

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



Shantanu Ghosh¹, Ke Yu², Kayhan Batmanghelich¹

¹BU ECE, ²Pitt ISP



Desiderata of Explainability in Healthcare

Trust

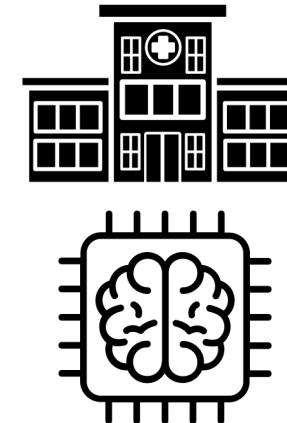
Causality

Transferability

Informativeness

Ethical Decision making

Training



Deployment



Desiderata of Explainability in Healthcare

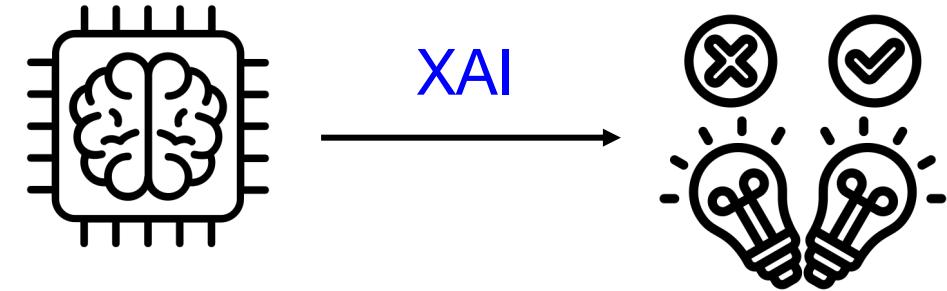
Trust

Causality

Transferability

Informativeness

Ethical Decision making



Obtaining genetics insights from
deep learning via explainable artificial
intelligence

Gherman Novakovsky ^{1,2,7}, Nick Dexter  ^{3,4,7}, Maxwell W. Libbrecht  ^{4,8}✉,
Wyeth W. Wasserman  ^{1,8}✉ and Sara Mostafavi  ^{5,6,8}✉

Desiderata of Explainability in Healthcare

Trust

Causality

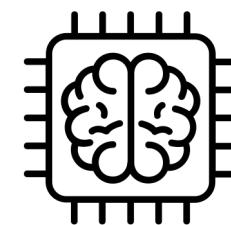
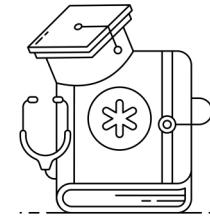
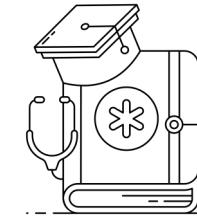
Transferability

Informativeness

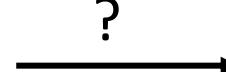
Ethical Decision making

Training

Deployment



?



