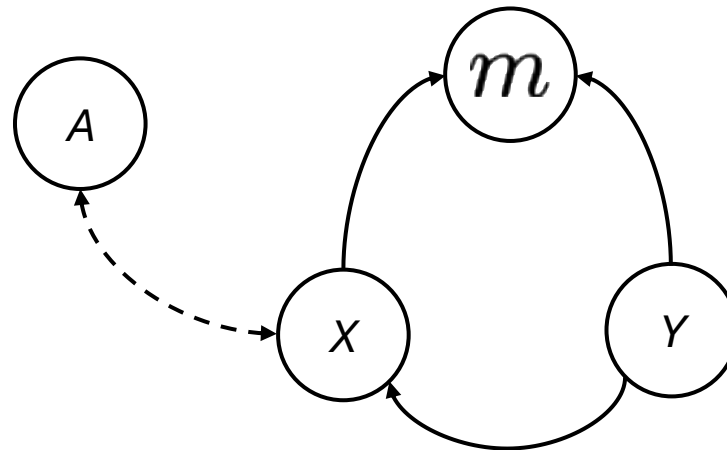
Performance of the subgroup of interest,
e.g. Age or younger patient

Avg Performance of the model

$$G_a = \overbrace{\mathbb{E}[m(R, Y) \mid A = a]}^{\text{Performance of the subgroup of interest, e.g. Age or younger patient}} - \overbrace{\mathbb{E}[m(R, Y)]}^{\text{Avg Performance of the model}}$$

Metric (accuracy, recall) Prediction from model Ground-truth label

Presentation
Shift



X: Covariates or images

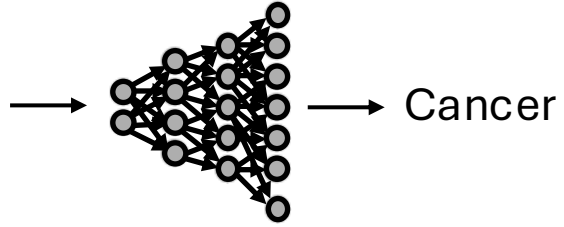
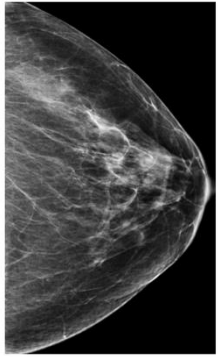
Y: Label

A: Subgroups

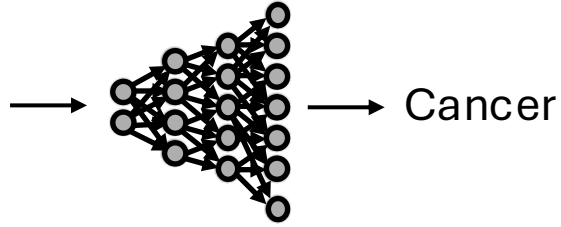
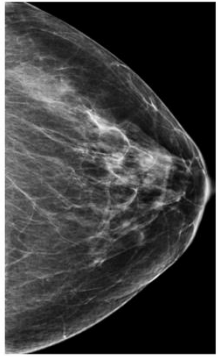
←-----→ Unobserved confounding

————→ Causal direction

A toy example



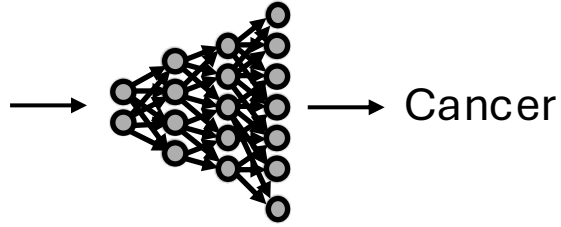
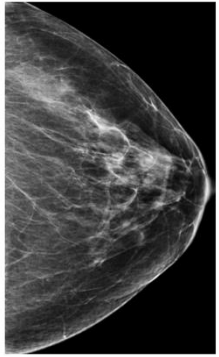
A toy example



Population (on cancer patients)

Avg Model error = 11 %

A toy example



Cancer

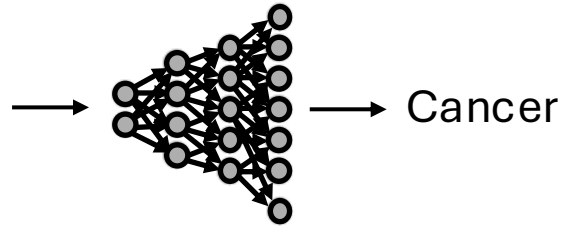
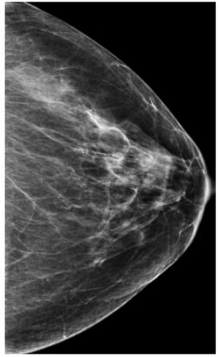
Population (on cancer patients)

Avg Model error = 11 %

Subgroup (on younger patients)

Model error = 29%

A toy example



Cancer

Population (on cancer patients)

Avg Model error = 11 %

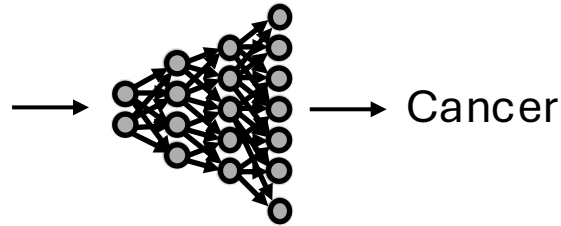
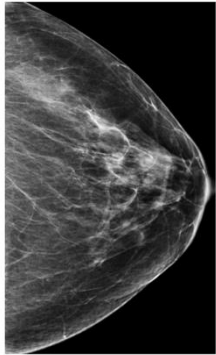
Subgroup (on younger patients)

