

# Divide and Conquer: Carving Out Symbolic Models out of BlackBox for More Efficient Domain Adaptation



Shantanu Ghosh<sup>1</sup>, Ke Yu<sup>2</sup>, Kayhan Batmanghelich<sup>1</sup>

<sup>1</sup>BU ECE, <sup>2</sup>Pitt ISP



# Desiderata of Explainability in Healthcare

Trust

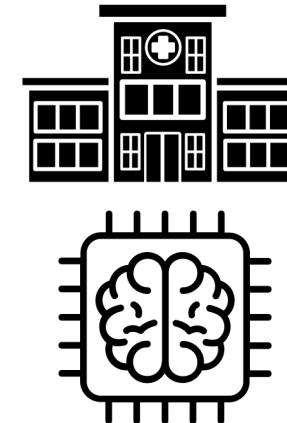
Causality

Transferability

Informativeness

Ethical Decision making

Training



Deployment



# Desiderata of Explainability in Healthcare

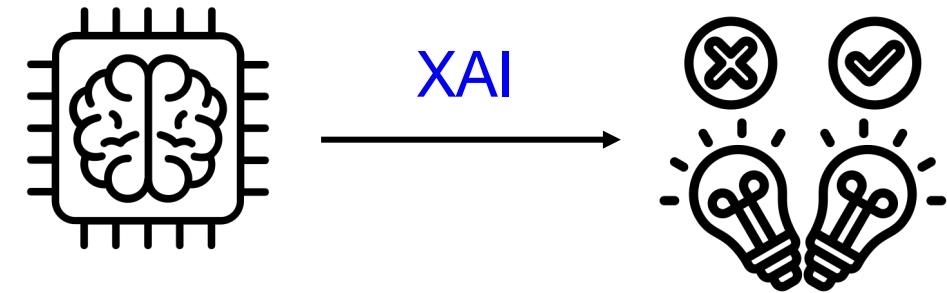
Trust

Causality

Transferability

Informativeness

Ethical Decision making



Obtaining genetics insights from  
deep learning via explainable artificial  
intelligence

Gherman Novakovsky <sup>1,2,7</sup>, Nick Dexter  <sup>3,4,7</sup>, Maxwell W. Libbrecht  <sup>4,8</sup>✉,  
Wyeth W. Wasserman  <sup>1,8</sup>✉ and Sara Mostafavi  <sup>5,6,8</sup>✉

# Desiderata of Explainability in Healthcare

Trust

Causality

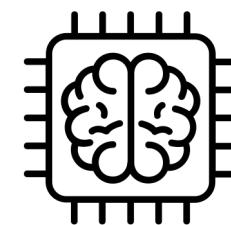
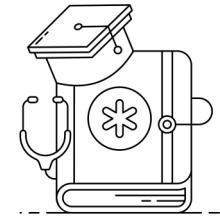
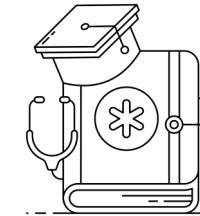
Transferability

Informativeness

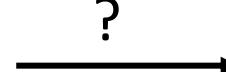
Ethical Decision making

Training

Deployment



?



# Desiderata of Explainability in Healthcare

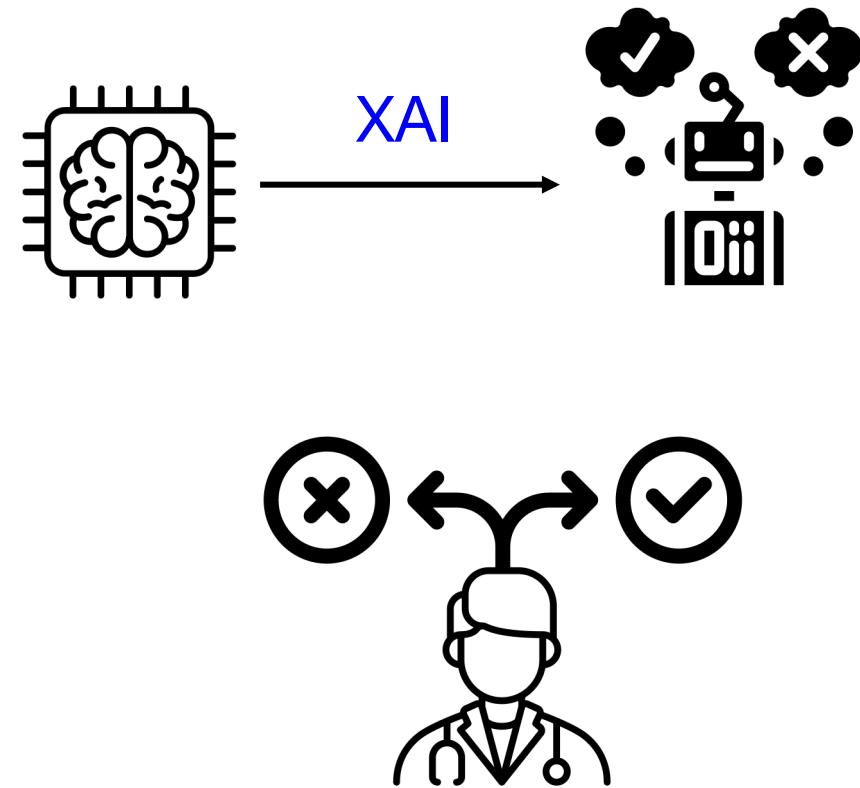
Trust

Causality

Transferability

Informativeness

Ethical Decision making



# Desiderata of Explainability in Healthcare

Trust

Causality

Transferability

Informativeness

Ethical Decision making

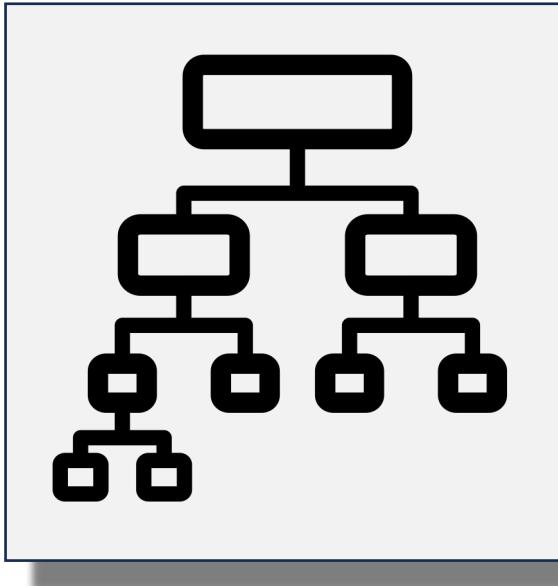
Training

Deployment

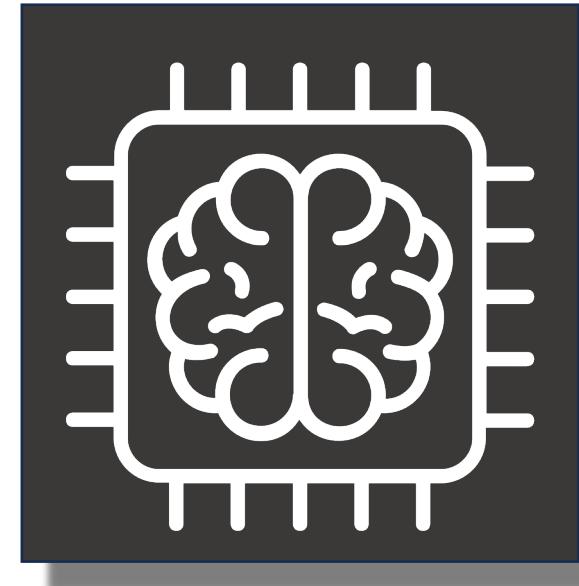


# General Design Choices

Transparent Models

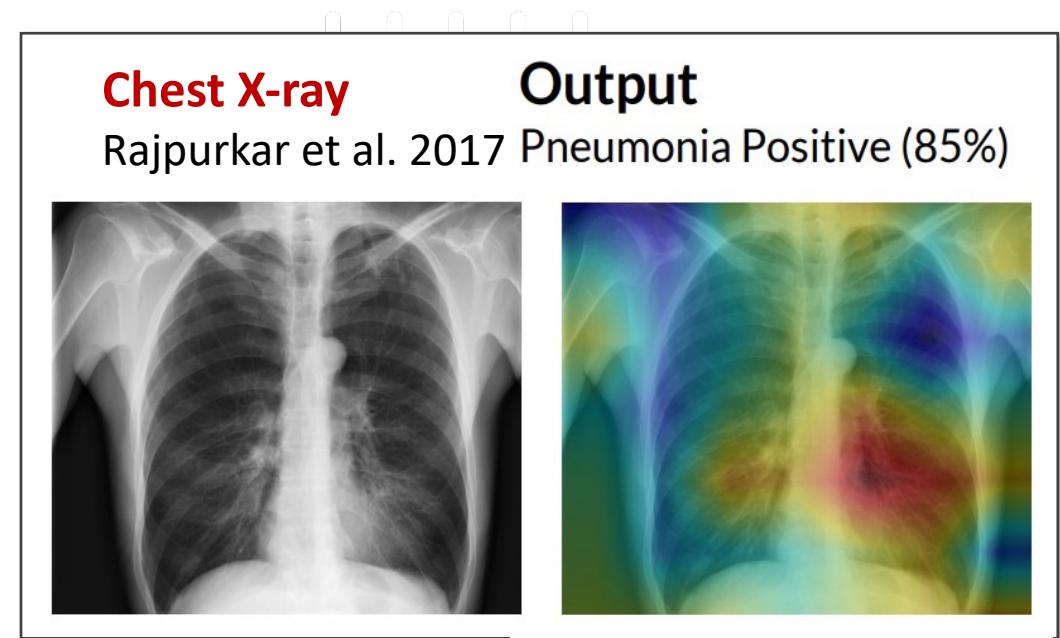


Post-hoc Explanation



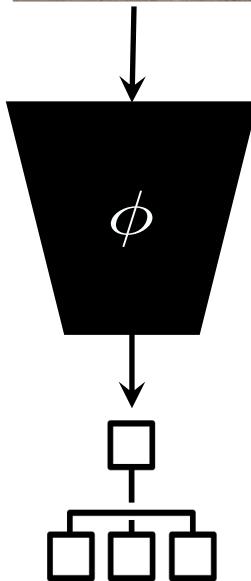
# General Design Choices

## Post-hoc Explanation

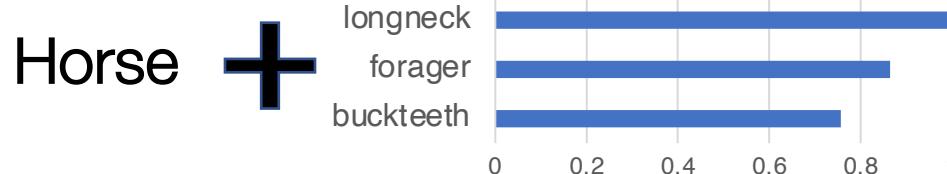


# General Design Choices

## Transparent Models



Top 3 concepts to identify Horse



## Post-hoc Explanation

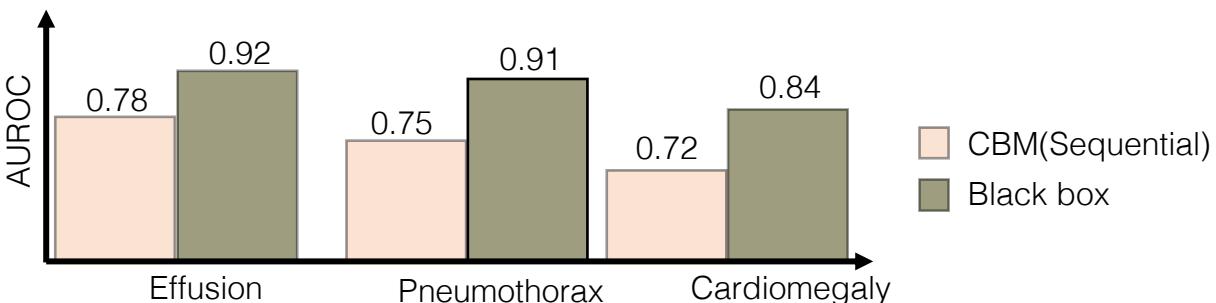
**Chest X-ray**  
Rajpurkar et al. 2017  $\triangleright$  pneumonia Positive (85%)

**Output**

# General Design Choices

## Transparent Models

Support concept intervention  
Limited design choices  
Difficult to get high-performance



## Post-hoc Explanation

Does not alter the Black box  
Leaves too much to human  
Inconsistent explanations  
No recourse