Desiderata of Explainability in Healthcare

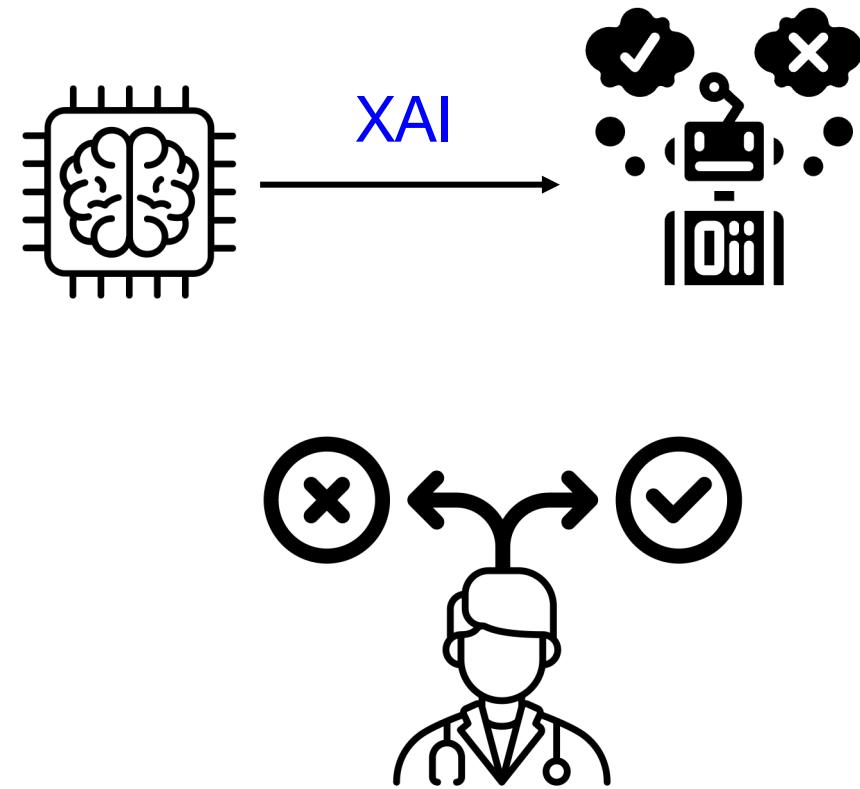
Trust

Causality

Transferability

Informativeness

Ethical Decision making



Desiderata of Explainability in Healthcare

Trust

Causality

Transferability

Informativeness

Ethical Decision making

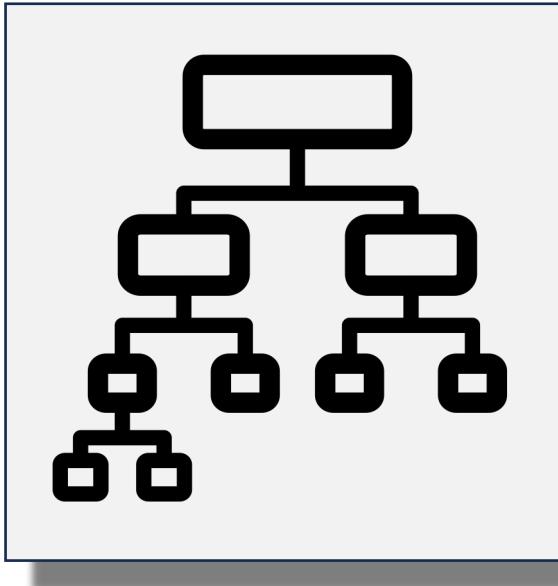
Training

Deployment

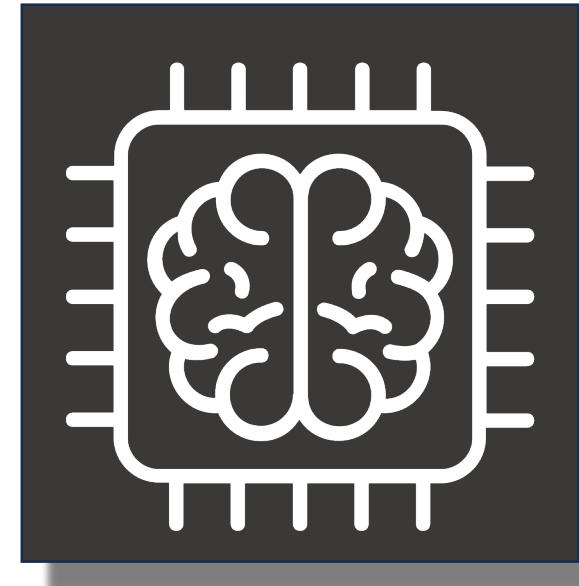


General Design Choices

Transparent Models

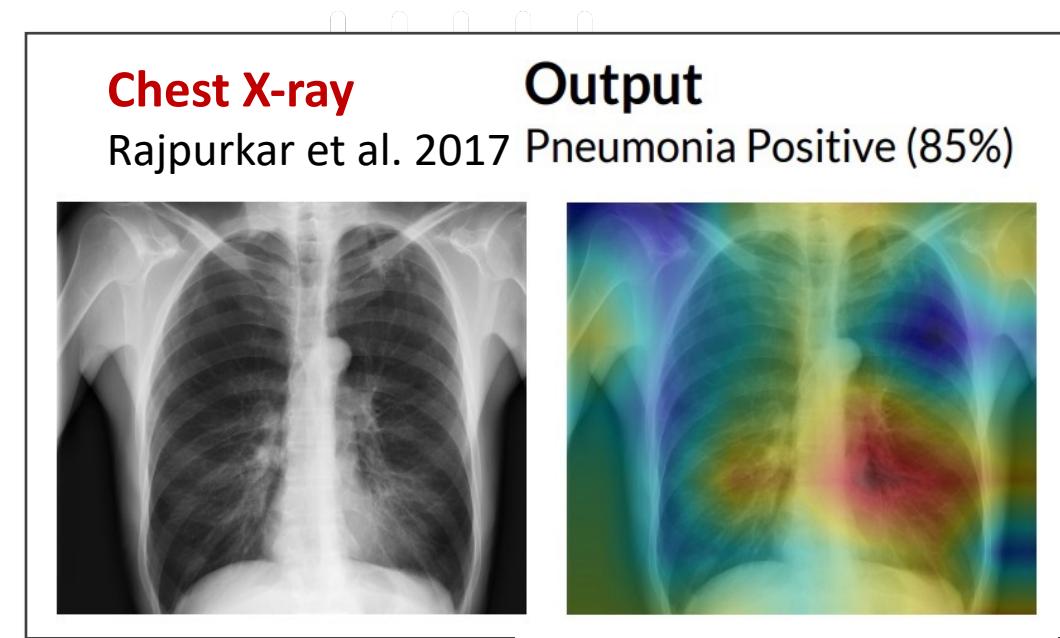


Post-hoc Explanation



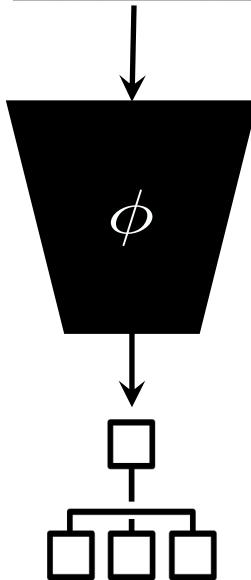
General Design Choices

Post-hoc Explanation

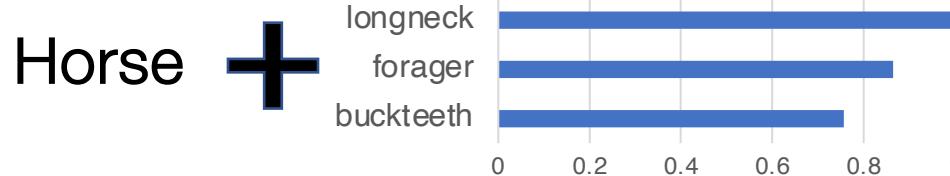


General Design Choices

Transparent Models



Top 3 concepts to identify Horse



Post-hoc Explanation

Output

Chest X-ray

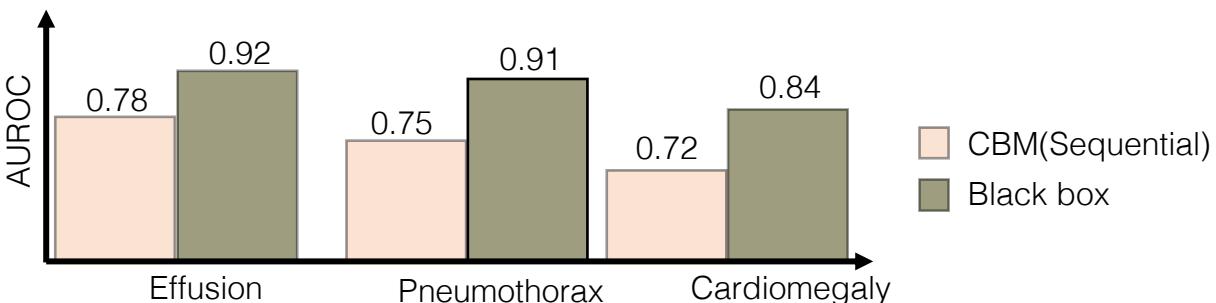
Rajpurkar et al. 2017

Pneumonia Positive (85%)

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



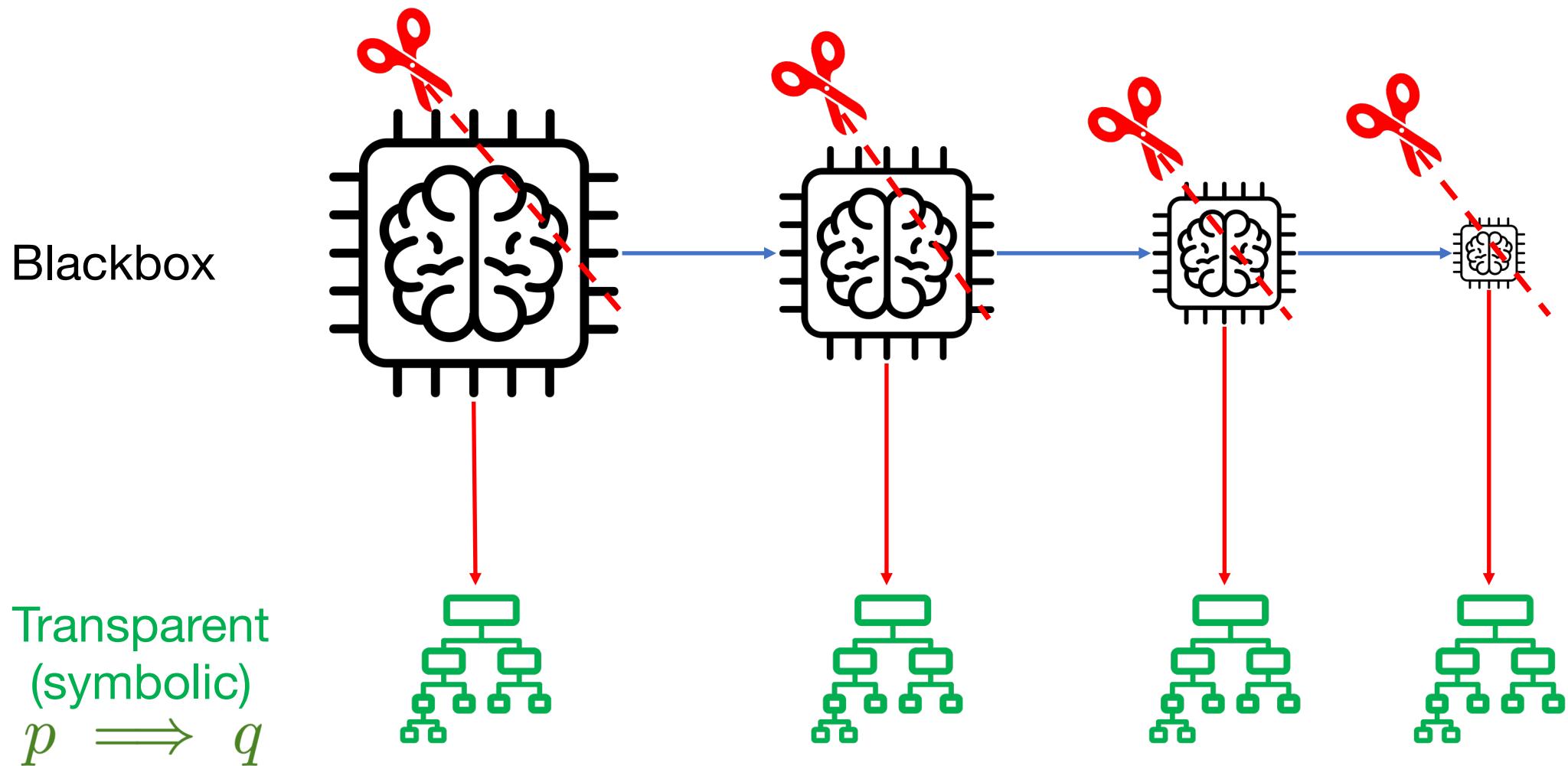
Desiderata

- 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

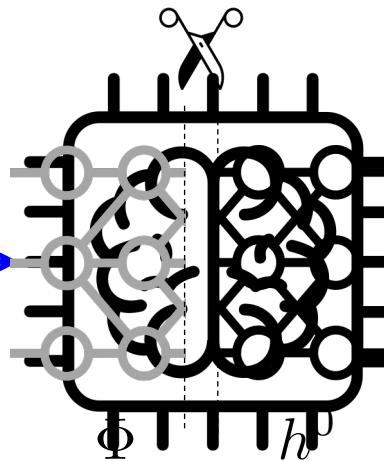


Problem Set Up

χ



γ



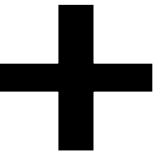
Dog

c

thin ears
shortened muzzle
round feet

....