Model error = 29%

Significant Gap

Conclusion: The model is struggling
on younger patients

A toy example



Cancer

Population (on cancer patients)

Avg Model error = **11 %**

Subgroup (on younger patients)

Model error = **29%**

Significant Gap

nature medicine

Article <https://doi.org/10.1038/s41591-024-03113-4>

The limits of fair medical imaging AI in real-world generalization

Received: 8 December 2023
Accepted: 5 June 2024
Published online: 28 June 2024
Check for updates

Yuzhe Yang^{1,2,3}, Haoran Zhang^{1,2}, Judy W. Gichoya², Dina Katabi¹ & Marzyeh Ghassemi^{1,3}

As artificial intelligence (AI) rapidly approaches human-level performance in medical imaging, it is crucial that it does not exacerbate or propagate healthcare disparities. Previous research established AI's capacity to

ARTICLES

<https://doi.org/10.1038/s41591-021-01595-0>

nature
medicine

Check for updates

OPEN

Underdiagnosis bias of artificial intelligence algorithms applied to chest radiographs in under-served patient populations

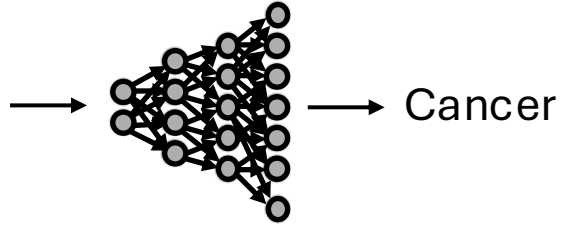
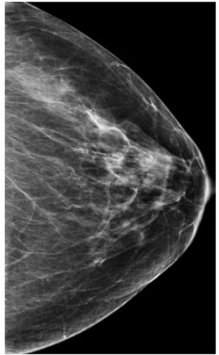
Laleh Seyyed-Kalantari^{1,2,3}, Haoran Zhang³, Matthew B. A. McDermott², Irene Y. Chen³ and Marzyeh Ghassemi^{1,3}

Artificial intelligence (AI) systems have increasingly achieved expert-level performance in medical imaging applications. However, there is growing concern that such AI systems may reflect and amplify human bias, and reduce the quality of their performance in historically under-served populations such as female patients, Black patients, or patients of low socioeconomic status. Such biases are especially troubling in the context of underdiagnosis, wherein the AI algorithm would incorrectly

Conclusion: The model is struggling on younger patients

Action: Balanced training, reweighting **w.r.t age** (GroupDRO, JTT, DFR)

A toy example



Cancer

Error

5%

35%

Fatty

Dense

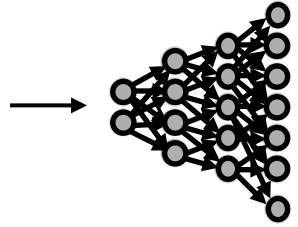
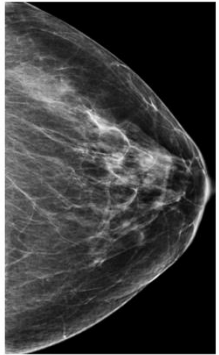
0.80

0.20

Population

$$\begin{aligned}\text{Avg Model error} &= (0.8 * 5) + (0.2 * 35) \\ &= \mathbf{11\%}\end{aligned}$$

A toy example



Cancer

Population

Error	5%	35%
	Fatty 0.80	Dense 0.20

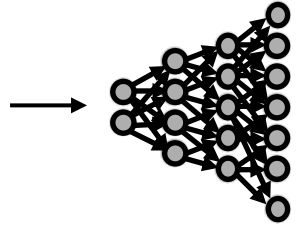
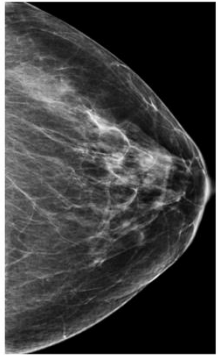
$$\begin{aligned}\text{Avg Model error} &= (0.8 * 5) + (0.2 * 35) \\ &= \mathbf{11\%}\end{aligned}$$

Young

Error	5%	35%
	Fatty 0.20	Dense 0.80

$$\begin{aligned}\text{Model error} &= (0.2 * 5) + (0.8 * 35) \\ &= \mathbf{29\%}\end{aligned}$$

A toy example



Cancer

Population

Error	5%	35%
	Fatty 0.80	Dense 0.20

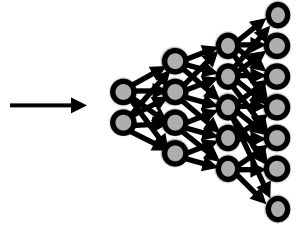
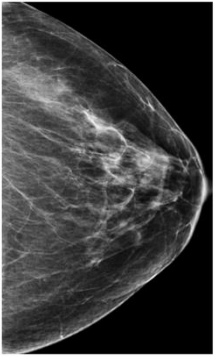
$$\text{Avg Model error} = (0.8 * 5) + (0.2 * 35) = 11\%$$