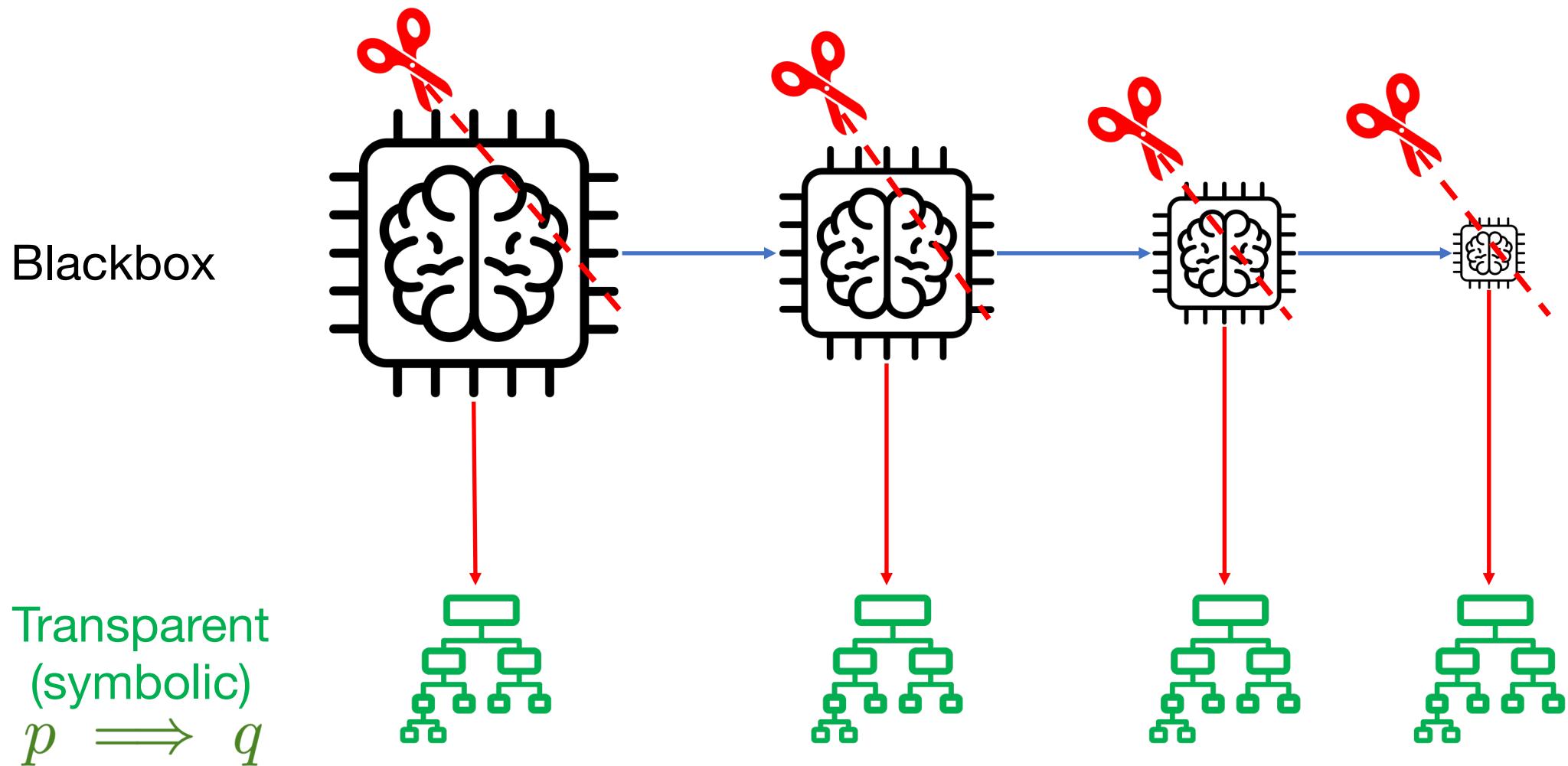
# Desiderata

- Not compromising performance
- Being able to intervene (fix) undesirable properties
- Transferability: Clinical rules are transferrable between domains
  - The transparent should look like clinical rules

$$p \implies q$$

- **R1:** fasting p-Glucose level > 126 mg/dL on two separate tests  $\implies$  may be diabetes.
- **R2:** 2-hour p-Glucose level during an Oral Glucose Tolerance Test > 200 mg/dL  $\implies$  diabetes (probabilistic).
- **R3:** random p-Glucose level > 200 mg/dL  $\wedge$  hyperglycemia (frequent urination  $\wedge$  increased thirst  $\wedge$  unexplained weight loss  $\implies$  diabetes (probabilistic).
- **R4:** A1C test > 6.5% (visit 1)  $\wedge$  A1C test > 6.5% (visit 2)  $\implies$  diabetes (probabilistic)

# General Idea

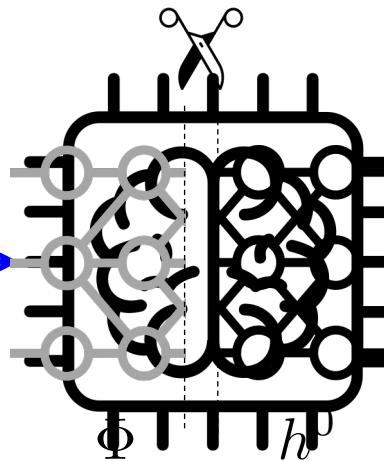


# Problem Set Up

$\chi$



$\gamma$

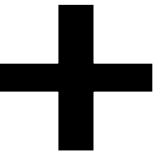


$c$

thin ears  
shortened muzzle  
round feet

....

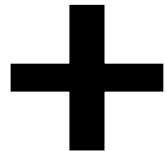
# Problem Set Up



## 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:

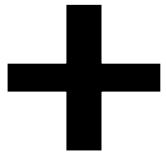
Right upper lobe consolidation with adjacent.  
While this may be infectious in nature, a CT scan is recommended for further clarification.

parse the reports to get the concepts

C

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

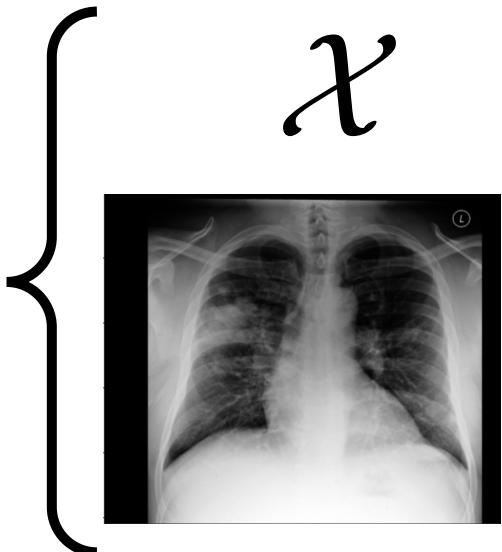
# Problem Set Up



## 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



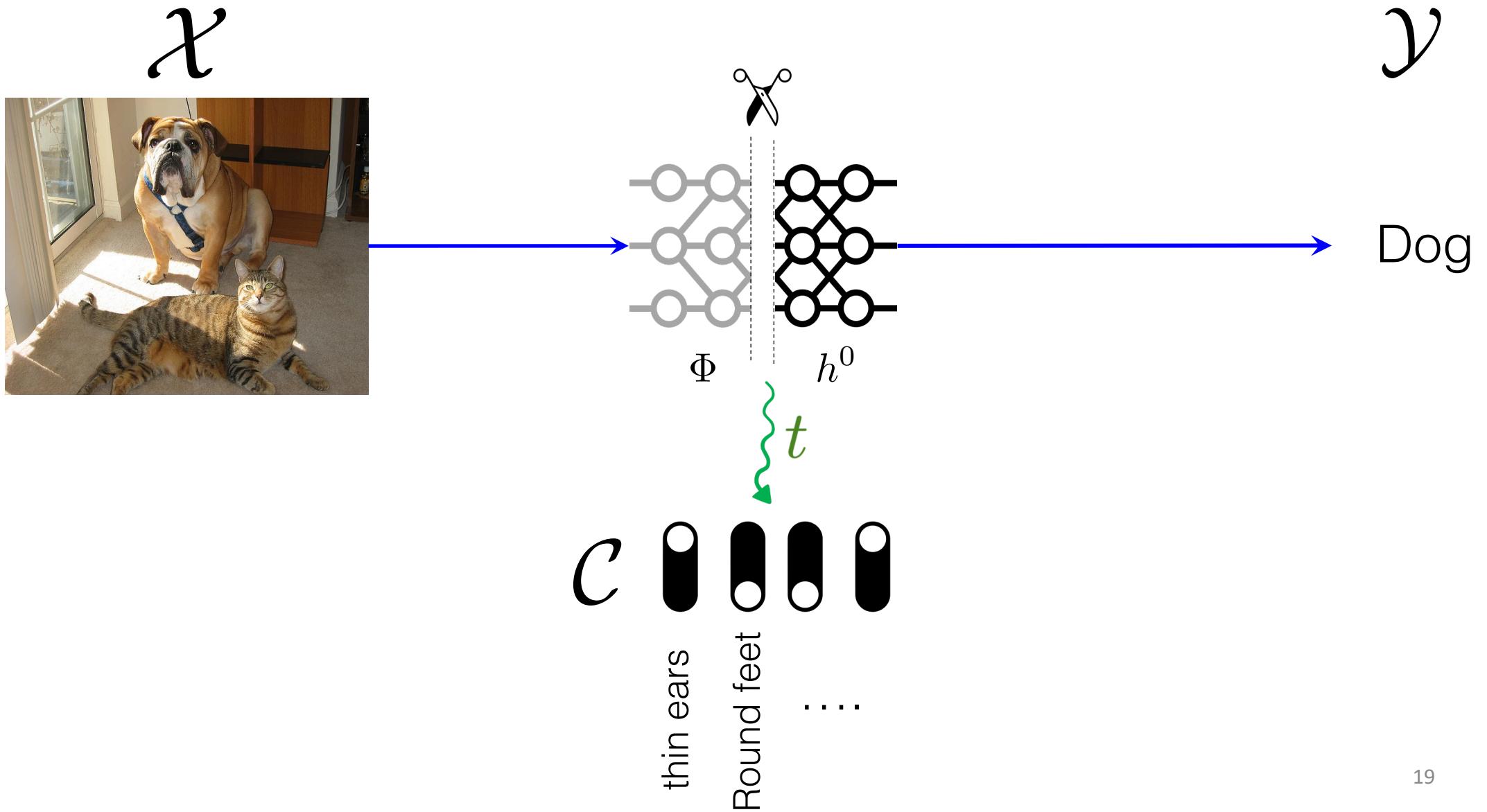
$C$

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

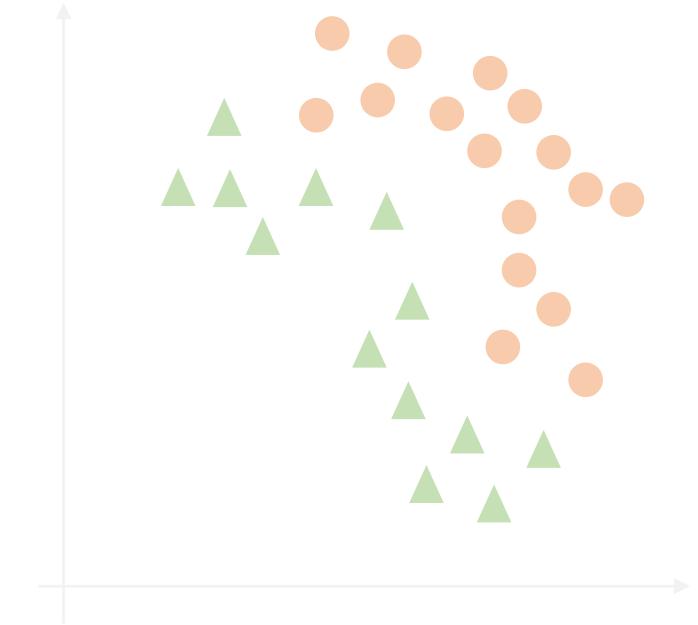
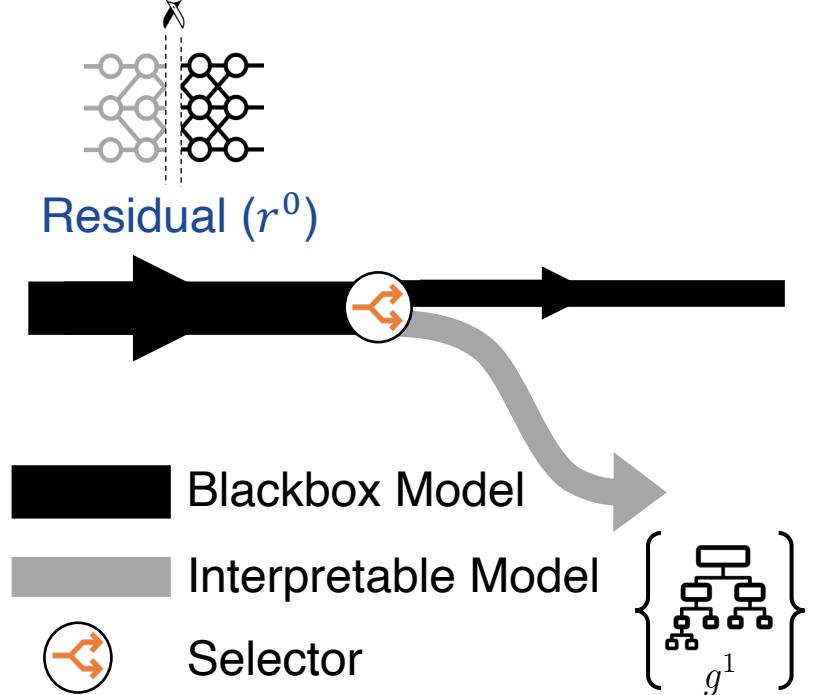
$\gamma$

Consolidation

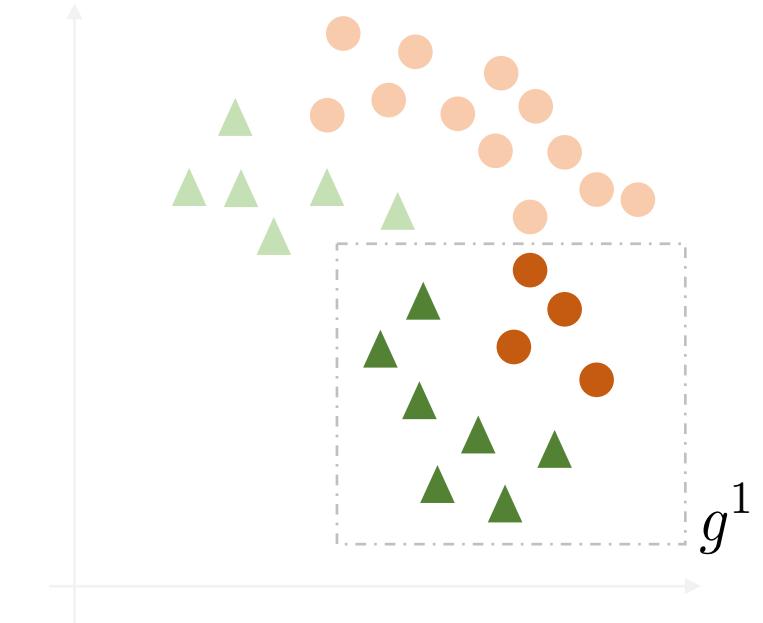
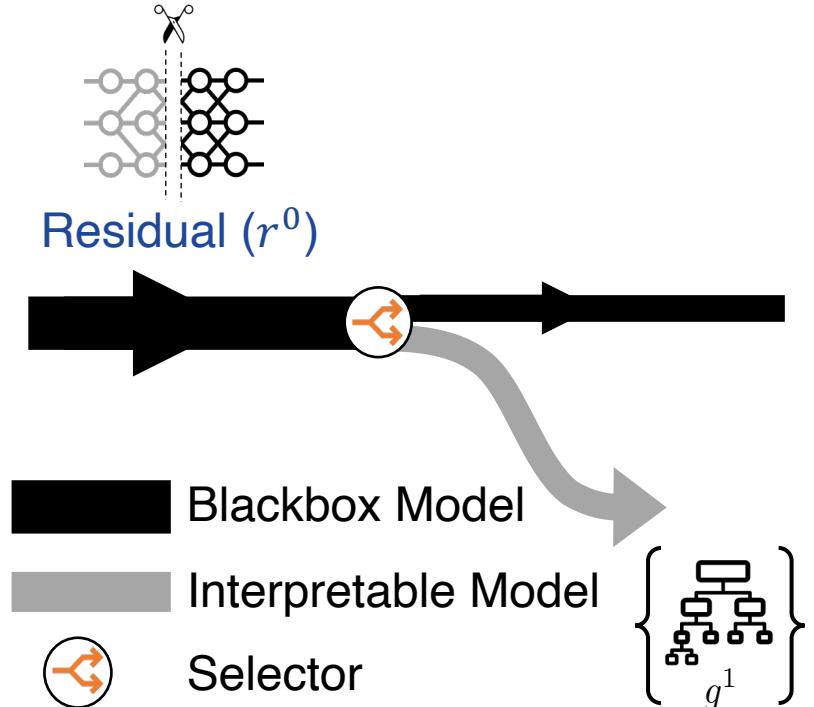
# Discovering Hidden Concepts