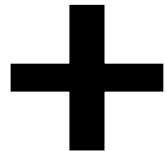
Problem Set Up



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:

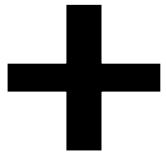
Right upper lobe consolidation with adjacent.
While this may be infectious in nature, a CT
scan is recommended for further clarification.

parse the reports to get the concepts

C

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

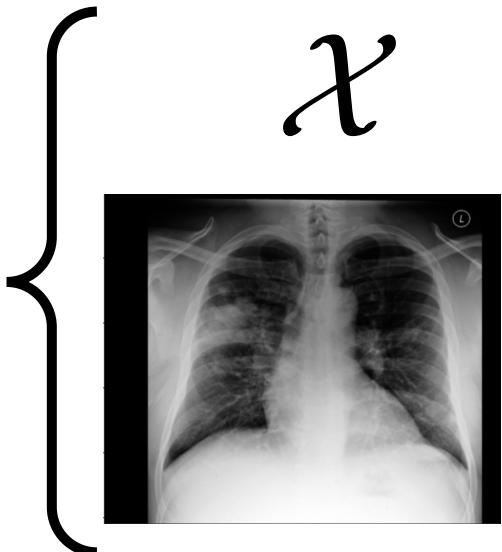
Problem Set Up



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

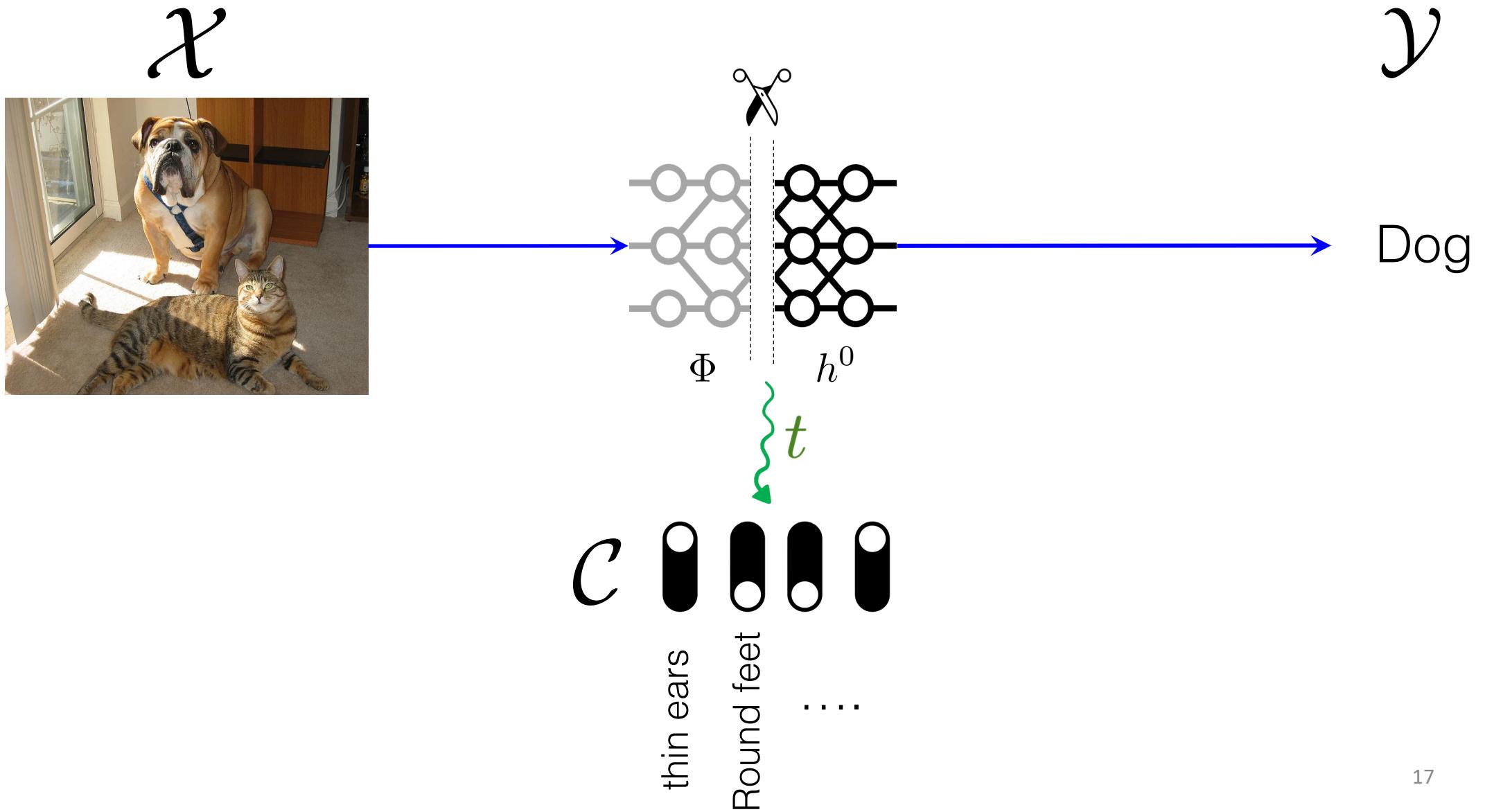
right upper lobe
left lower lobe
heart size
....

γ

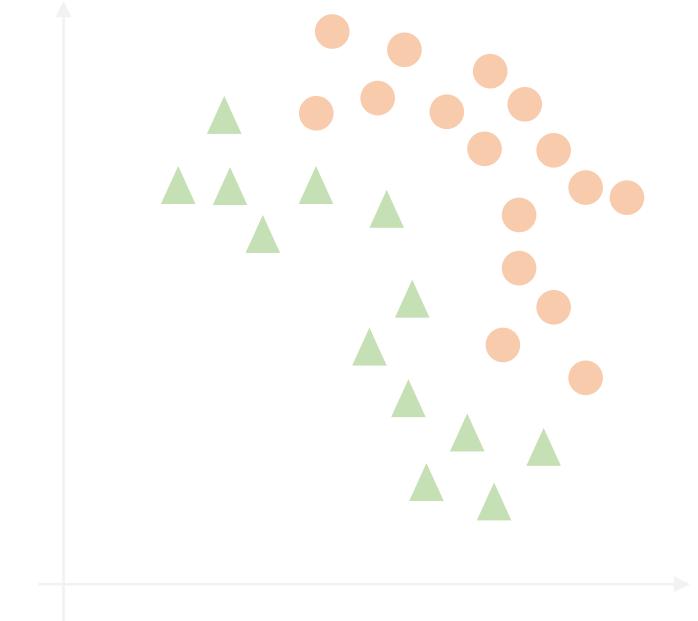
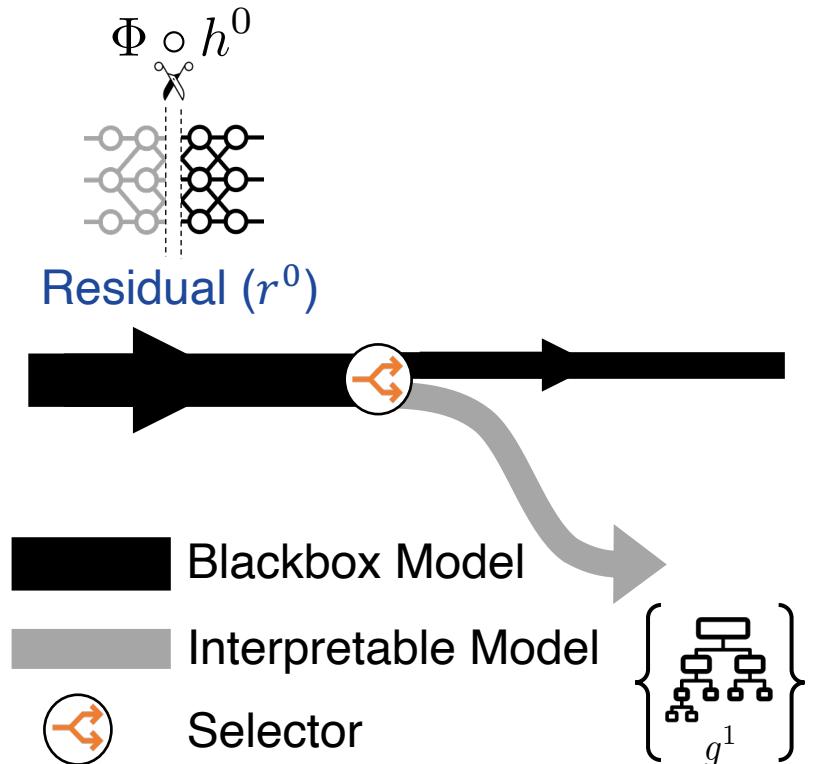
Consolidation