Young		
Error	5%	35%
	Fatty 0.20	Dense 0.80

$$\text{Model error} = (0.2 * 5) + (0.8 * 35) = 29\%$$

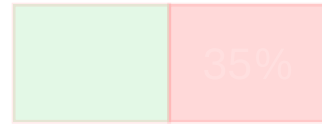
Significant Gap

A toy example



Cancer

Error

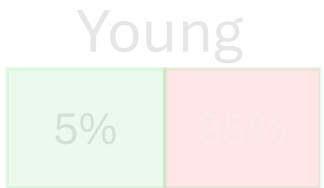


Fatty 0.80
Dense 0.20

Population

$$\text{Avg Model error} = (0.8 * 5) + (0.2 * 35) = 11\%$$

Error

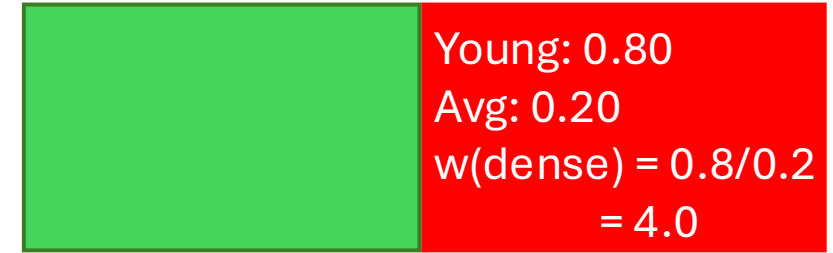


Fatty 0.20
Dense 0.80

Young

$$\text{Model error} = (0.2 * 5) + (0.8 * 35) = 29\%$$

Density

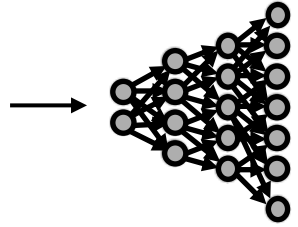
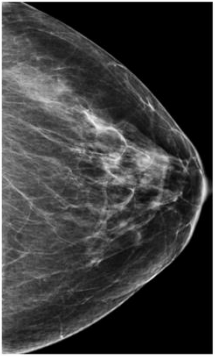


Fatty

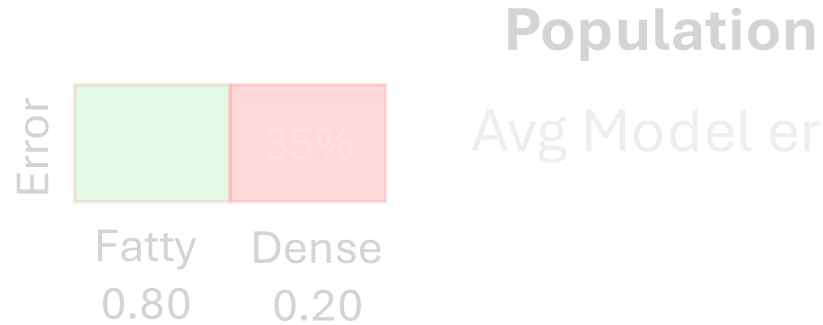
Dense

Young: 0.80
Avg: 0.20
 $w(\text{dense}) = 0.8/0.2 = 4.0$

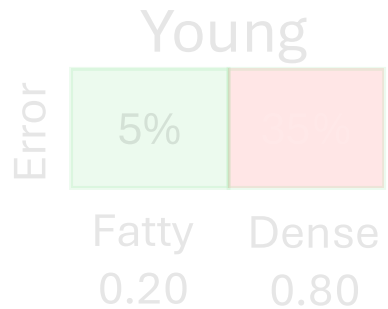
A toy example



Cancer



$$\text{Avg Model error} = (0.8 * 5) + (0.2 * 35) = 11\%$$



$$\text{Model error} = (0.2 * 5) + (0.8 * 35) = 29\%$$

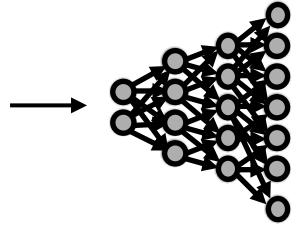
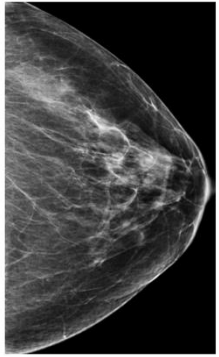
Density

Young: 0.20	Young: 0.80
Avg: 0.80	Avg: 0.20
w(fatty) = 0.2/0.8 = 0.25	w(dense) = 0.8/0.2 = 4.0
Fatty	Dense

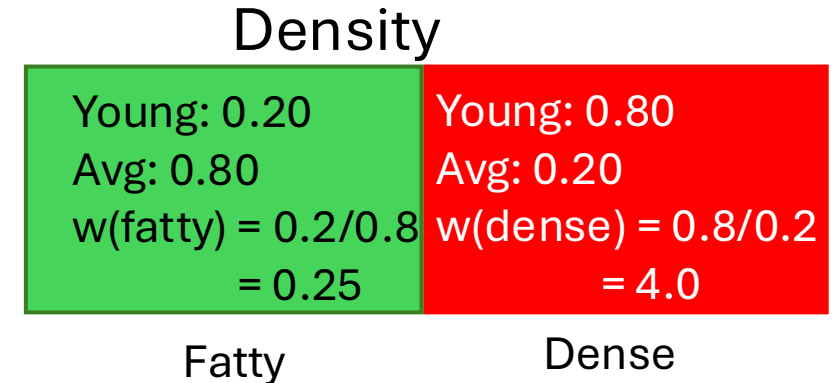
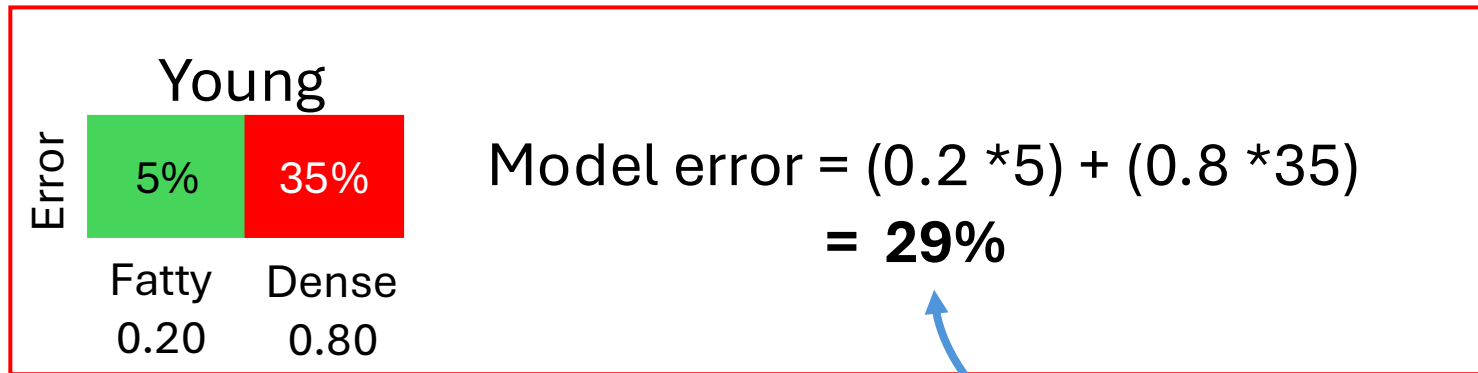
Model error Stratified by density