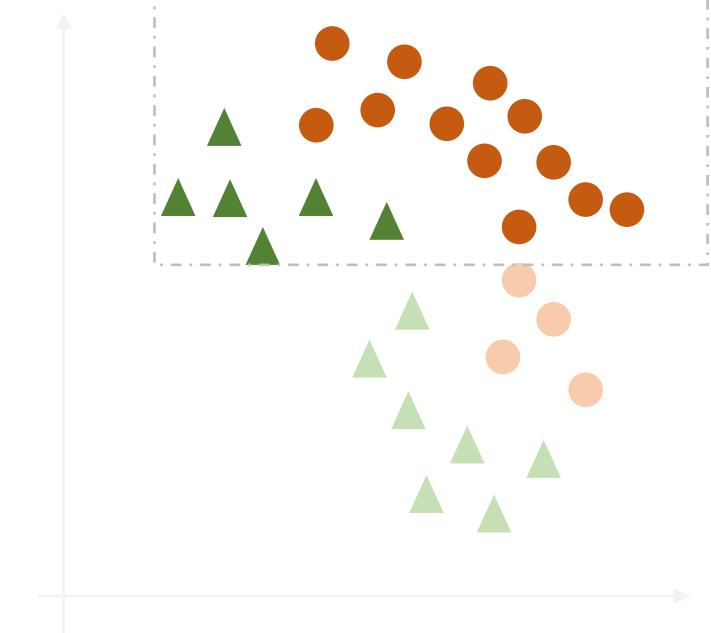
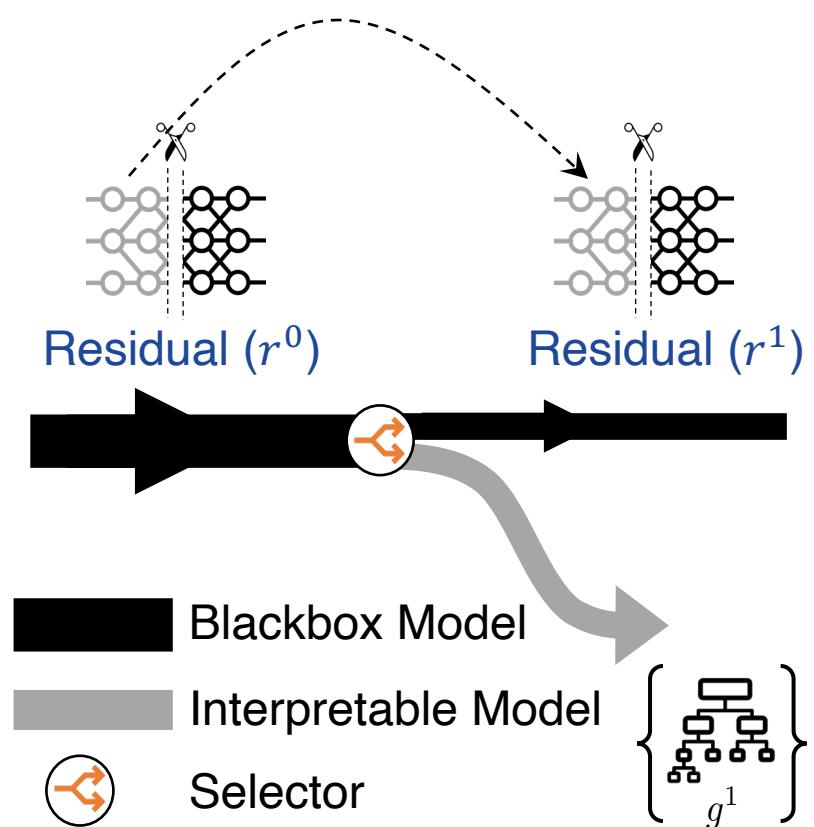
# Carving out Interpretable Models



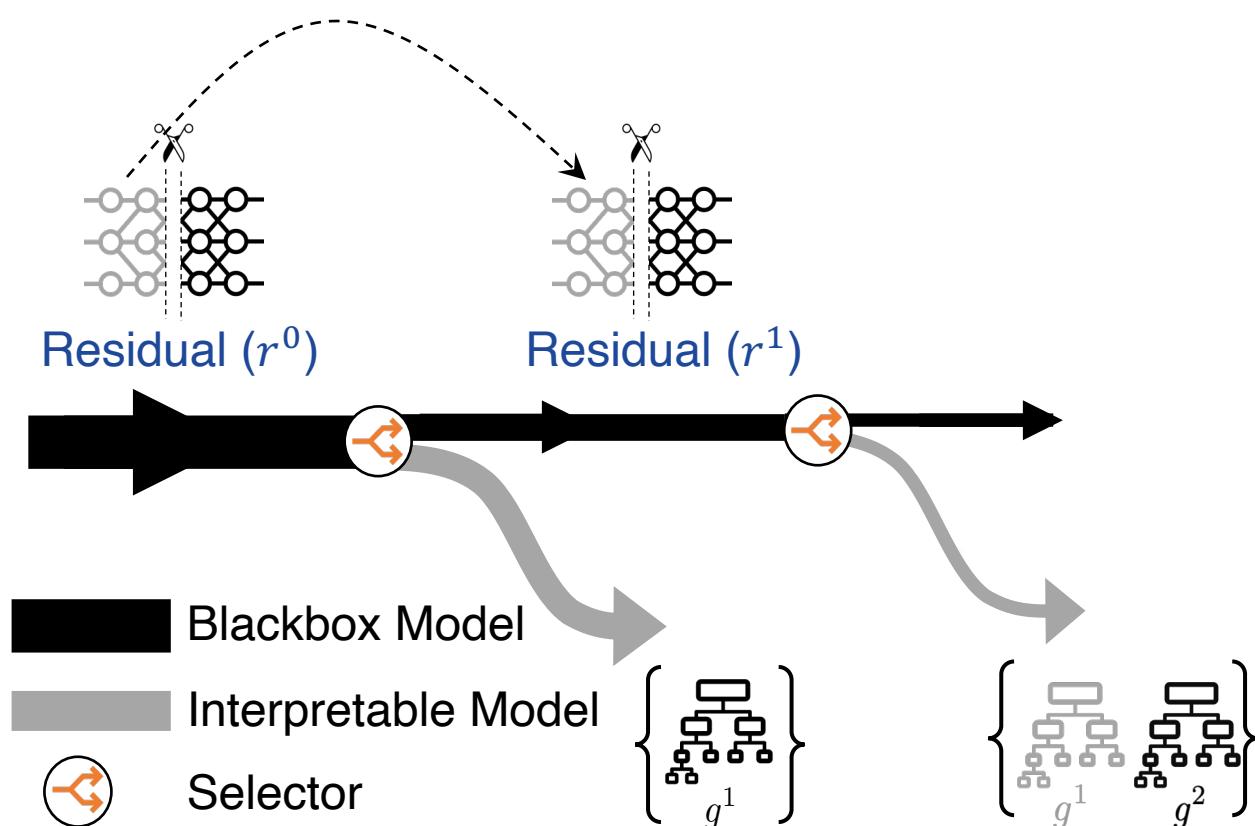
# Carving out Interpretable Models



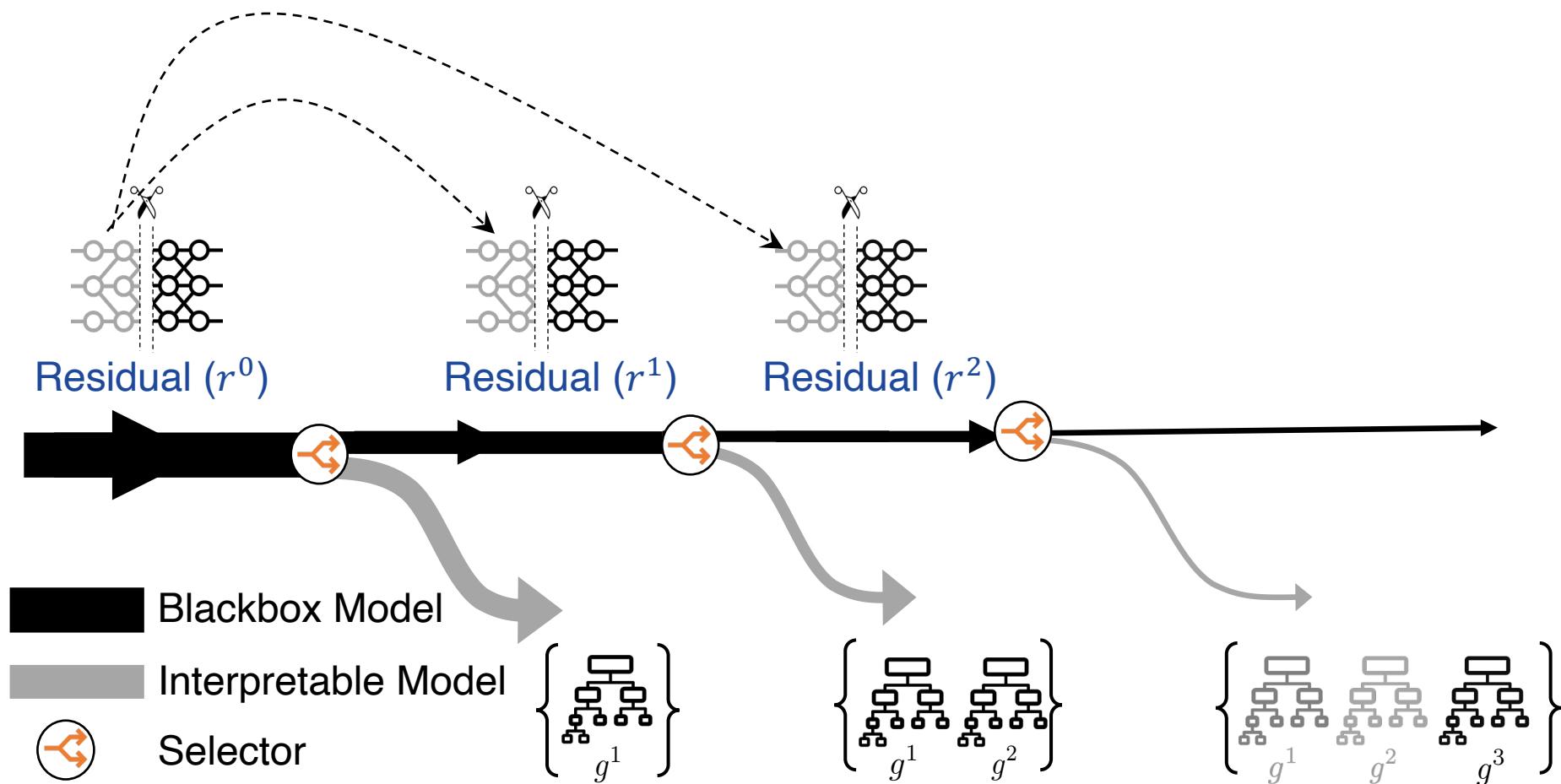
# Carving out Interpretable Models



# Carving out Interpretable Models



# Carving out Interpretable Models



# Examples on Bird Dataset



# Examples on Bird Dataset



Expert 1

Olive sided Flycatcher  $\leftrightarrow$  breast\_color\_grey  $\wedge$   
tail\_pattern\_solid



# Examples on Bird Dataset



Expert 1

Olive sided Flycatcher  $\leftrightarrow$  breast\_color\_grey  $\wedge$   
tail\_pattern\_solid



Expert 2

Olive sided Flycatcher  $\leftrightarrow$  underparts\_color\_grey  $\wedge$   
wing\_color\_grey

# Examples on Bird Dataset



Expert 1

Olive sided Flycatcher  $\leftrightarrow$  breast\_color\_grey  $\wedge$   
tail\_pattern\_solid



Expert 2

Olive sided Flycatcher  $\leftrightarrow$  underparts\_color\_grey  $\wedge$   
wing\_color\_grey

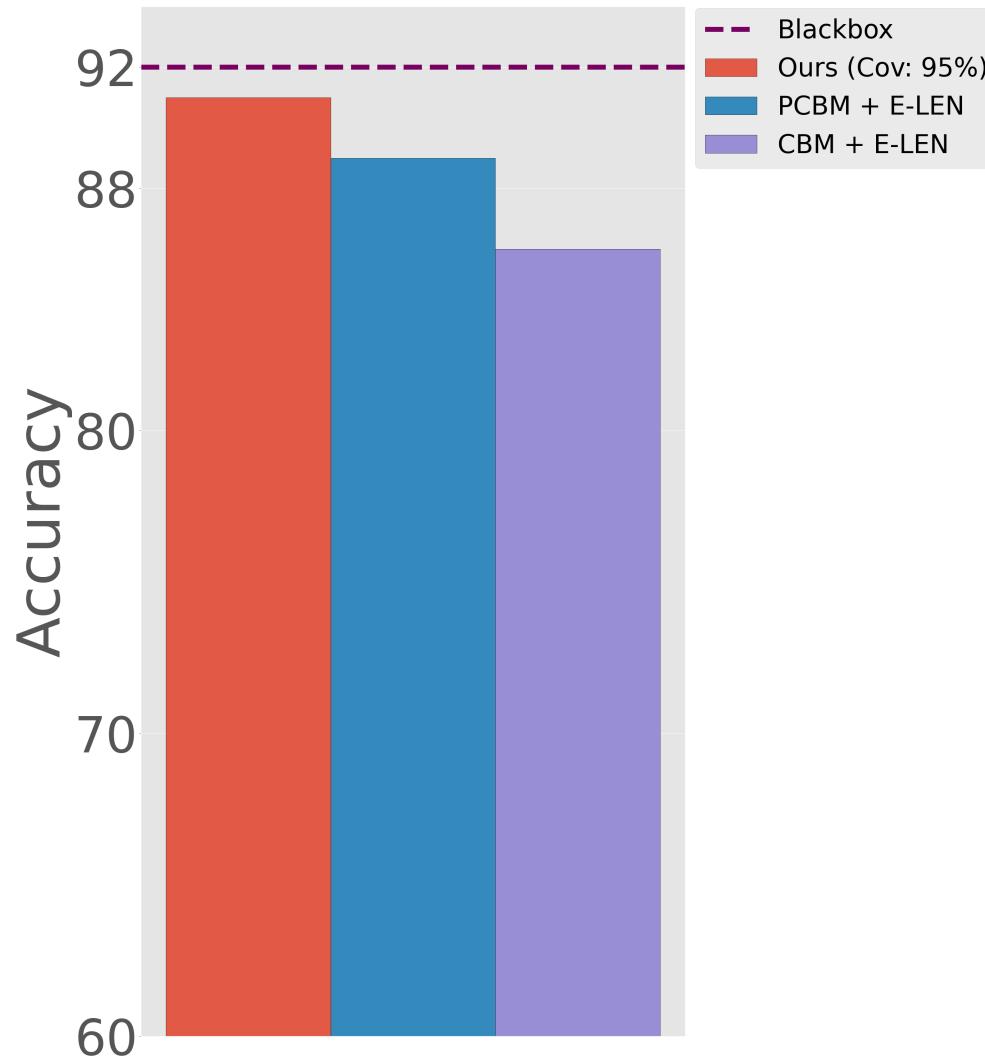


Final Residual  
(Unexplained)