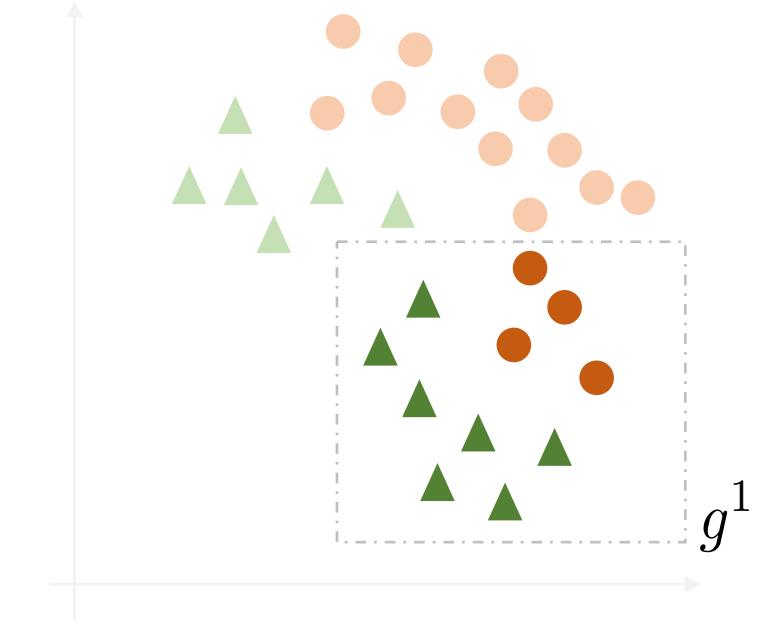
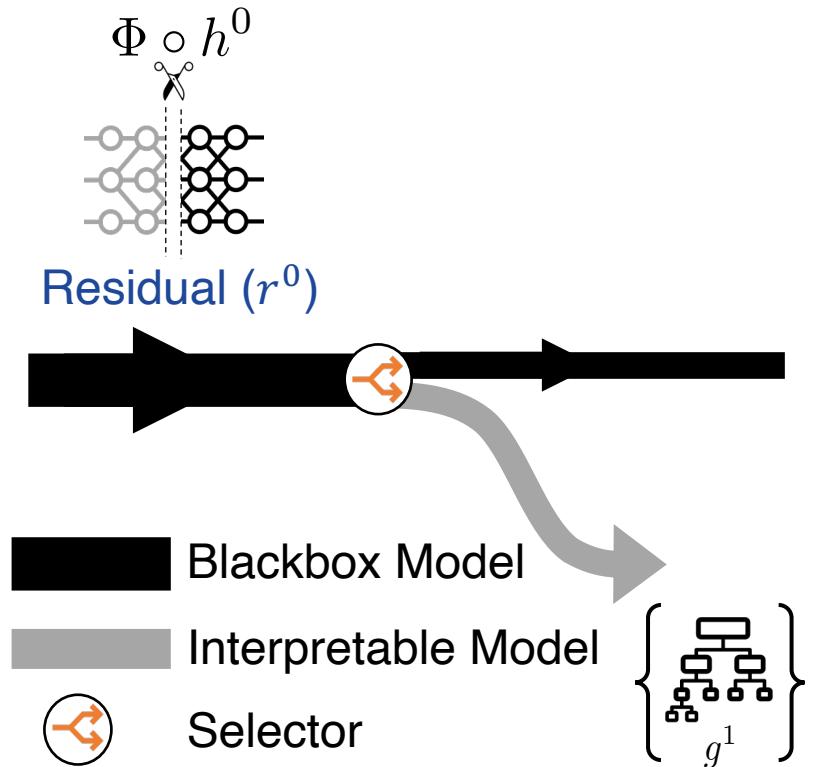
Discovering Hidden Concepts



Carving out Interpretable Models

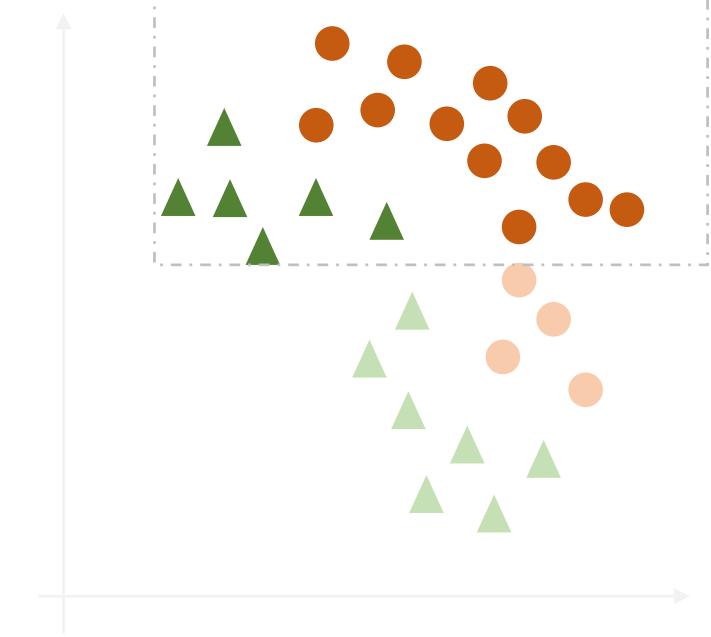
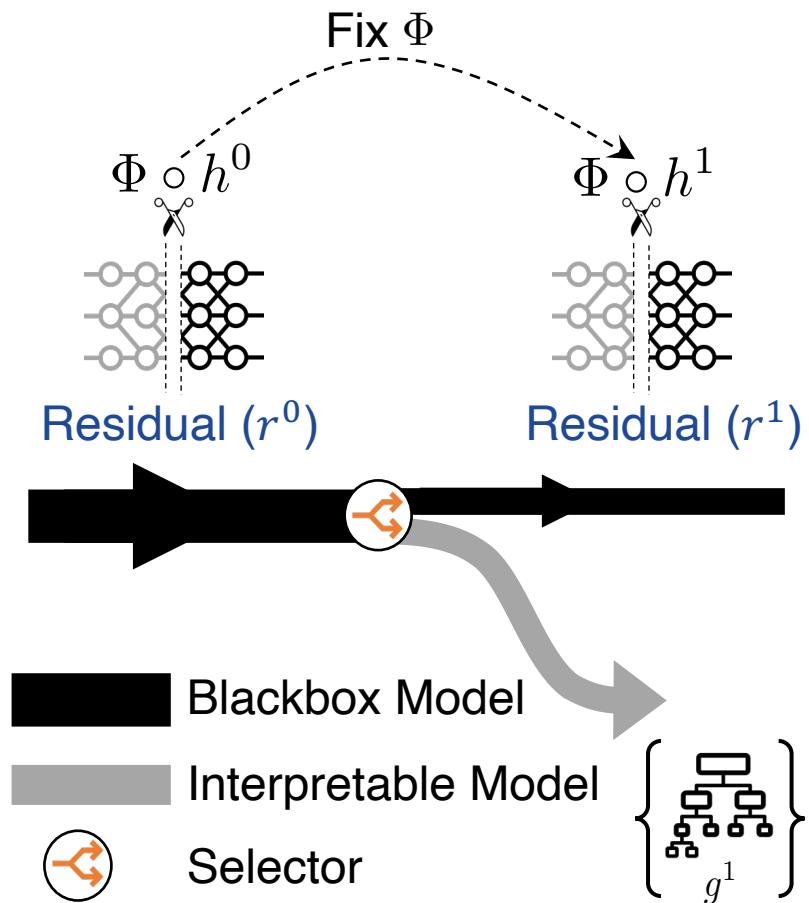