$$= (0.25 * 0.8 * 5 + 4.0 * 0.2 * 35) = 29\%$$

A toy example



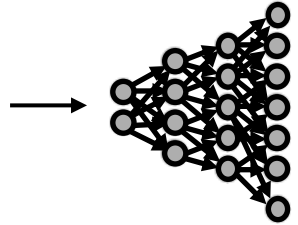
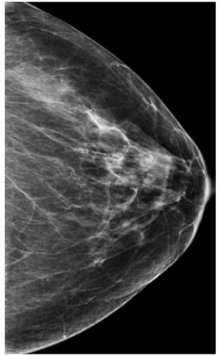
Cancer



Model error Stratified by density
= **29%**

Matches

A toy example



Cancer

Population

Error	5%	35%
	Fatty 0.80	Dense 0.20

$$\text{Avg Model error} = (0.8 * 5) + (0.2 * 35) = 11\%$$

Error	Young	
	5%	35%
	Fatty 0.20	Dense 0.80

$$\text{Model error} = (0.2 * 5) + (0.8 * 35) = 29\%$$

Density

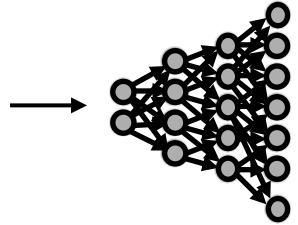
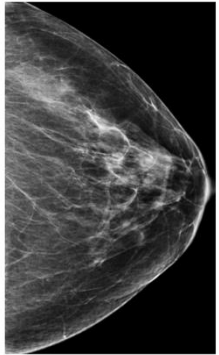
Young: 0.20 Avg: 0.80 $w(\text{fatty}) = 0.2/0.8 = 0.25$	Young: 0.80 Avg: 0.20 $w(\text{dense}) = 0.8/0.2 = 4.0$
--	---

Fatty
Down weight easy
one

Dense
up weight hard
one

$$\text{Model error Stratified by density} = 29\%$$

A toy example



Cancer

Population

Error	5%	35%
	Fatty	Dense
	0.80	0.20

$$\text{Avg Model error} = (0.8 * 5) + (0.2 * 35) = 11\%$$

Pseudo young Population

Density (V)

Young: 0.20	Young: 0.80
Avg: 0.80	Avg: 0.20
w(fatty) = 0.2/0.8 = 0.25	w(dense) = 0.8/0.2 = 4.0

Fatty
Down weight easy
one

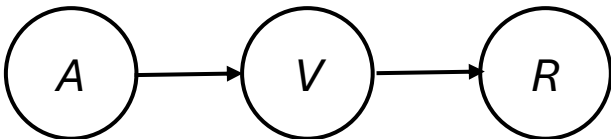
Dense
up weight hard
one

$$\text{Model error Stratified by density} = 29\%$$

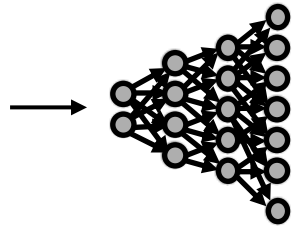
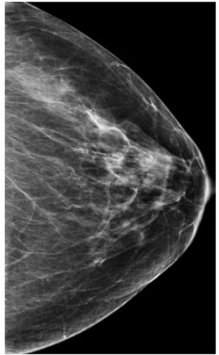
Young

Error	5%	35%
	Fatty	Dense
	0.20	0.80

$$\text{Model error} = (0.2 * 5) + (0.8 * 35) = 29\%$$



A toy example



Cancer

Population

Error	5%	35%
	Fatty	Dense
	0.80	0.20

$$\text{Avg Model error} = (0.8 * 5) + (0.2 * 35) = 11\%$$

Pseudo young Population

Density (V)

Young: 0.20	Young: 0.80
Avg: 0.80	Avg: 0.20
w(fatty) = 0.2/0.8 = 0.25	w(dense) = 0.8/0.2 = 4.0