# Comparing Performance

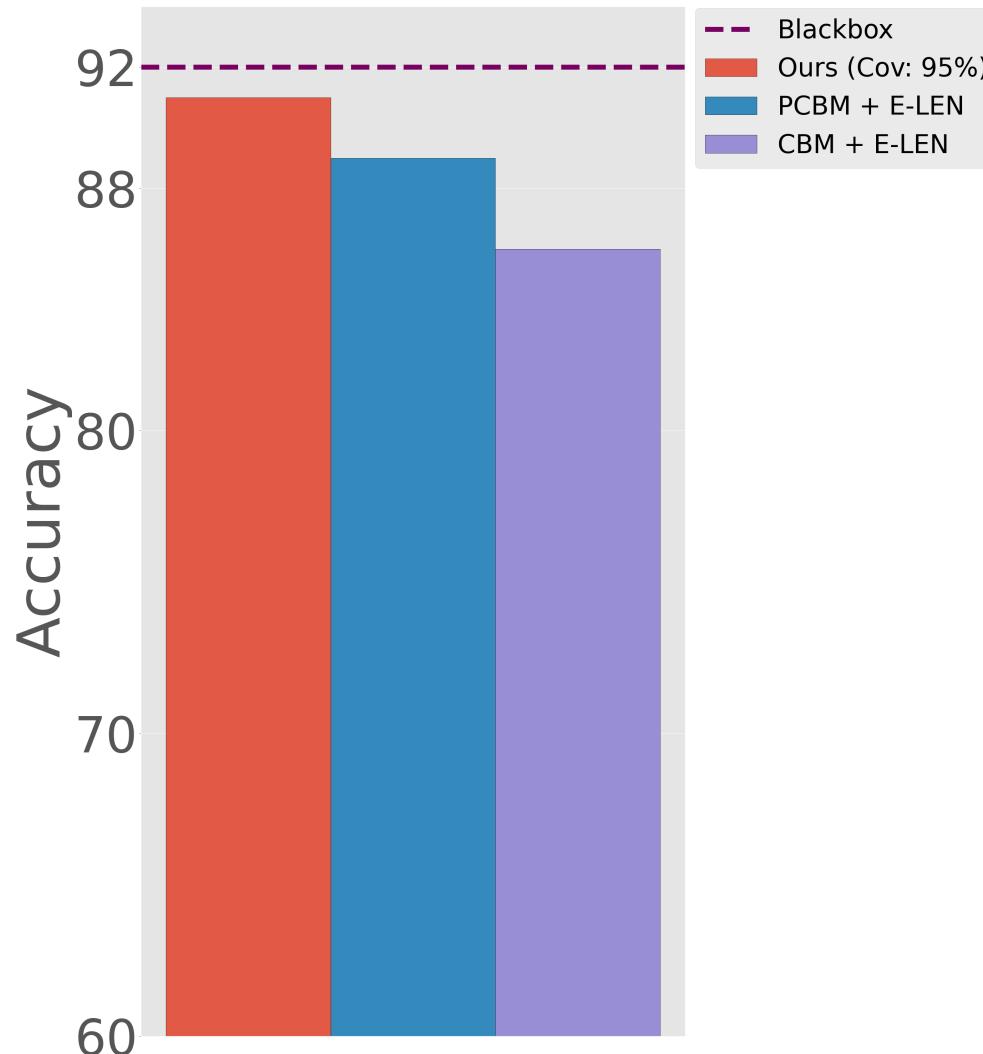
# Comparing Performance

## CUB-200 with ViT

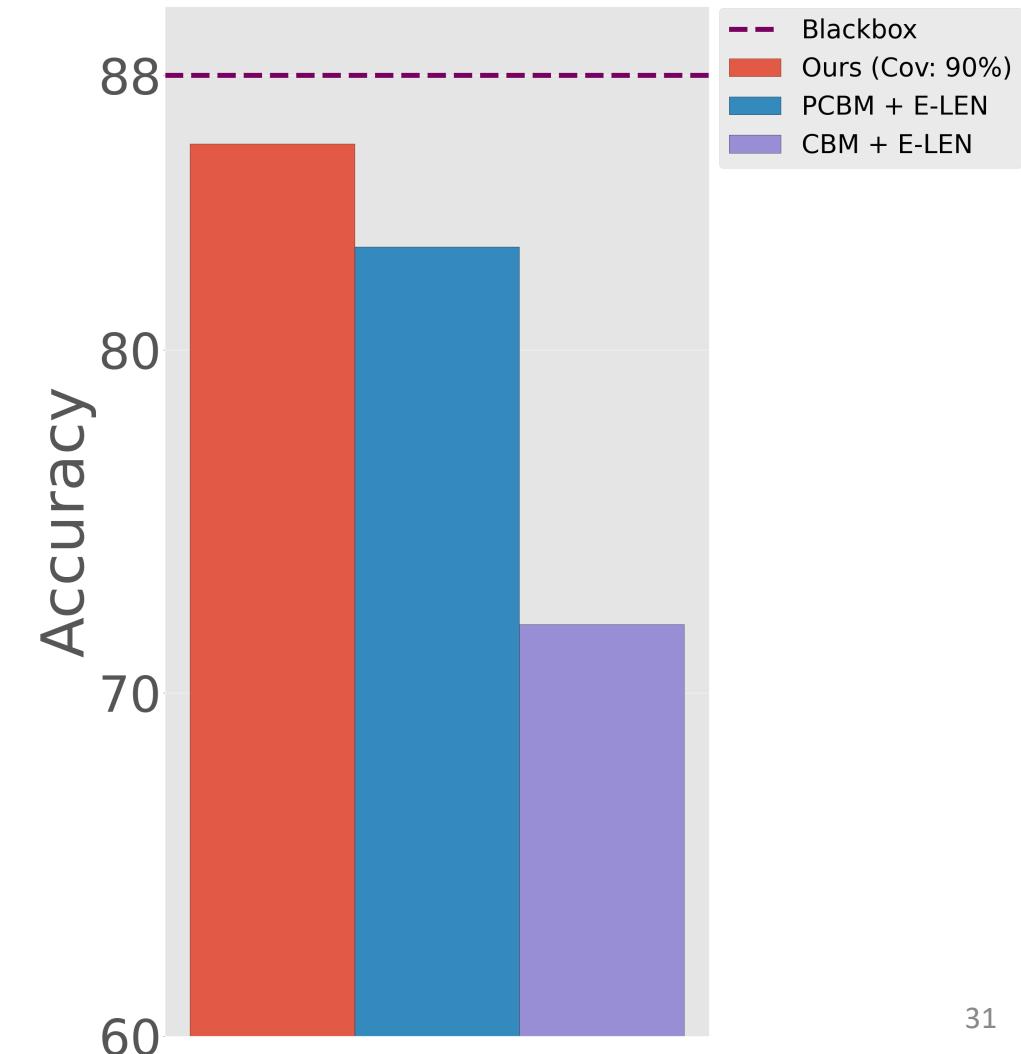


# Comparing Performance

CUB-200 with ViT

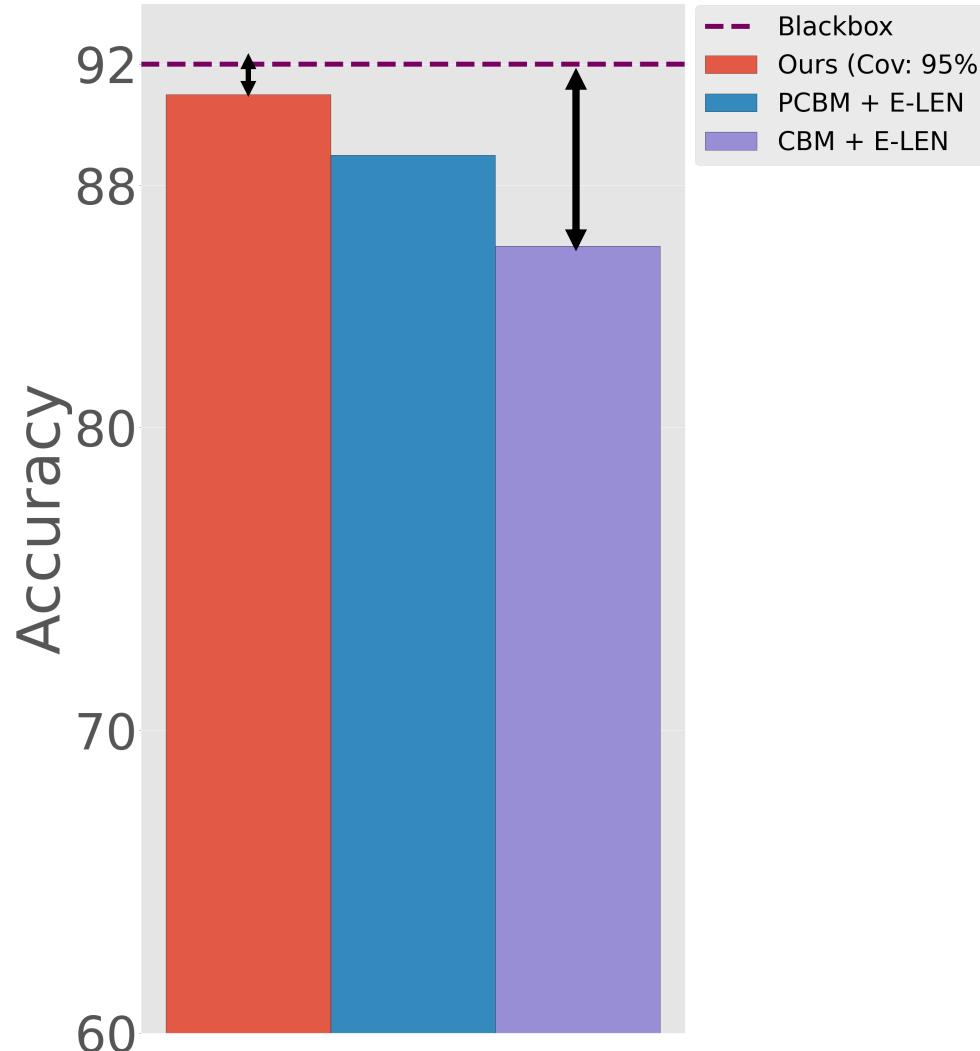


CUB-200 with ResNet101

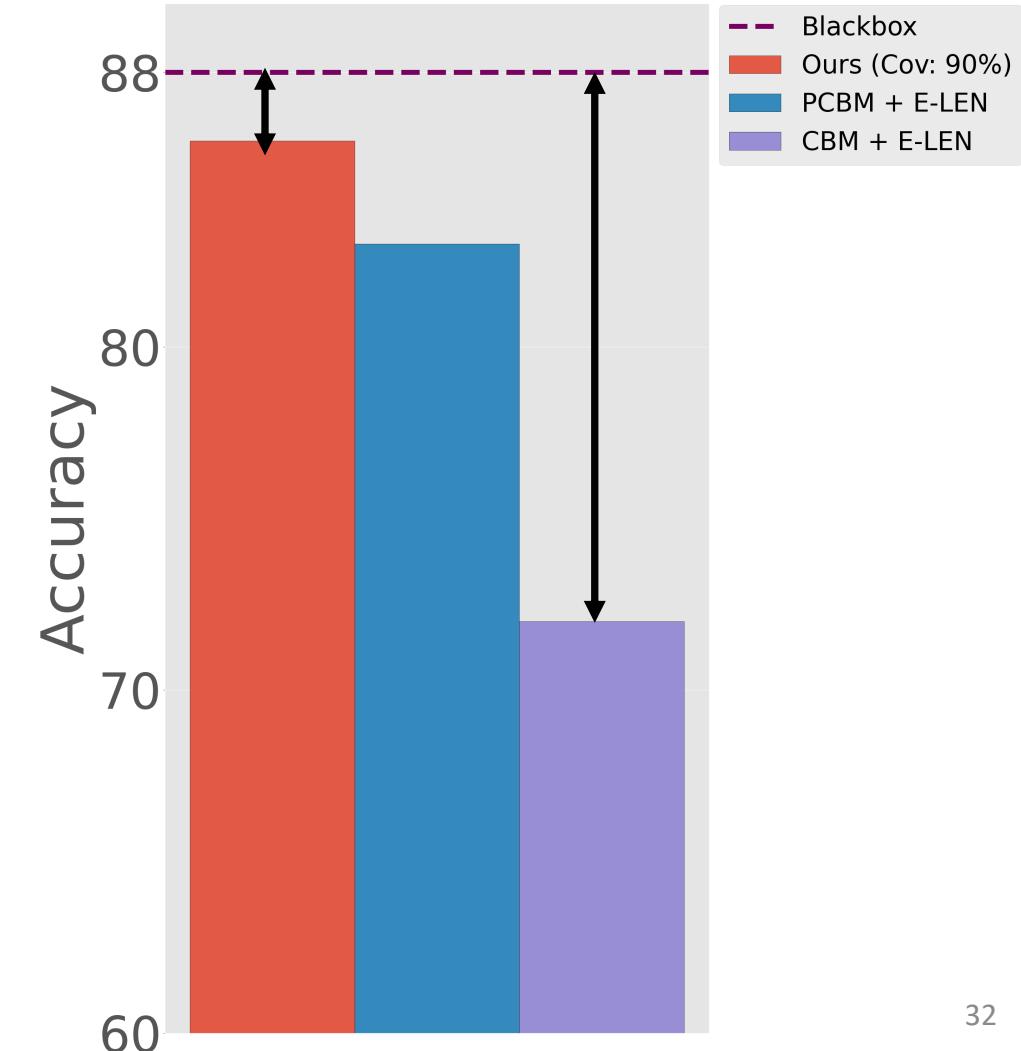


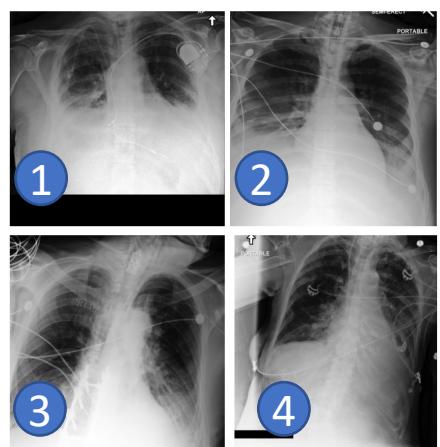
# Comparing Performance

CUB-200 with ViT

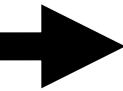


CUB-200 with ResNet101



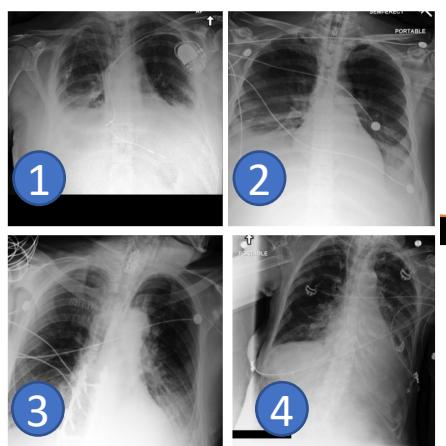


## Examples on Chest X-ray



Pleural unspec is “unspecified pleural effusion” referred to as “hydrothorax”.