Carving out Interpretable Models



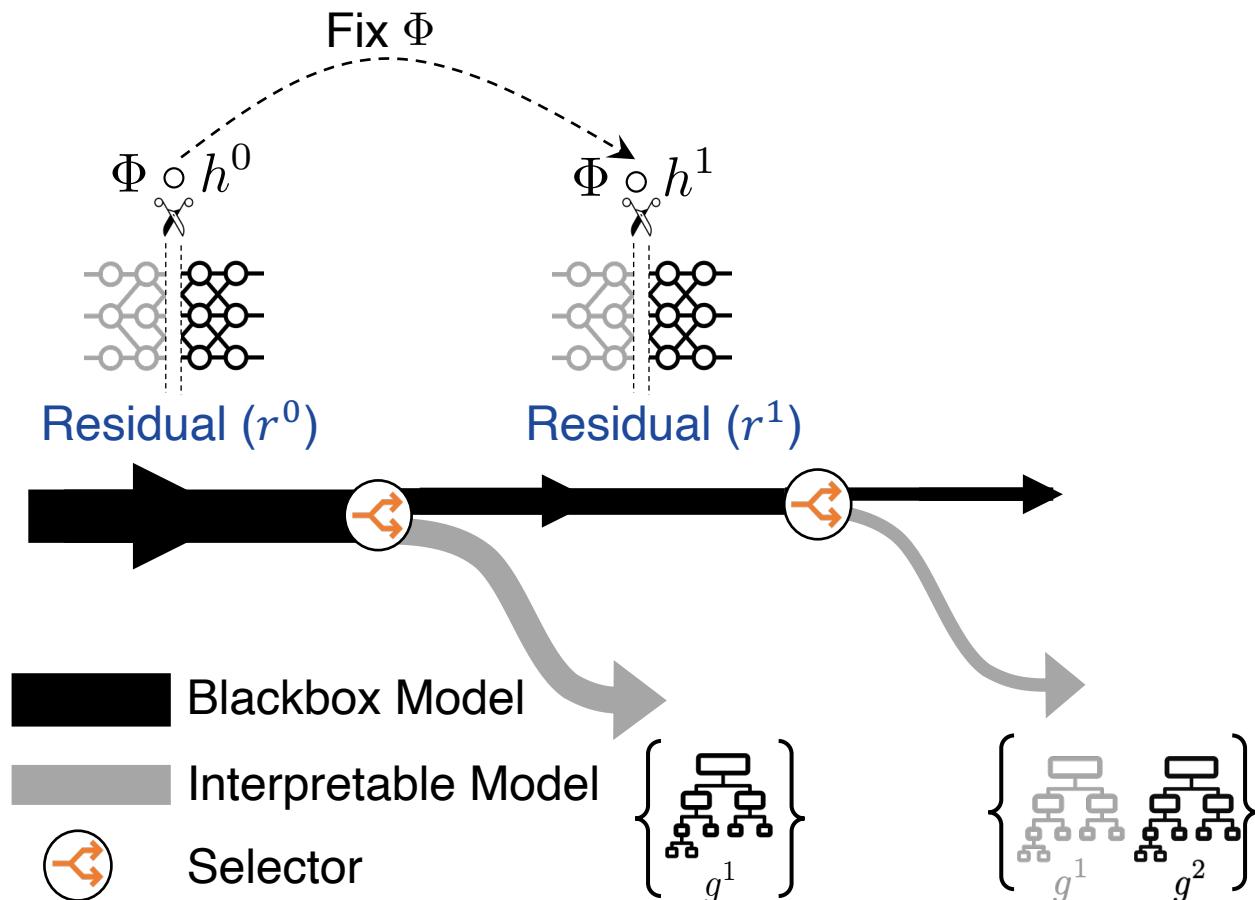
Each g is E-LENs (Barberio et al. AAAI 2022) to produce sample specific FOLs.

Carving out Interpretable Models



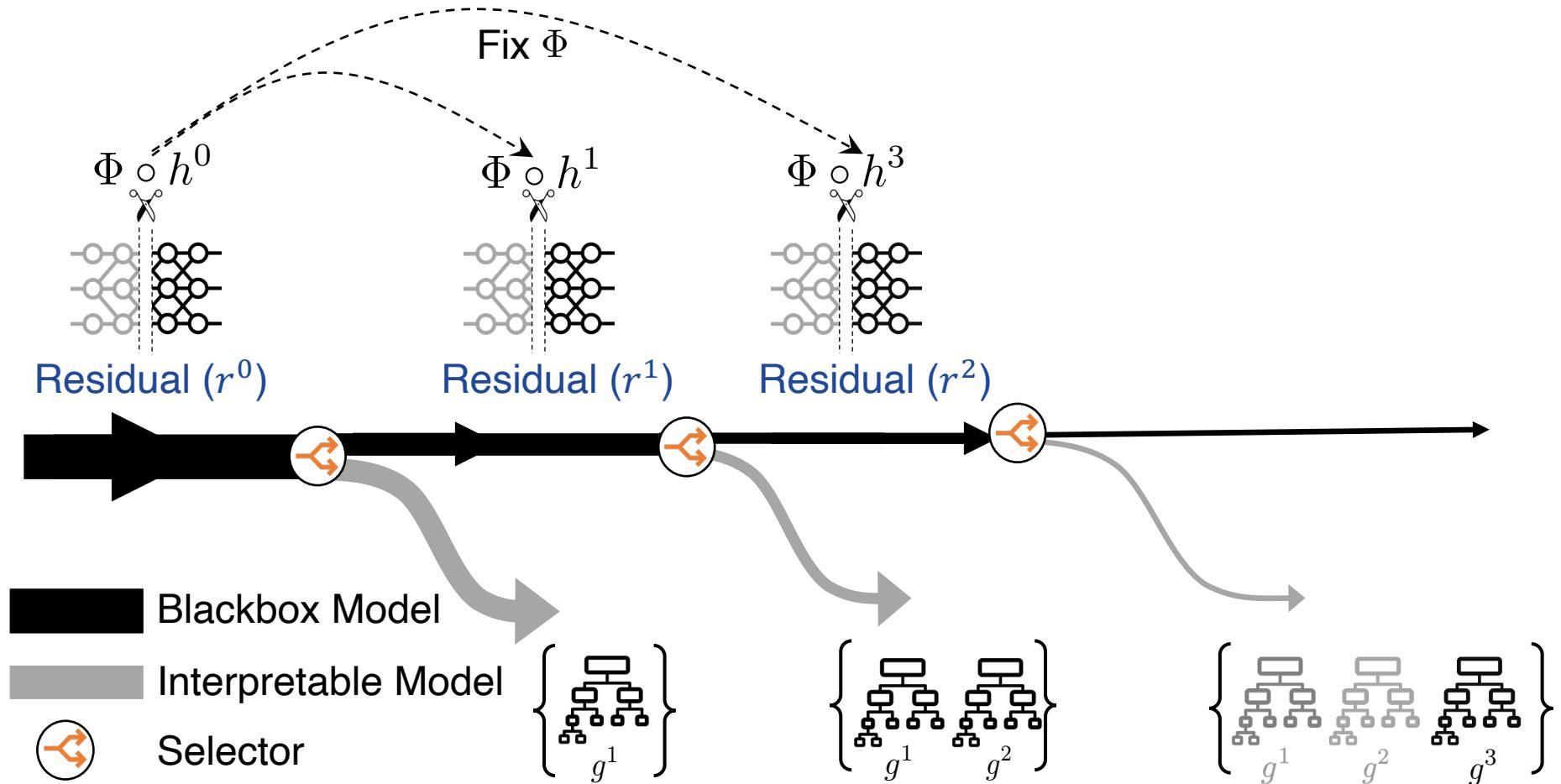
Each g is E-LENs (Barberio et al. AAAI 2022) to produce sample specific FOLs.