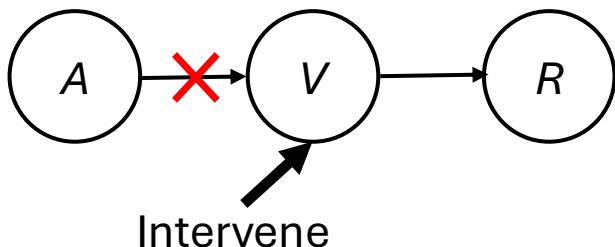
Fatty
Down weight easy one

Dense
up weight hard one

$$\text{Model error Stratified by density} = 29\%$$

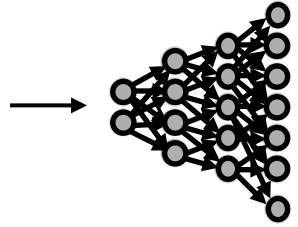
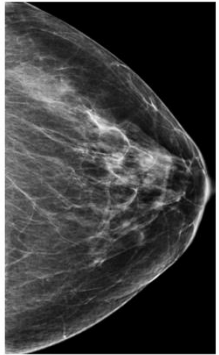
	Young
Error	5% 35%
	Fatty Dense
	0.20 0.80

$$\text{Model error} = (0.2 * 5) + (0.8 * 35) = 29\%$$



What would happen if population is same as the subgroup?

A toy example



Cancer

Population

Error	5%	35%
	Fatty 0.80	Dense 0.20

$$\text{Avg Model error} = (0.8 * 5) + (0.2 * 35) = 11\%$$

Young

Error	5%	35%
	Fatty 0.20	Dense 0.80

$$\text{Model error} = (0.2 * 5) + (0.8 * 35) = 29\%$$

Pseudo young Population

Density (V)

Young: 0.20 Avg: 0.80 w(fatty) = 0.2/0.8 = 0.25	Young: 0.80 Avg: 0.20 w(dense) = 0.8/0.2 = 4.0
--	---

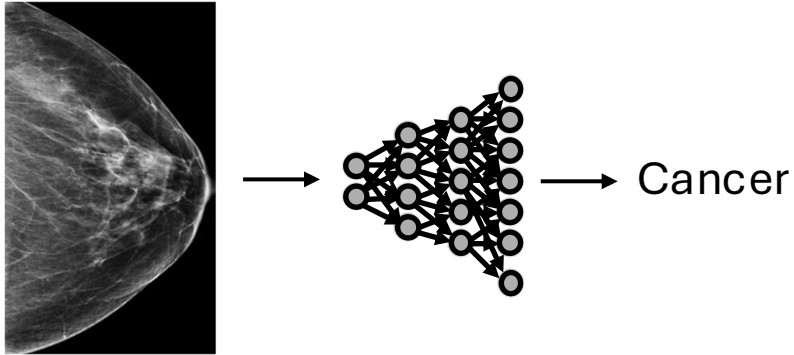
Fatty
Down weight easy
one

Dense
up weight hard
one

$$\text{Model error Stratified by density} = 29\%$$

Conclusion: The model is struggling
on younger patients **dense breasts**

A toy example



$$\mathbb{E}[m(R, Y) \mid A = a]$$

$$M_a(V) = \mathbb{E}[w_a(V)m(R, Y)]$$

Young		Error	Model error = $(0.2 * 5) + (0.8 * 35)$ = 29%
5%	35%		
Fatty 0.20	Dense 0.80		

Density (V)	
Young: 0.20 Avg: 0.80 $w(\text{fatty}) = 0.2/0.8 = 0.25$	Young: 0.80 Avg: 0.20 $w(\text{dense}) = 0.8/0.2 = 4.0$
Fatty Down weight easy one	Dense up weight hard one
Model error Stratified by density = 29%	

$$T_a(V) = \mathbb{E}[m(R, Y) \mid A = a] - \underbrace{\mathbb{E}[w_a(V)m(R, Y)]}_{M_a(V)}$$

If $T_a \approx 0 \Rightarrow V$ explains the performance gap
 $\Rightarrow (R, Y) \perp A \mid V$

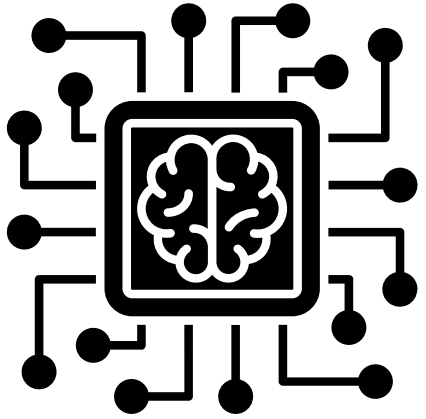
So, out of many hypotheses proposed by Ladder,
 we reject them by choosing appropriate V 's and doing this conditional
 Independence test.

$$T_a(V) = \mathbb{E}[m(R, Y) \mid A = a] - \underbrace{\mathbb{E}[w_a(V)m(R, Y)]}_{M_a(V)}$$

If $T_a \approx 0 \Rightarrow V$ explains the performance gap
 $\Rightarrow (R, Y) \perp A \mid V$