Hydrothorax is a noninflammatory collection of serous fluid within the pleural cavities



# Examples on Chest X-ray



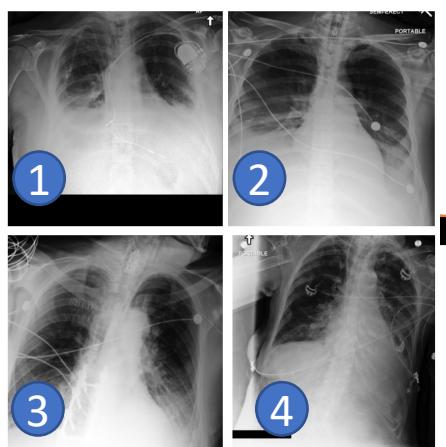
Expert 1

Effusion ↔  
left\_pleural  
^ right\_pleural  
^ pleural\_unspec

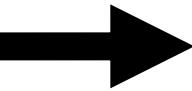
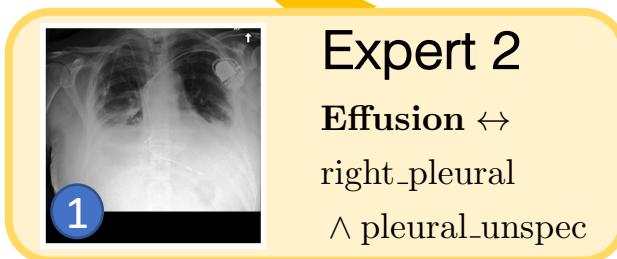


Pleural\_unspec is “unspecified pleural effusion” referred to as “hydrothorax”.

Hydrothorax is a noninflammatory collection of serous fluid within the pleural cavities



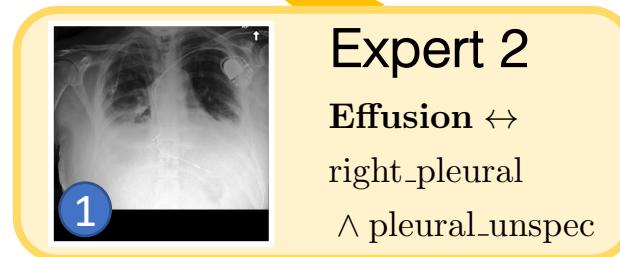
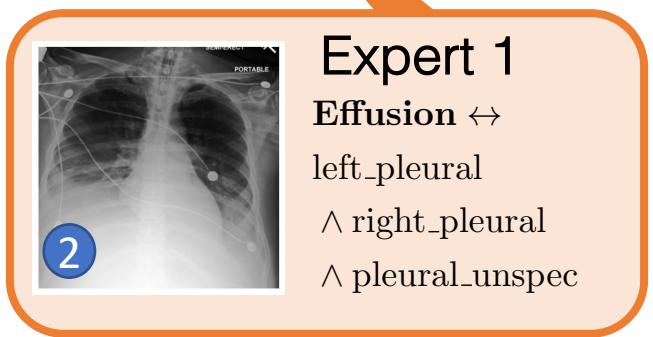
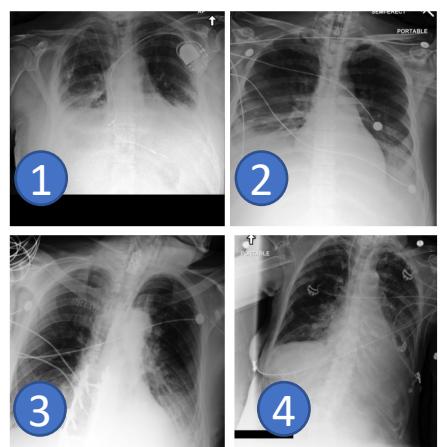
# Examples on Chest X-ray



Pleural\_unspec is “unspecified pleural effusion” referred to as “hydrothorax”.

Hydrothorax is a noninflammatory collection of serous fluid within the pleural cavities

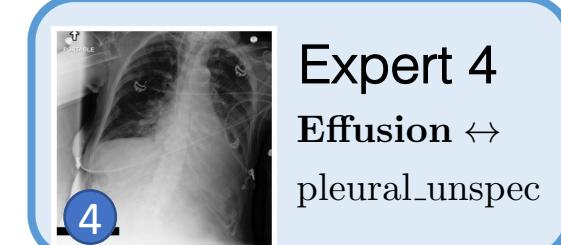
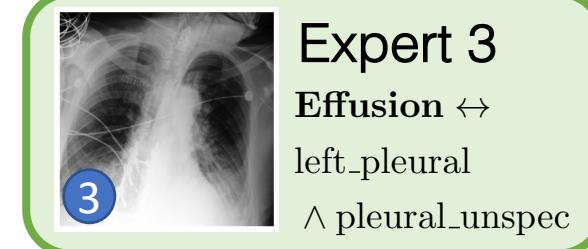
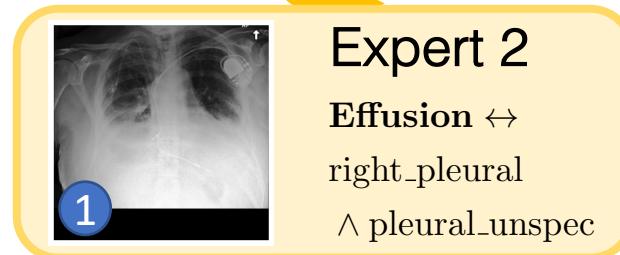
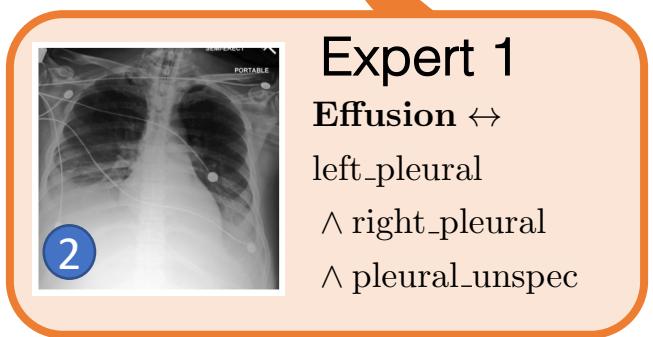
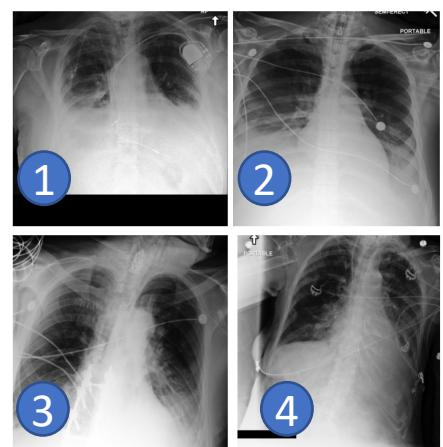
# Examples on Chest X-ray



Pleural\_unspec is “unspecified pleural effusion” referred to as “hydrothorax”.

Hydrothorax is a noninflammatory collection of serous fluid within the pleural cavities

# Examples on Chest X-ray

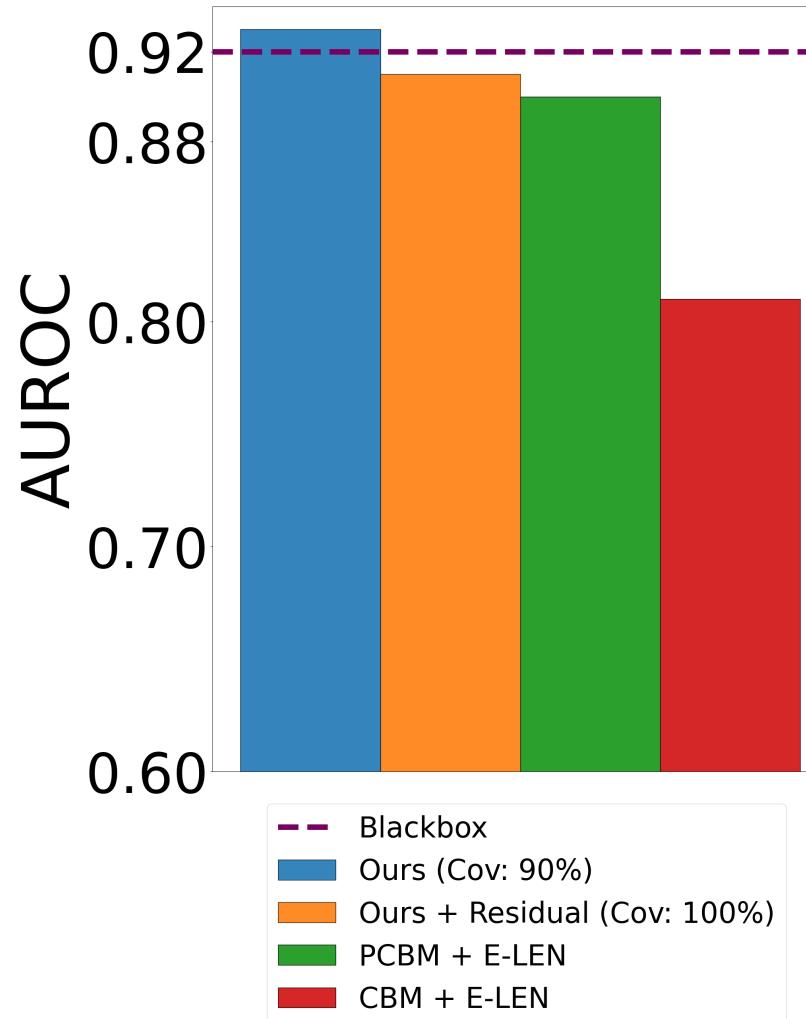


Pleural\_unspec is “unspecified pleural effusion” referred to as “hydrothorax”.

Hydrothorax is a noninflammatory collection of serous fluid within the pleural cavities