Carving out Interpretable Models



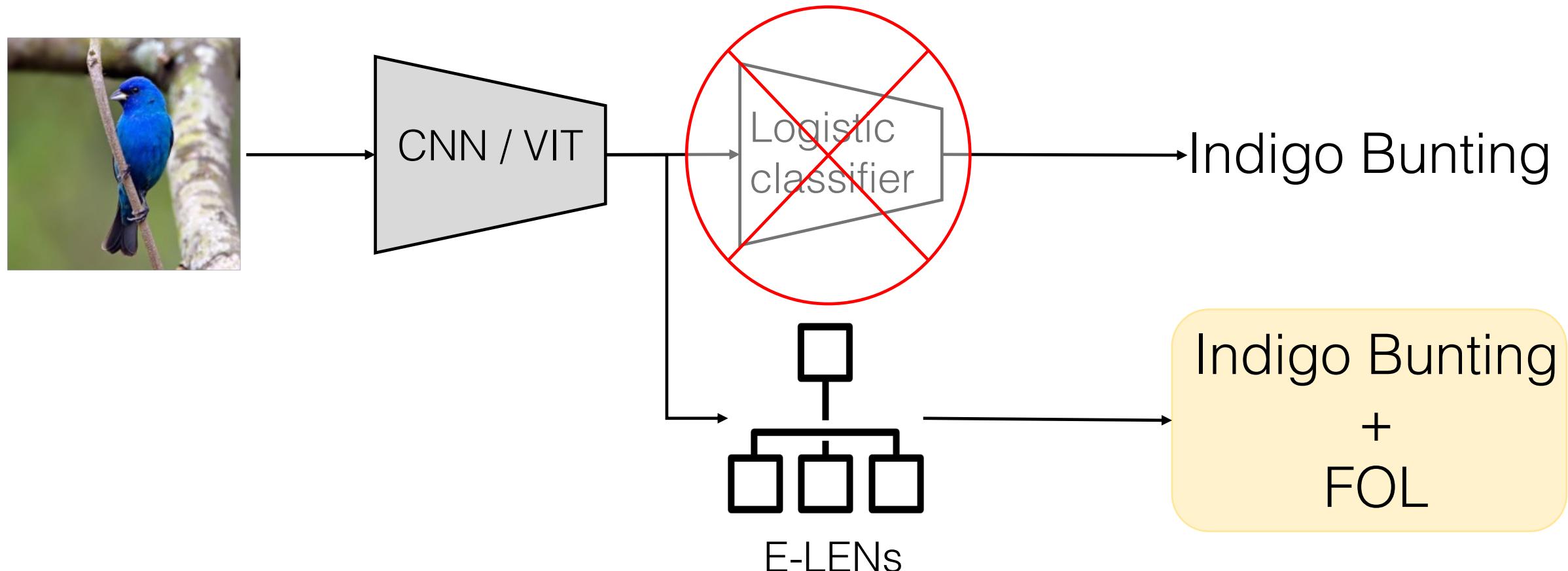
Each g is E-LENs (Barberio et al. AAAI 2022) to produce sample specific FOLs.

Carving out Interpretable Models



Each g is E-LENs (Barberio et al. AAAI 2022) to produce sample specific FOLs.

Baselines



Replace the standard logistic classifiers in CBM (Koh et al. ICML 2020) and post-hoc CBM (Yuksekgonul et al. ICLR 2023) with E-LENs to compare the FOLs

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

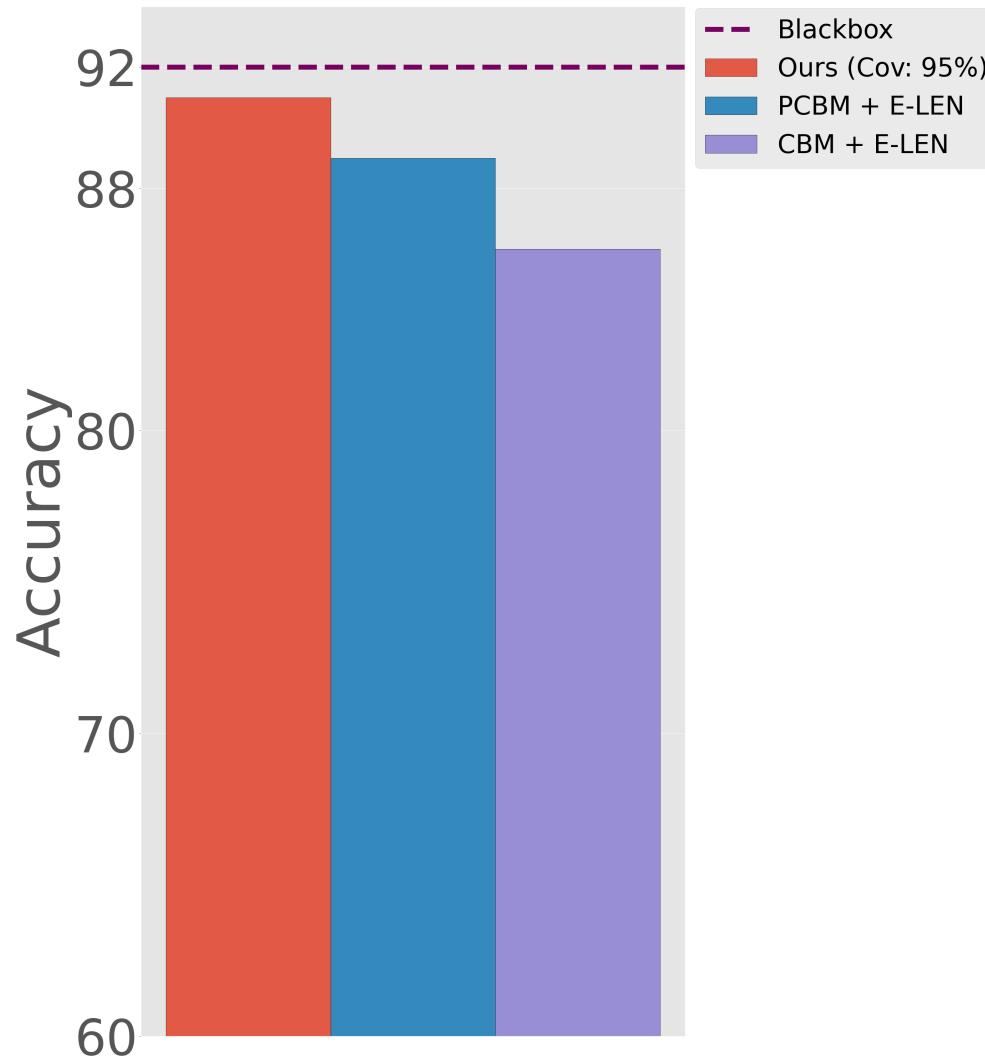


Final Residual
(Unexplained)