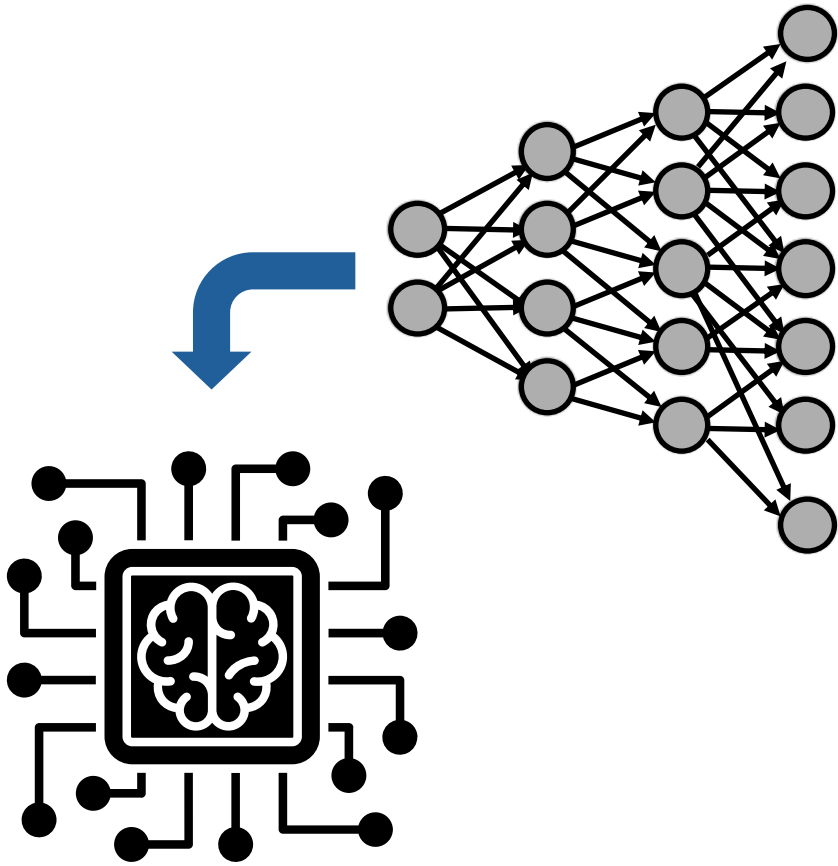
How to find the V ?

Proposal: Where we are going?



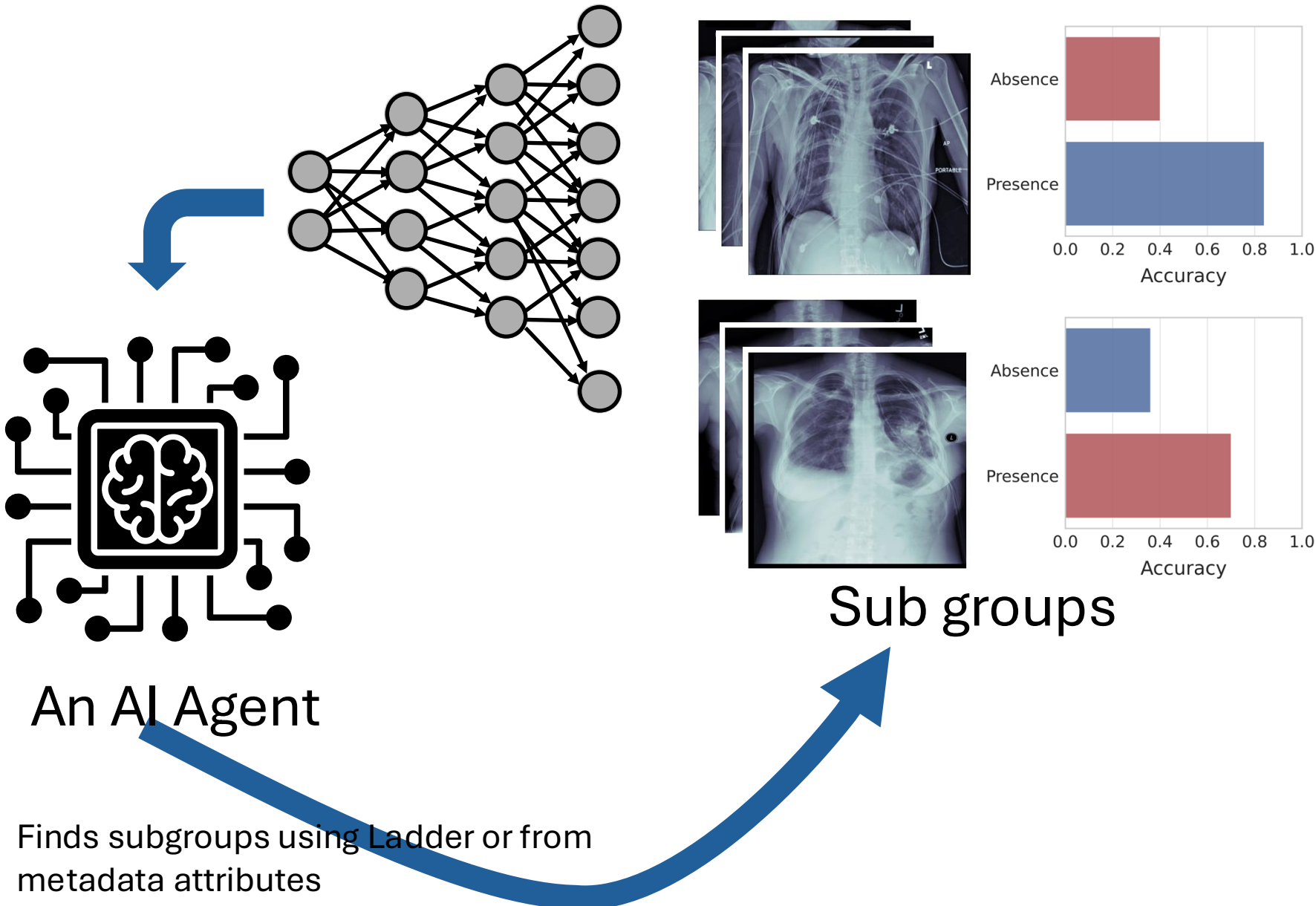
An AI Agent

Proposal: Where we are going?