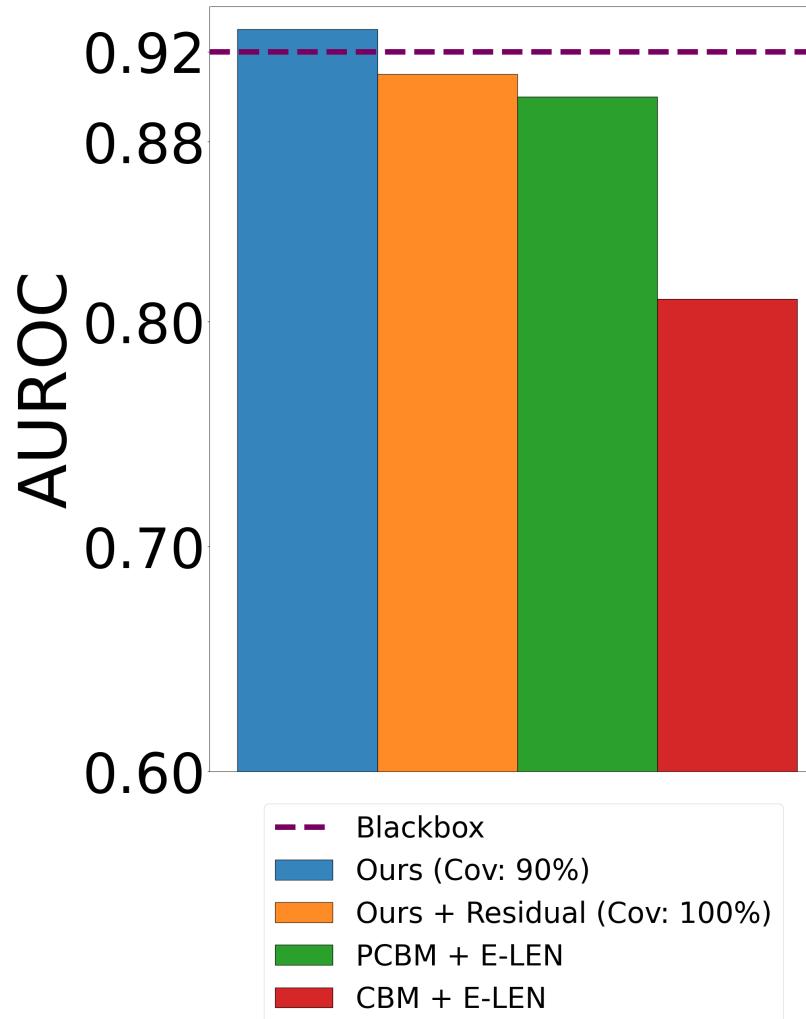
# Comparing Performance

## EFFUSION

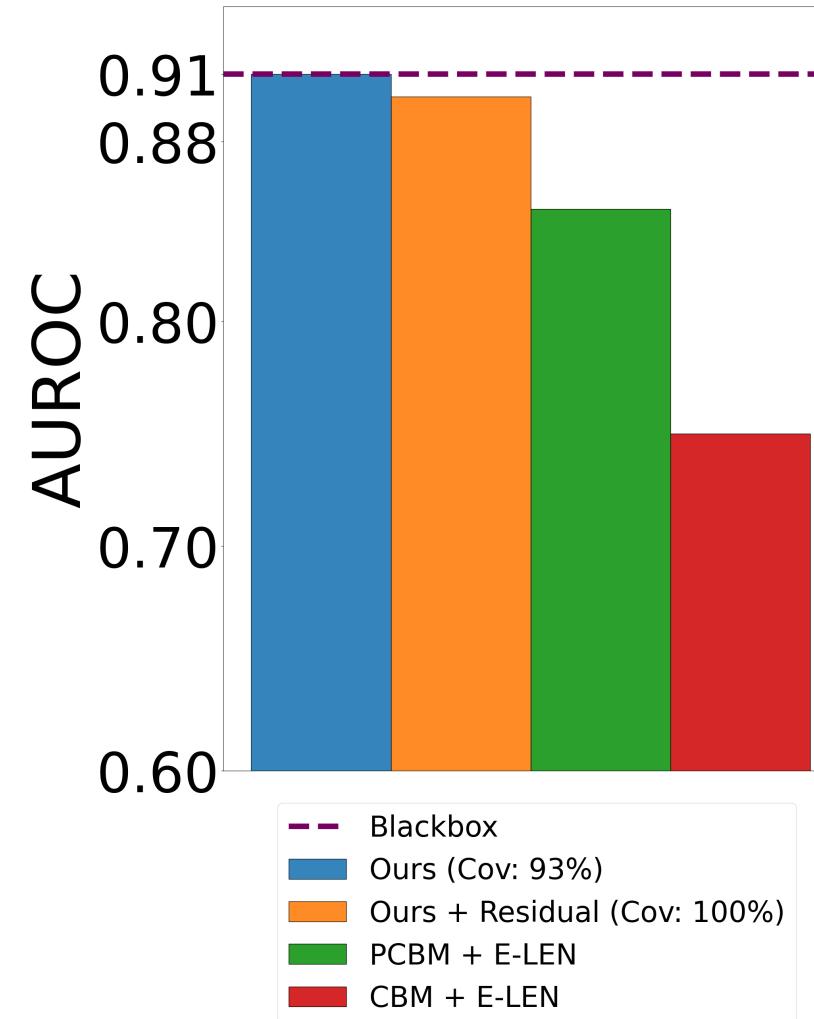


# Comparing Performance

## EFFUSION

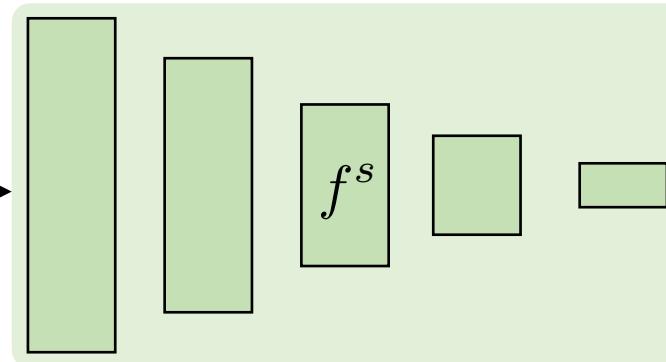


## PNEUMOTHORAX



# Data-Efficient Fine-tuning

MIMIC-CXR

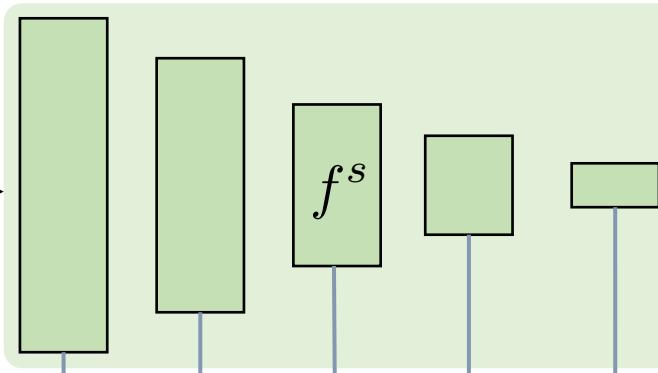


→ Pneumothorax

A horizontal arrow pointing from the model output towards the classification result.

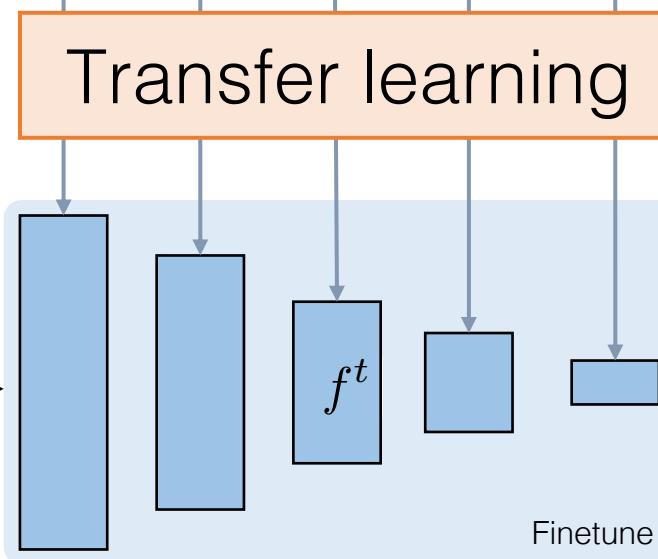
# Data-Efficient Fine-tuning

MIMIC-CXR



Pneumothorax

Stanford-CXR



Pneumothorax

# Data-Efficient Fine-tuning

MIMIC-CXR



Pneumothorax

Data and Computationally inefficient

Stanfo



Finetune

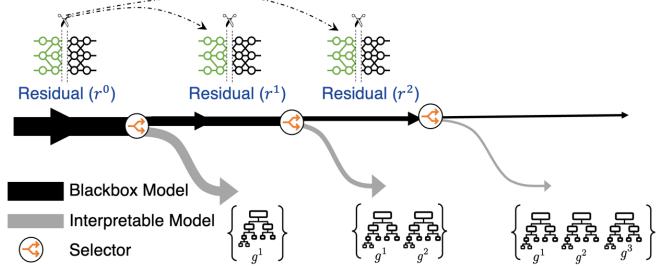
Pneumothorax

The clinical rules are “invariant”

# Fine-tune to a New Domain

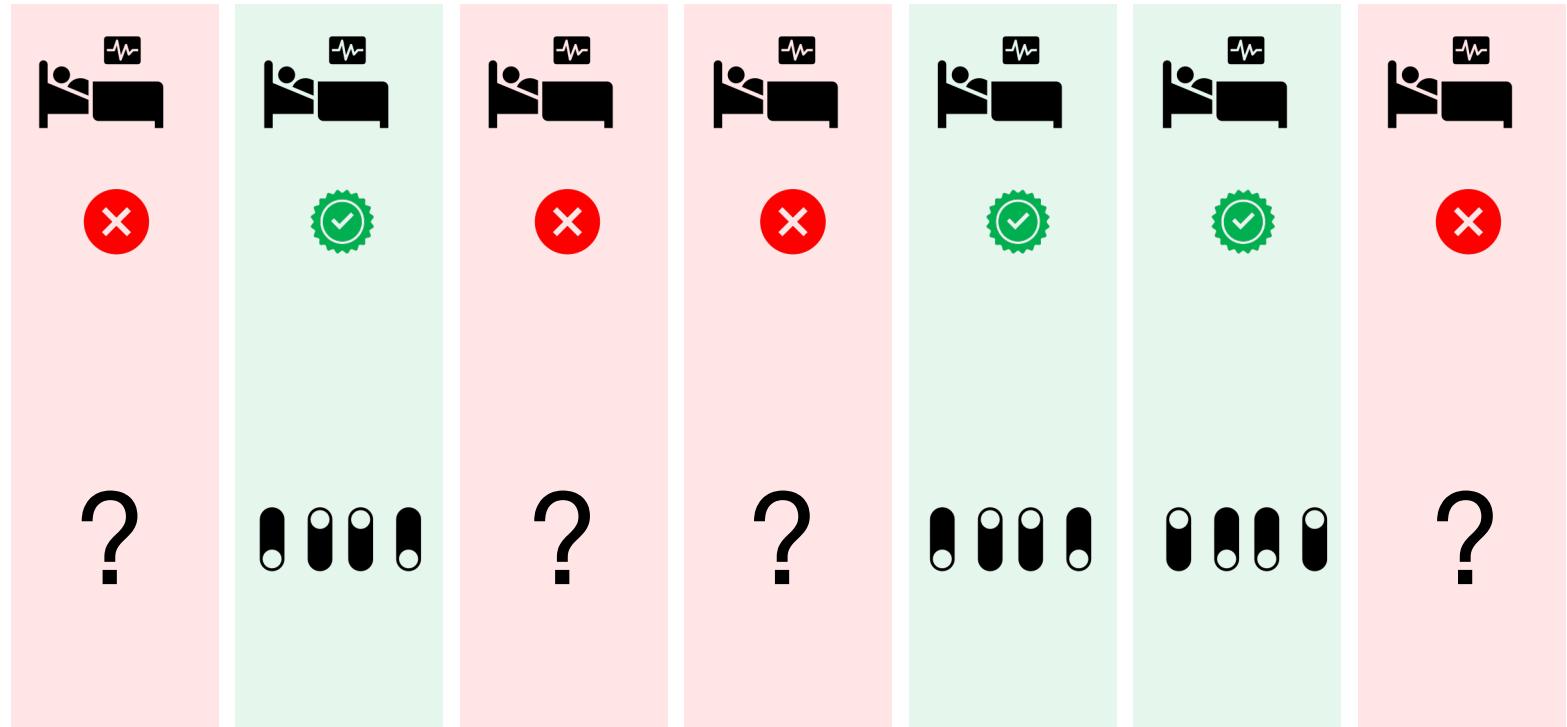
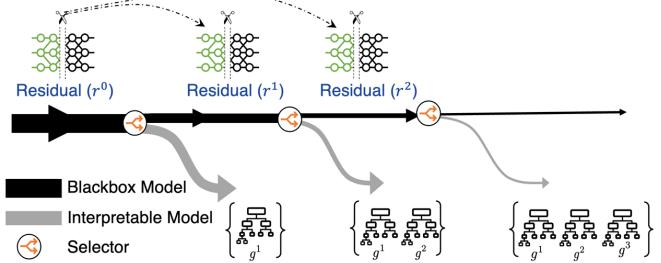
1

Apply source model



# Fine-tune to a New Domain

## 1 Apply source model

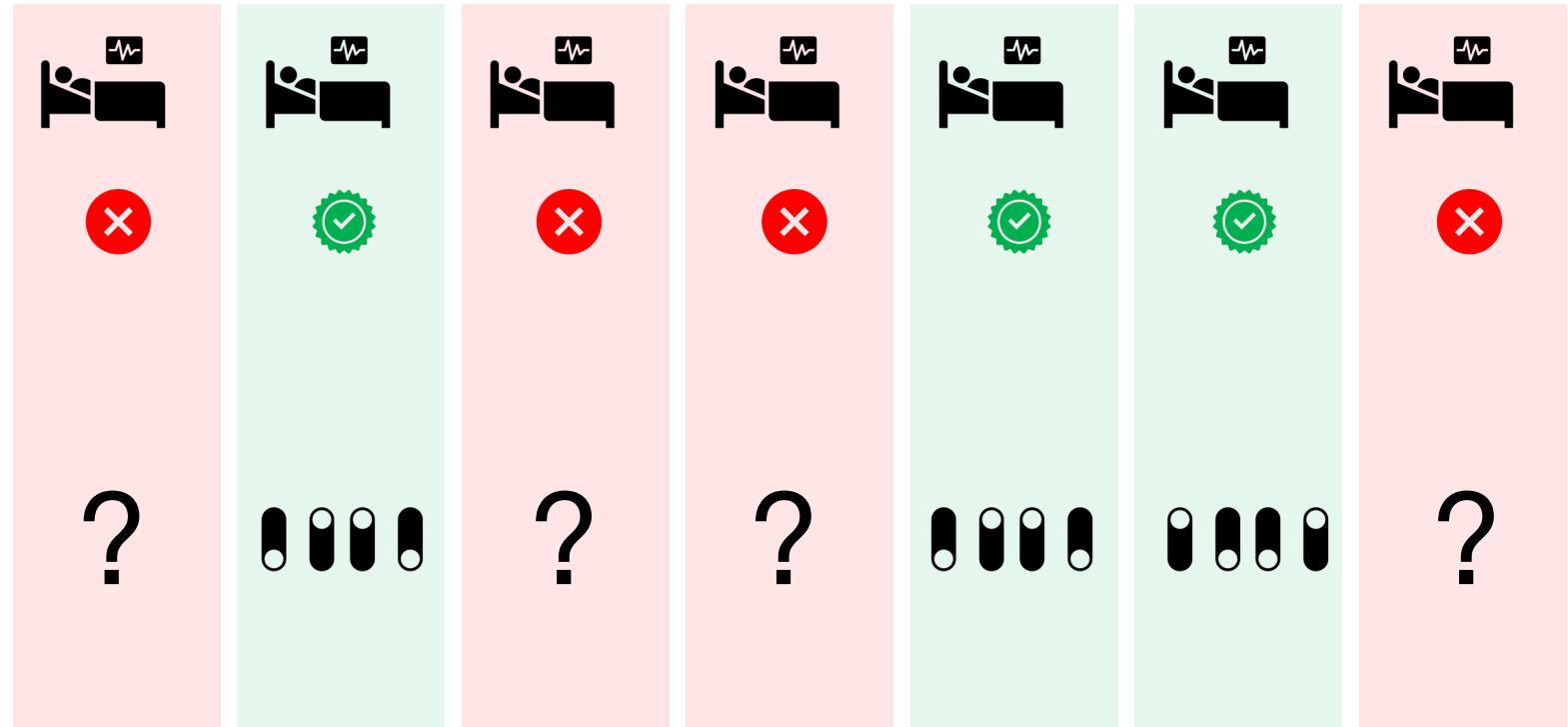
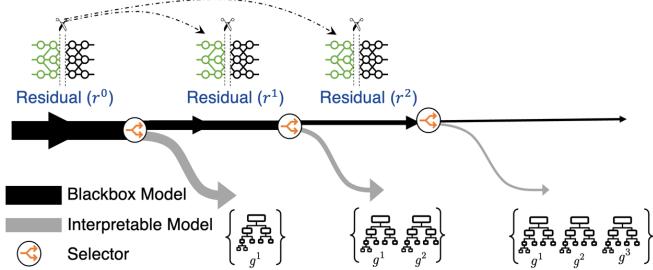