Comparing Performance

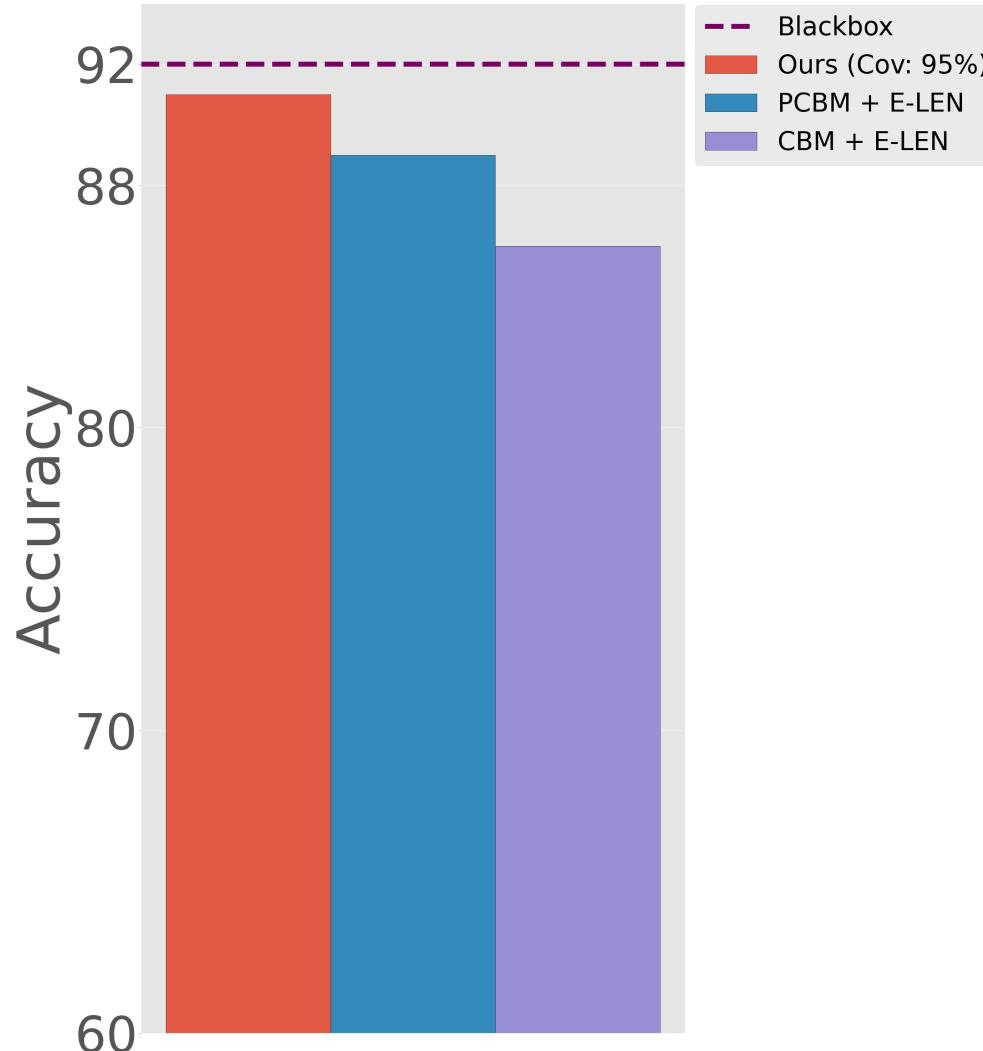
Comparing Performance

CUB-200 with ViT

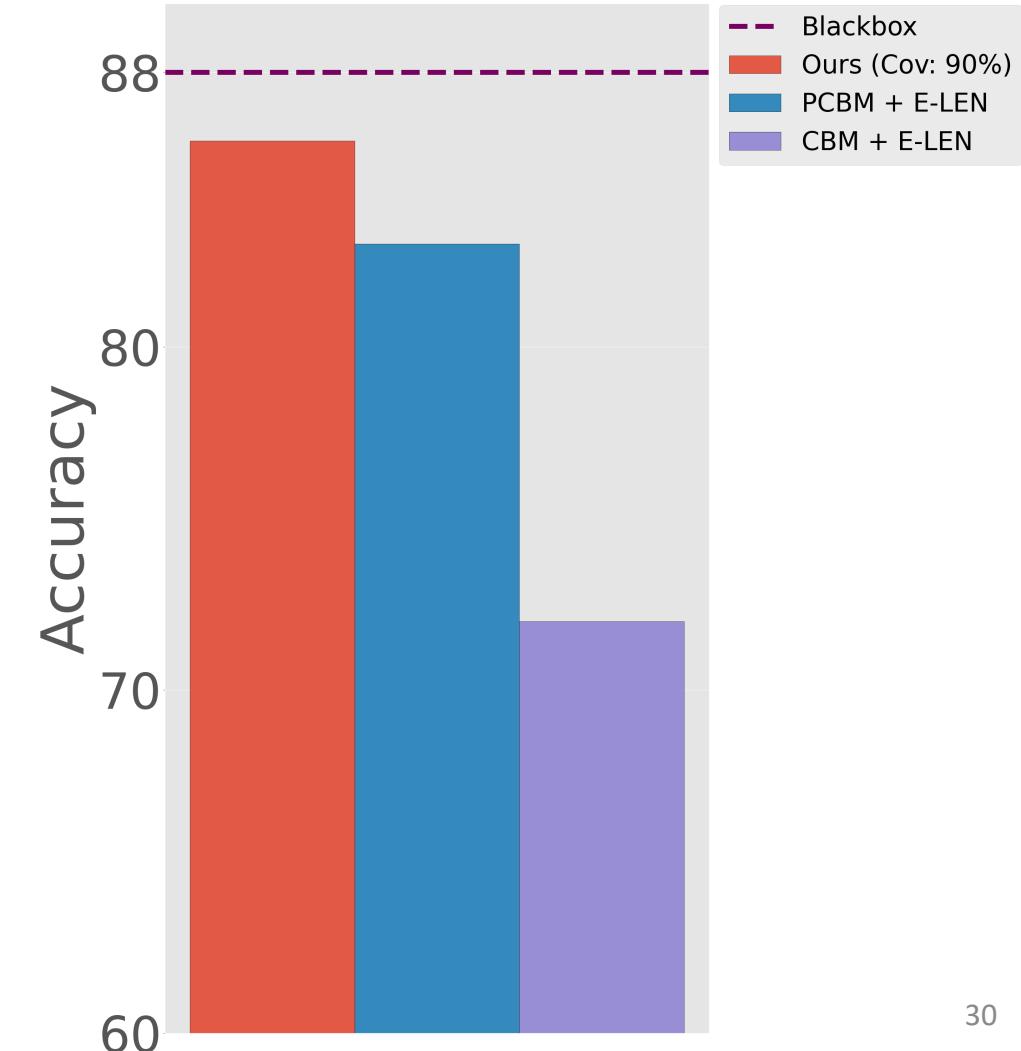


Comparing Performance

CUB-200 with ViT

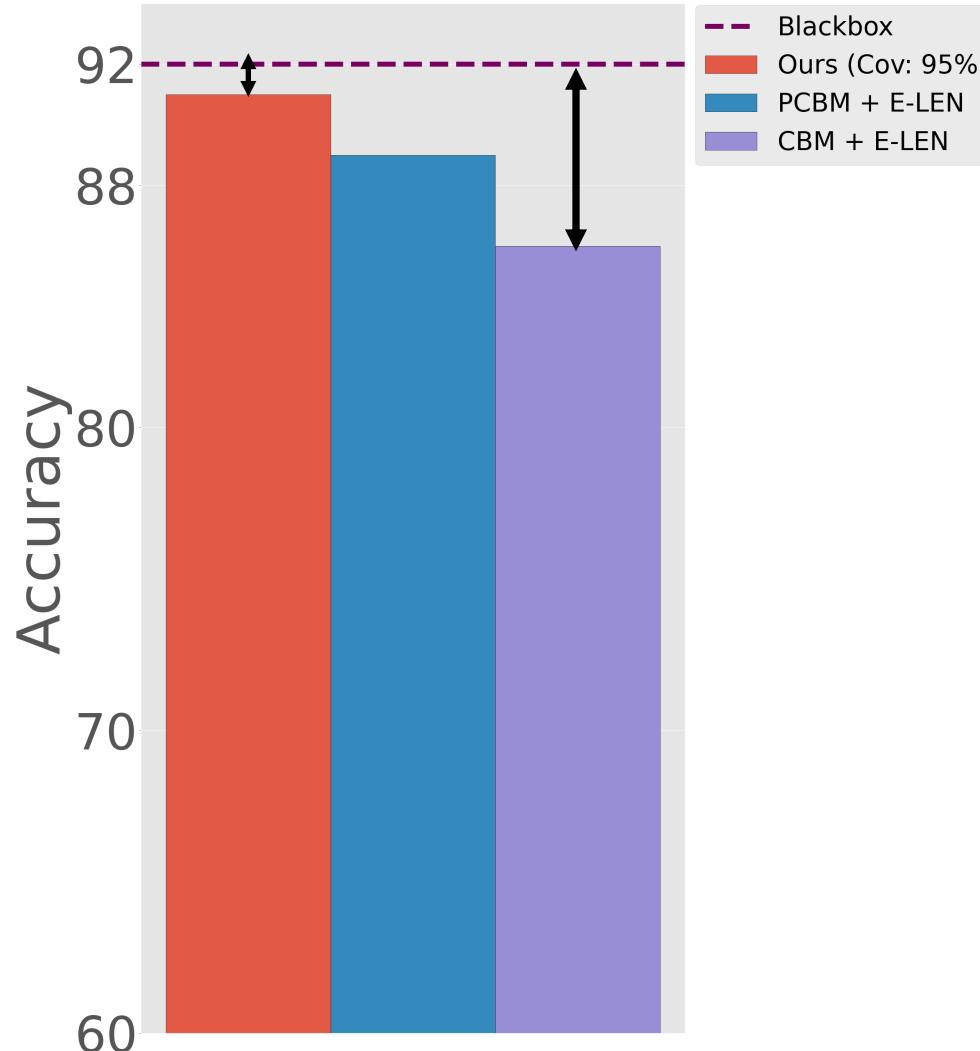


CUB-200 with ResNet101

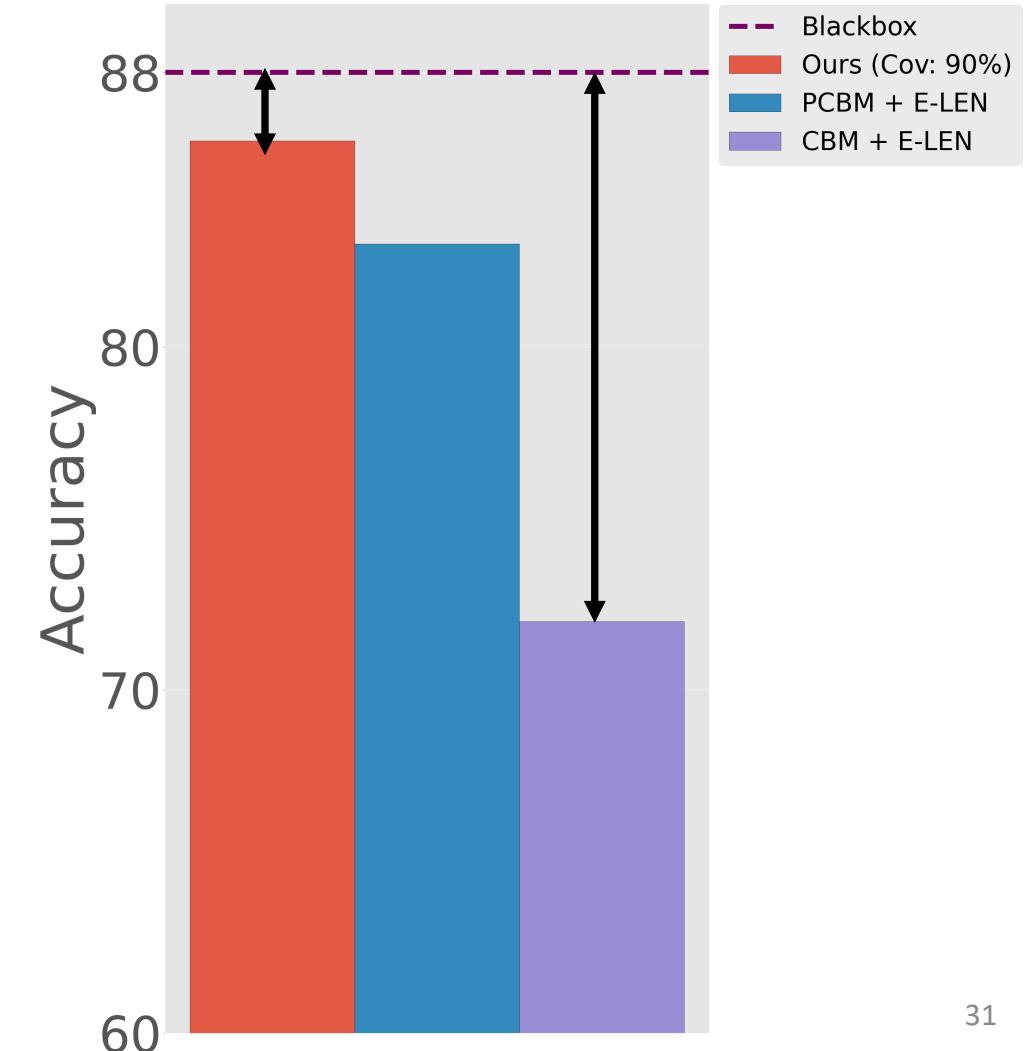


Comparing Performance