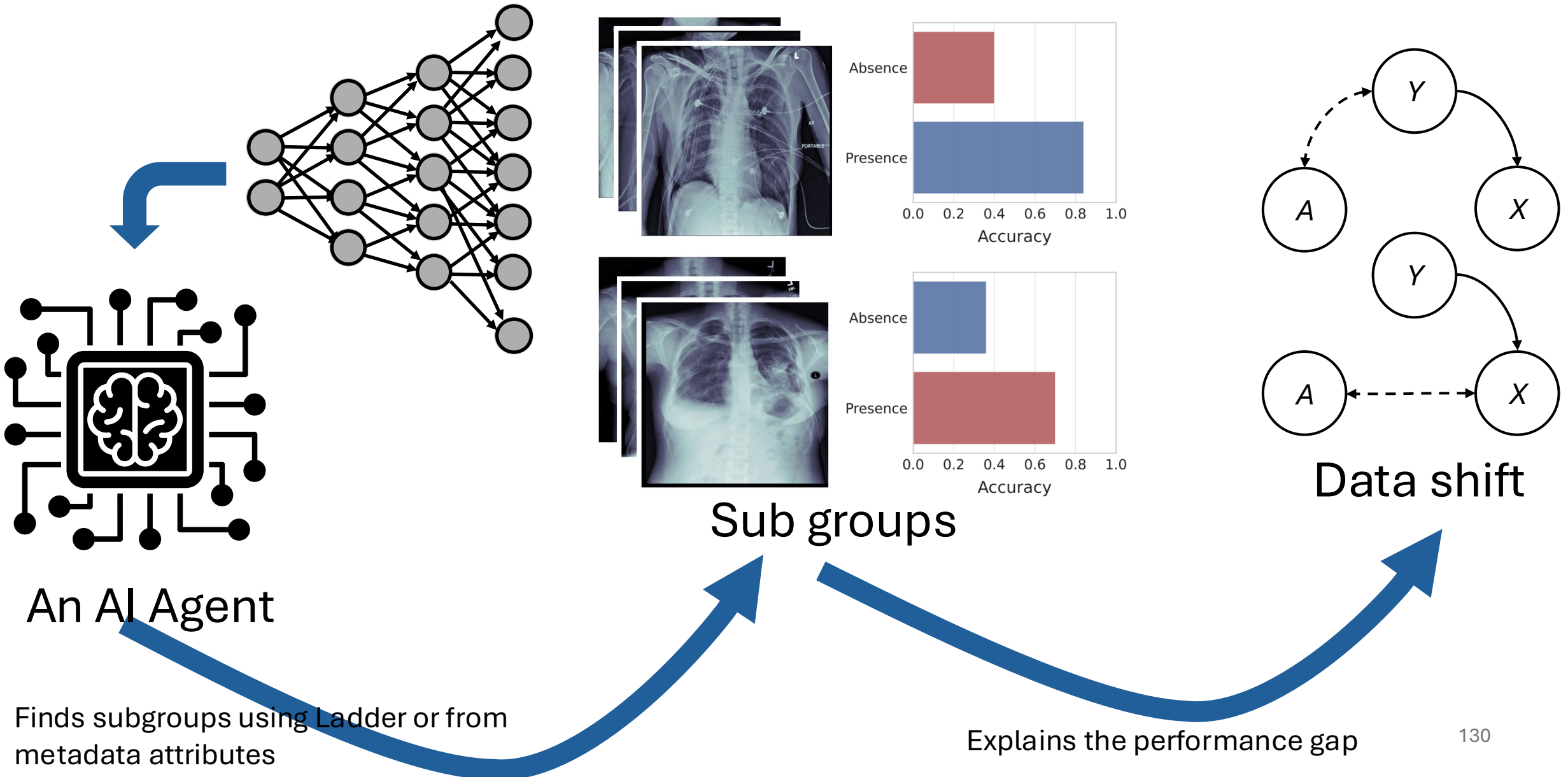


An AI Agent

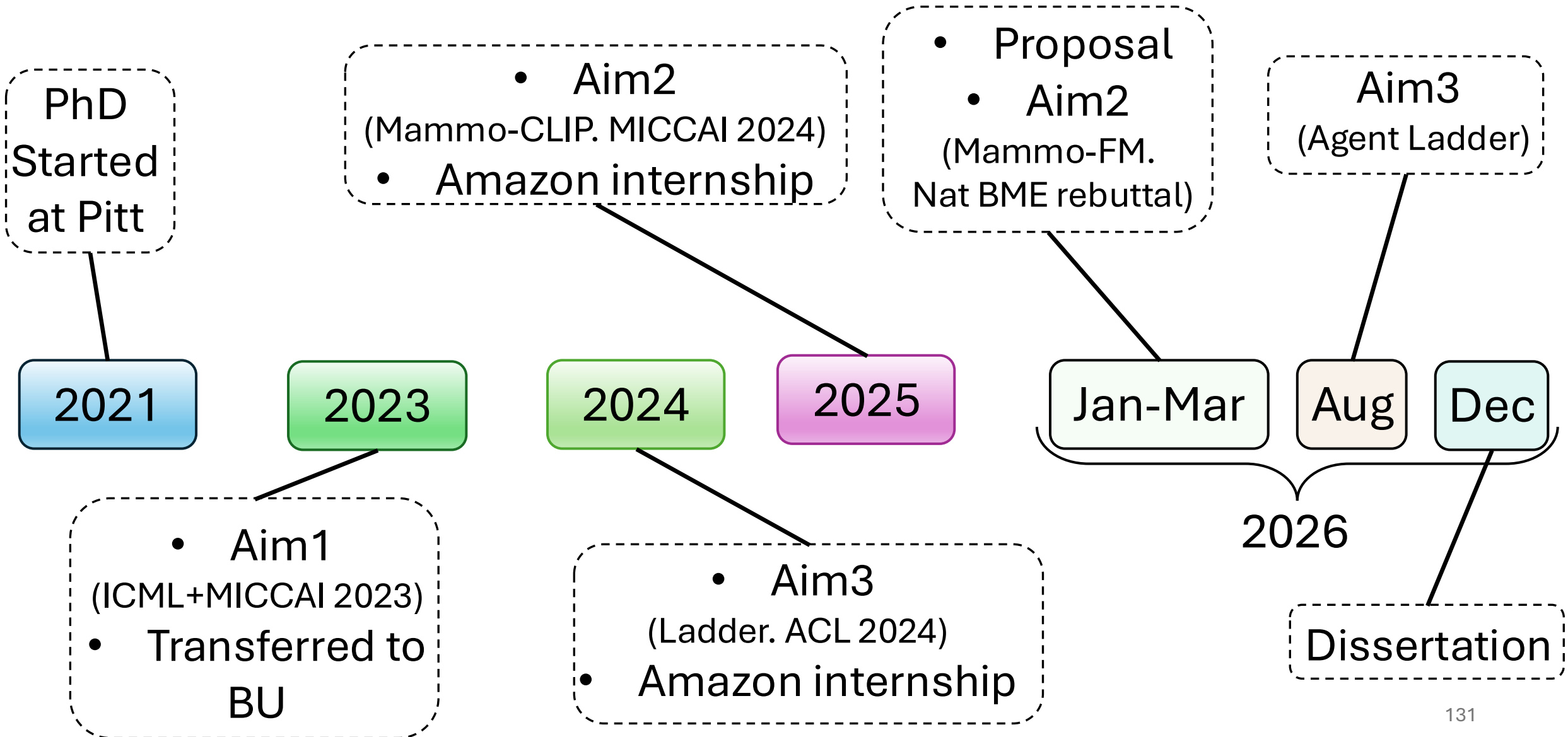
Proposal: Where we are going?



Proposal: Where we are going?



Timeline

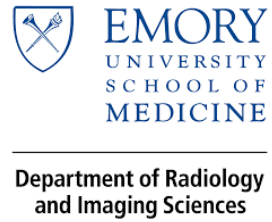


Papers

- Anatomy-Guided Weakly-Supervised Abnormality Localization in Chest X-rays. Ke Yu, [Shantanu Ghosh](#), Zhexiong Liu, Christopher Deible, Kayhan Batmanghelich. **MICCAI 2022**
- Anatomy-specific Progression Classification in Chest Radiographs via Weakly Supervised Learning. Ke Yu, [Shantanu Ghosh](#), Zhexiong Liu, Christopher Deible, Clare B. Poynton, Kayhan Batmanghelich. **RAD: AI**
- Dividing and Conquering a BlackBox to a Mixture of Interpretable Models: Route, Interpret, Repeat. [Shantanu Ghosh](#), Ke Yu, Forough Arabshahi, Kayhan Batmanghelich. **ICML 2023**
- Tackling Shortcut Learning in Deep Neural Networks: An Iterative Approach with Interpretable Models. [Shantanu Ghosh](#), Ke Yu, Forough Arabshahi, Kayhan Batmanghelich. **SCIS (w)@ICML 2023**
- Distilling BlackBox to Interpretable models for Efficient Transfer Learning. [Shantanu Ghosh](#), Ke Yu, Kayhan Batmanghelich. **MICCAI 2023 (Top 14%)**
- Bridging the Gap: From Post Hoc Explanations to Inherently Interpretable Models for Medical Imaging. [Shantanu Ghosh](#), Ke Yu, Forough Arabshahi, Kayhan Batmanghelich. **IMLH (w)@ICML 2023**