## 2 Use concepts from matching patients

C

# Fine-tune to a New Domain

## 1 Apply source model



## 2 Use concepts from matching patients

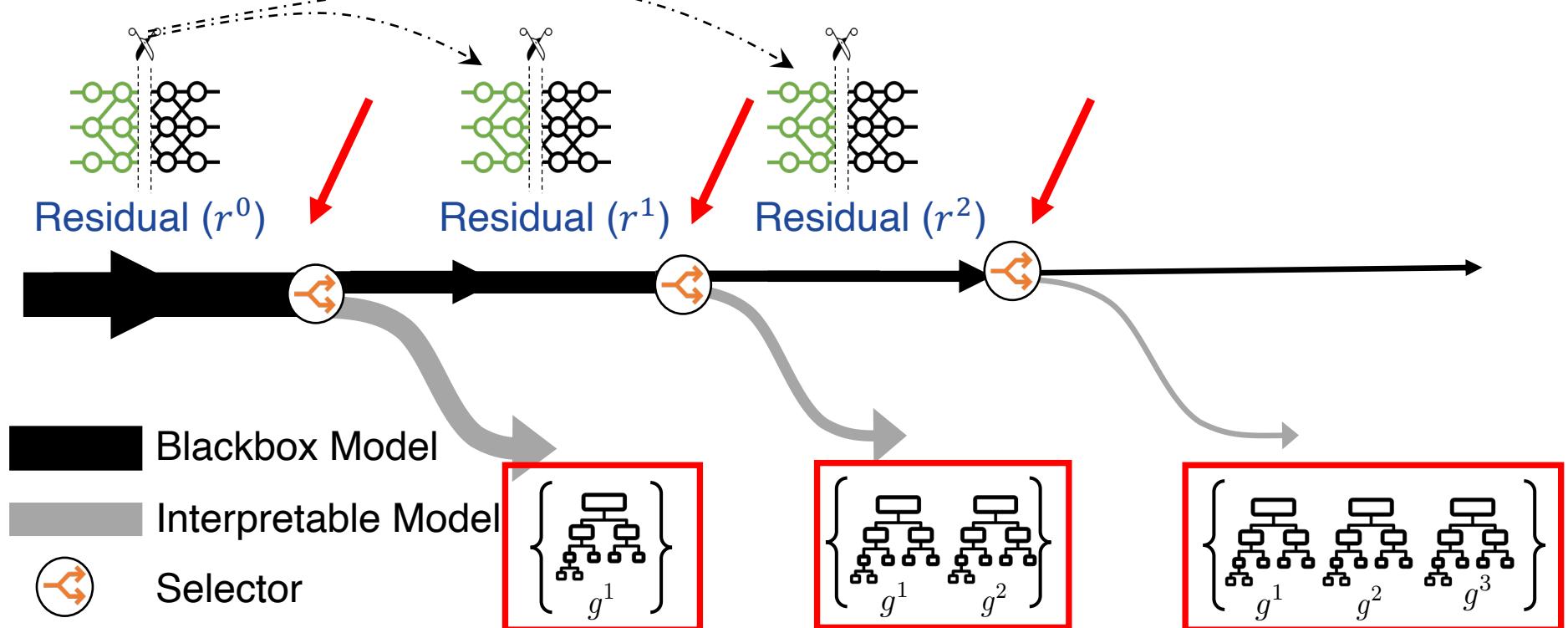
C

## 3 Propagate the concepts and update the concept extractor

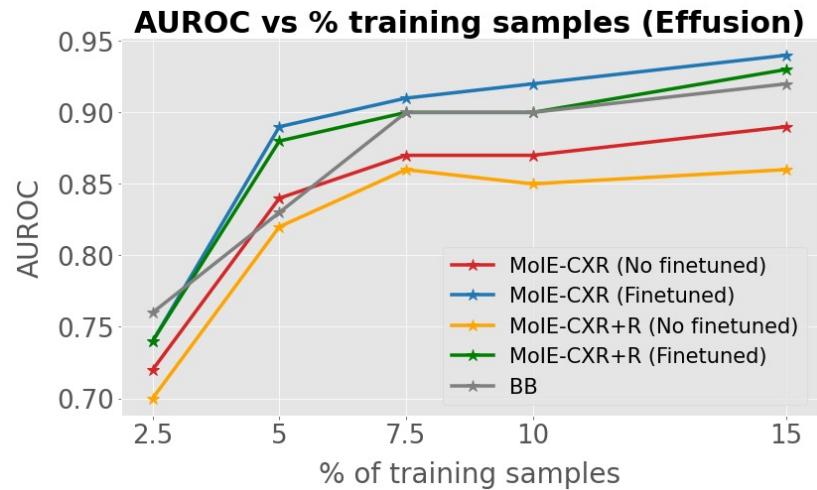
# Fine-tune to a New Domain

4

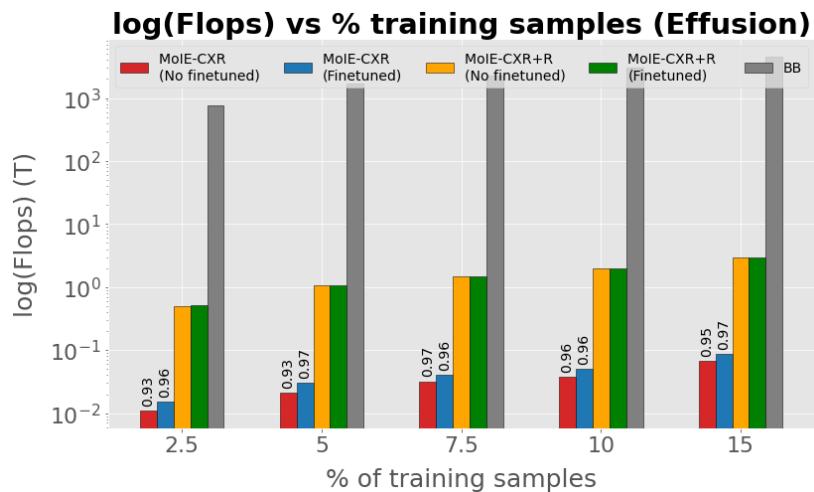
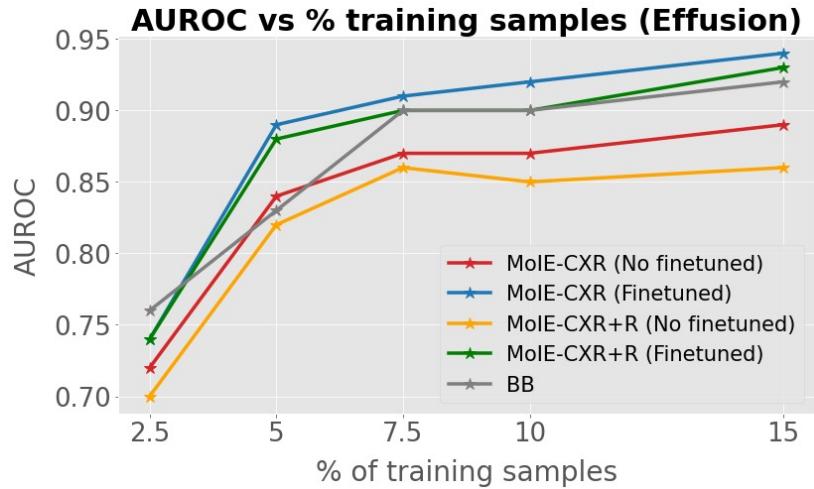
Update the selectors and interpretable models for 5 epochs



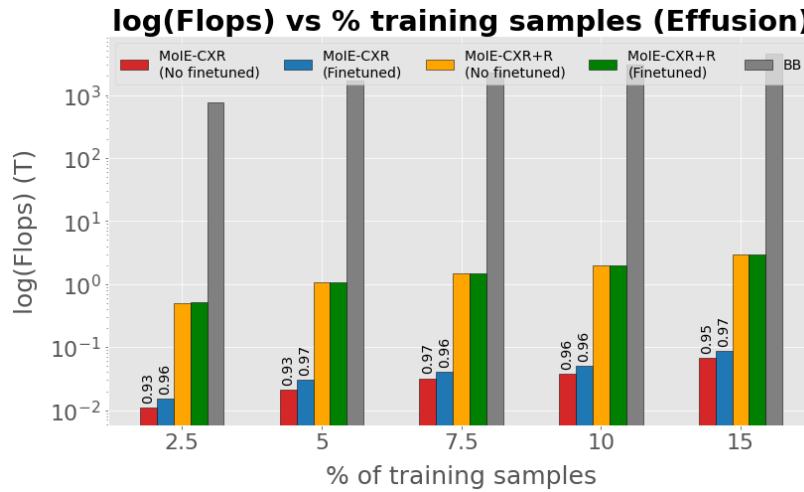
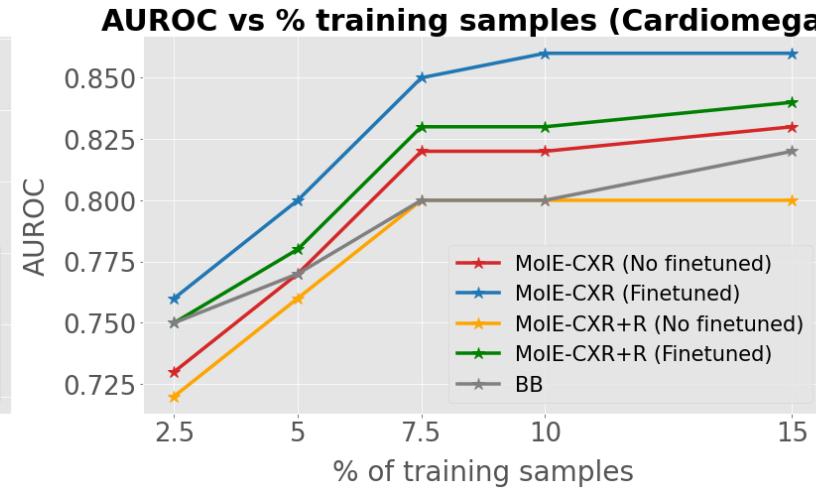
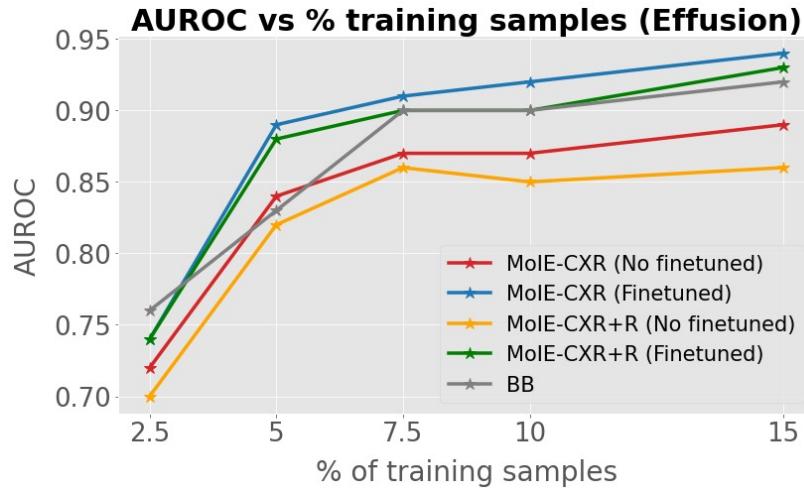
# Transferring to Stanford-CXR



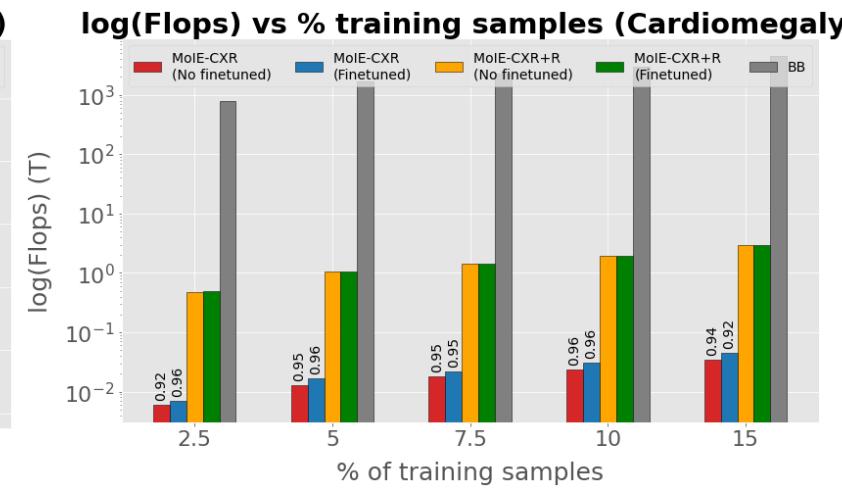
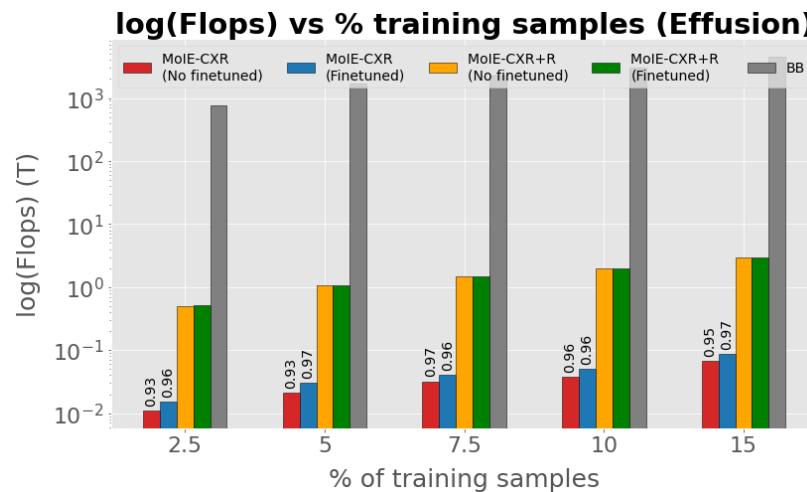
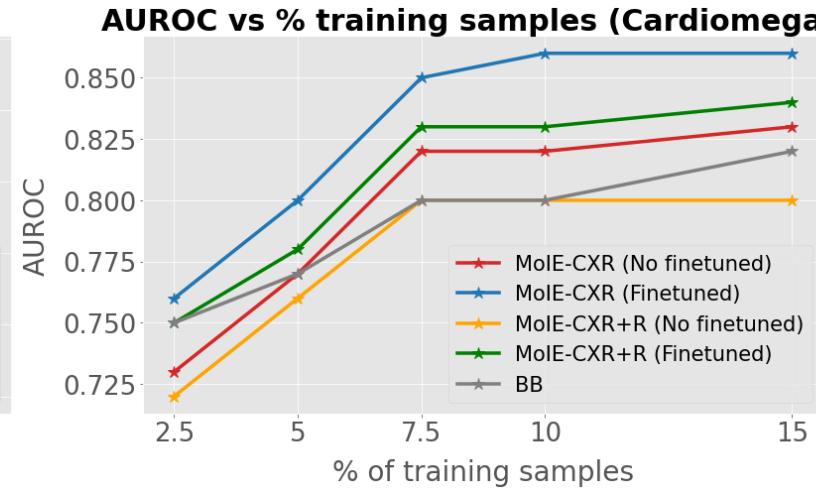
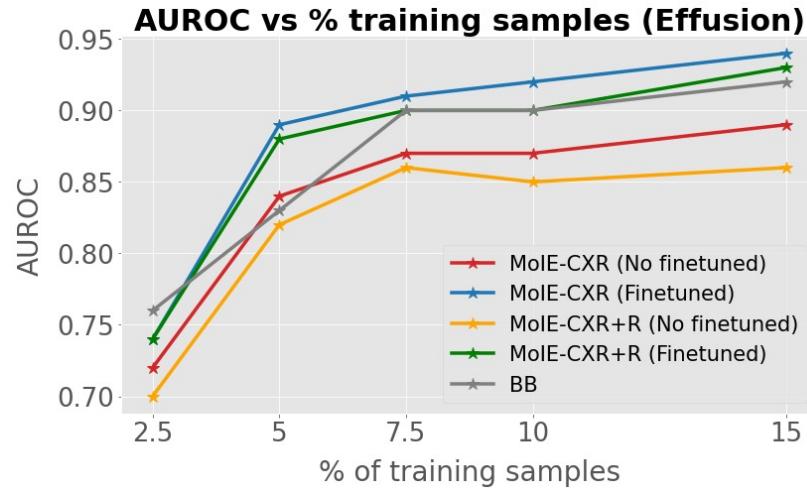
# Transferring to Stanford-CXR



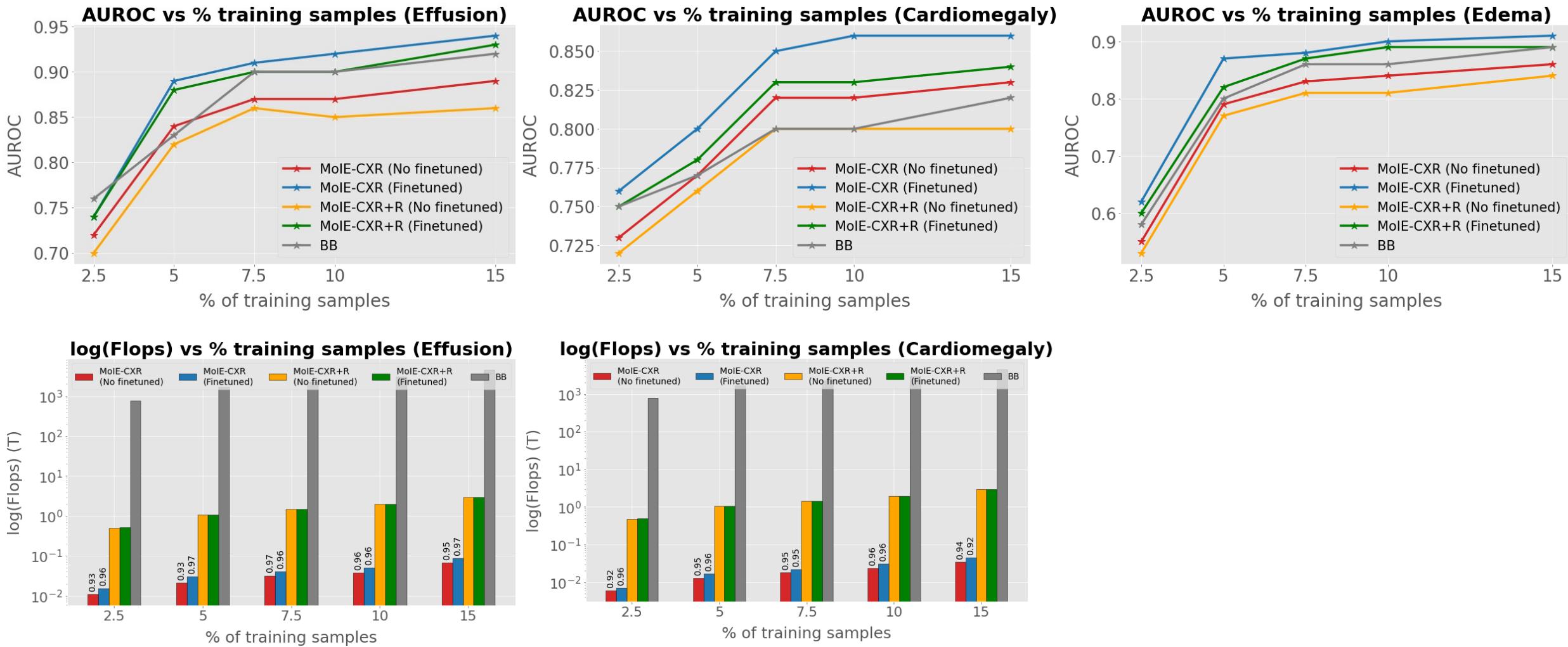
# Transferring to Stanford-CXR



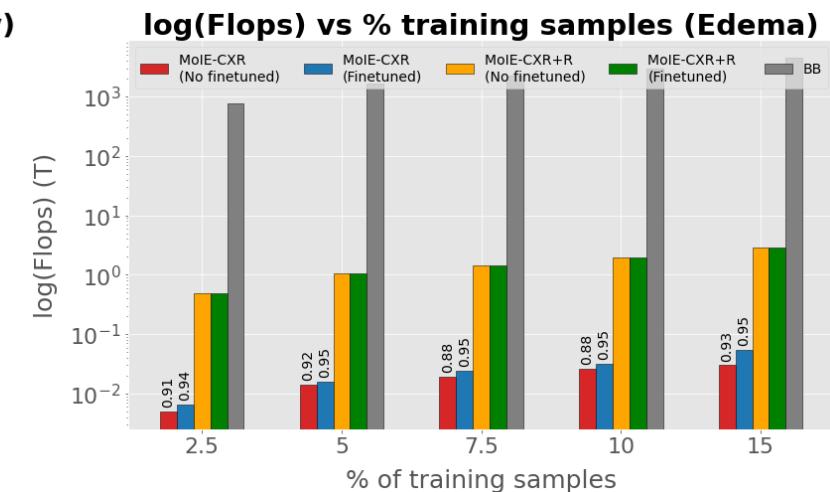
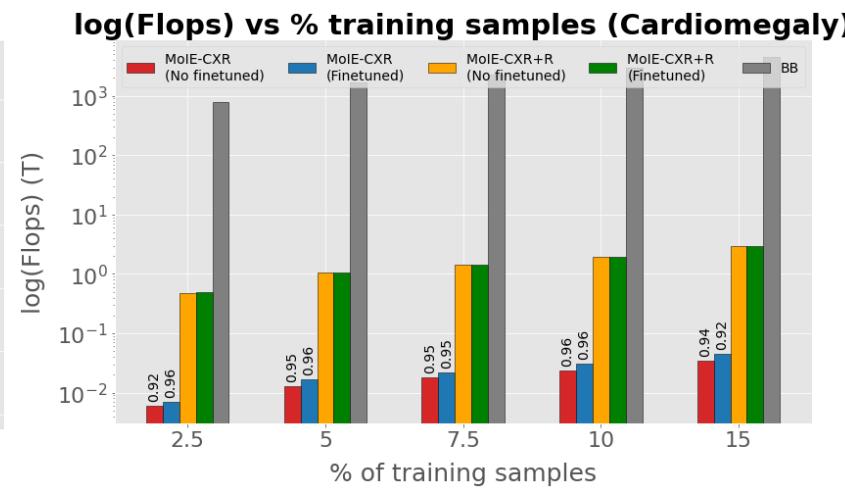
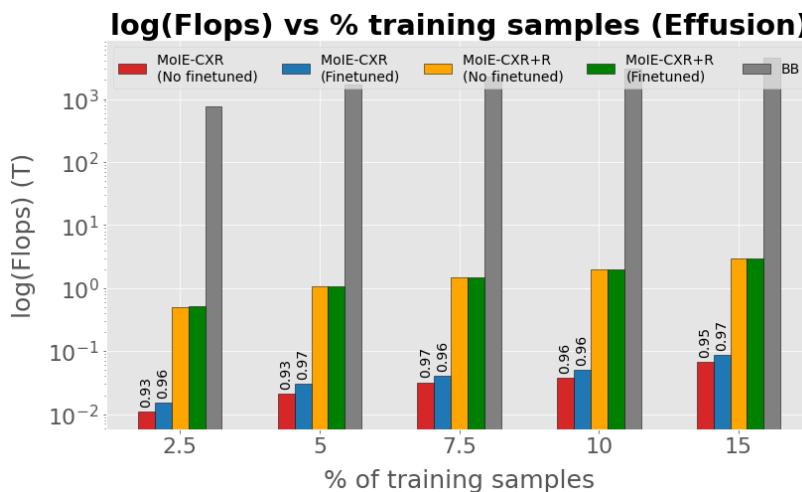
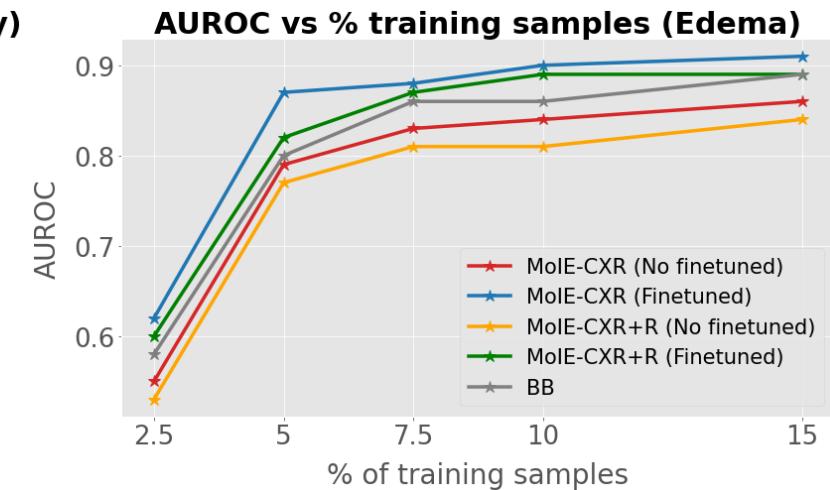
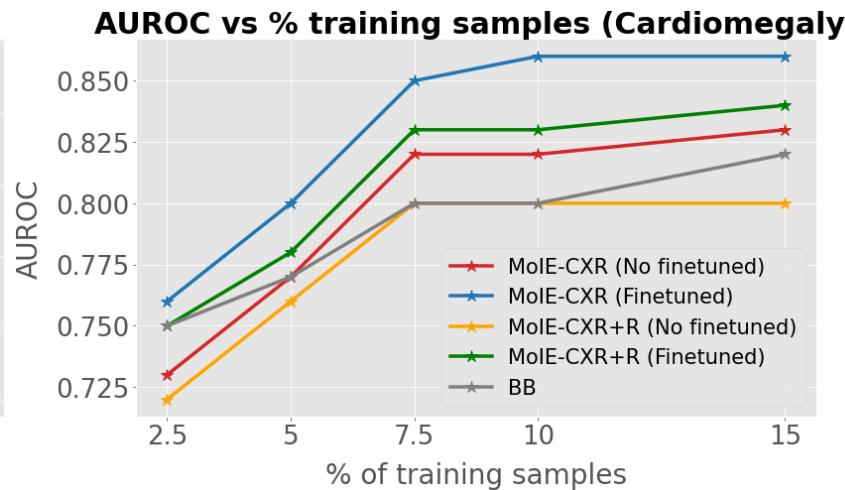
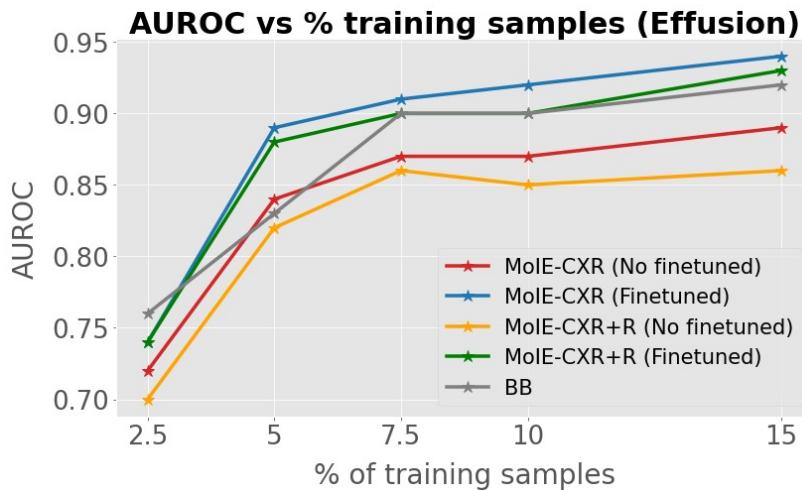
# Transferring to Stanford-CXR



# Transferring to Stanford-CXR



# Transferring to Stanford-CXR



# Conclusion

- Post-hoc approach allows the use of flexible Blackbox models, but post-hoc explainability is problematic and doesn't allow recourse.
- Interpretable-by-design allows recourse but suffer in performance.
- A mixture of interpretable models are carved out of a Blackbox model offering best of both worlds.
- Each interpretable model (expert) is modeled as First Order Logic (FOL), but other choices are possible for the interpretable model.
- Transfer Learning is more efficient (data & compute) with the new model.
- Intervention is possible with the new interpretable model.

# References

- Ghosh S, Yu K, Arabshahi F, Batmanghelich K, “Dividing and Conquering a BlackBox to a Mixture of Interpretable Models: Route, Interpret, Repeat,” **ICML 2023**
- Ghosh S, Yu K, Batmanghelich K, “Distilling BlackBox to Interpretable models for Efficient Transfer Learning,” **MICCAI 2023, Early accept ~ 13%**
- Ghosh S, Yu K, Arabshahi F, Batmanghelich K, “Tackling Shortcut Learning in Deep Neural Networks: An Iterative Approach with Interpretable Models,” **SCIS workshop, ICML 2023**
- Ghosh S, Yu K, Batmanghelich K, “Bridging the Gap: From Post Hoc Explanations to Inherently Interpretable Models for Medical Imaging,” **IMLH workshop, ICML 2023**



Shantanu Ghosh<sup>1</sup>, Ke Yu<sup>2</sup>, Kayhan Batmanghelich<sup>1</sup>

<sup>1</sup>BU ECE, <sup>2</sup>Pitt ISP



# Thank you