- Mammo-CLIP: A Vision Language Foundation Model to Enhance Data Efficiency and Robustness in Mammography. [Shantanu Ghosh](#), Clare B. Poynton, Shyam Visweswaran, Kayhan Batmanghelich. **MICCAI 2024 (Top 11%)**
- Mammo-FM: Breast-specific foundational model for Integrated Mammographic Diagnosis, Prognosis, and Reporting. [Shantanu Ghosh](#), Vedant Parthesh Joshi, Rayan Syed, Aya Kassem, Abhishek Varshney, Payel Basak, Weicheng Dai, Judy Wawira Gichoya, Hari M. Trivedi, Imon Banerjee, Shyam Visweswaran, Clare B. Poynton, Kayhan Batmanghelich. **ArXiv 2025**
- Distributionally robust self-supervised learning for tabular data. [Shantanu Ghosh](#), Tiankang Xie, Mikhail Kuznetsov. **TRL (w)@NeurIPS 2024**
- LADDER: Language-Driven Slice Discovery and Error Rectification in Vision Classifiers. [Shantanu Ghosh](#), Rayan Syed, Chenyu Wang, Vaibhav Choudhary, Binxu Li, Clare B. Poynton, Shyam Visweswaran, Kayhan Batmanghelich. **ACL 2025**
- Semantic Consistency-Based Uncertainty Quantification for Factuality in Radiology Report Generation. Chenyu Wang, Weichao Zhou, [Shantanu Ghosh](#), Kayhan Batmanghelich, Wenchao Li. **NAACL 2025**
- PhyDiCT: Training-Free 3D CT Reconstruction from Sparse X-Rays via Differentiable Rendering and Strong Priors. Weicheng Dai, [Shantanu Ghosh](#), Kayhan Batmanghelich. **CVPR 2025 submission**



Advisor



Acknowledgments



Collaborators



Thank you!

Questions?