CUB-200 with ViT



CUB-200 with ResNet101



Residuals gradually performs poorly

Coverages of various experts for CUB-200 (ViT)

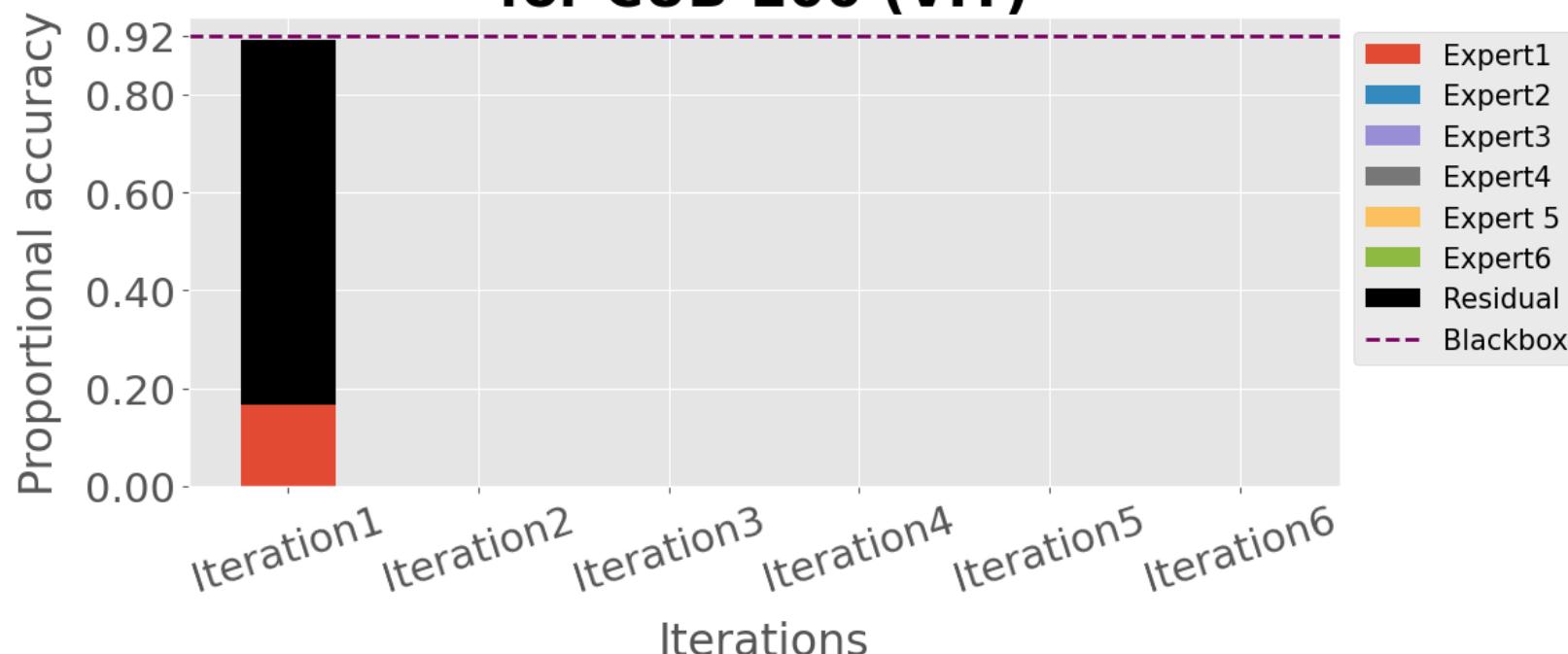


Residuals gradually performs poorly

Coverages of various experts for CUB-200 (ViT)



Performance across different iterations for CUB-200 (ViT)

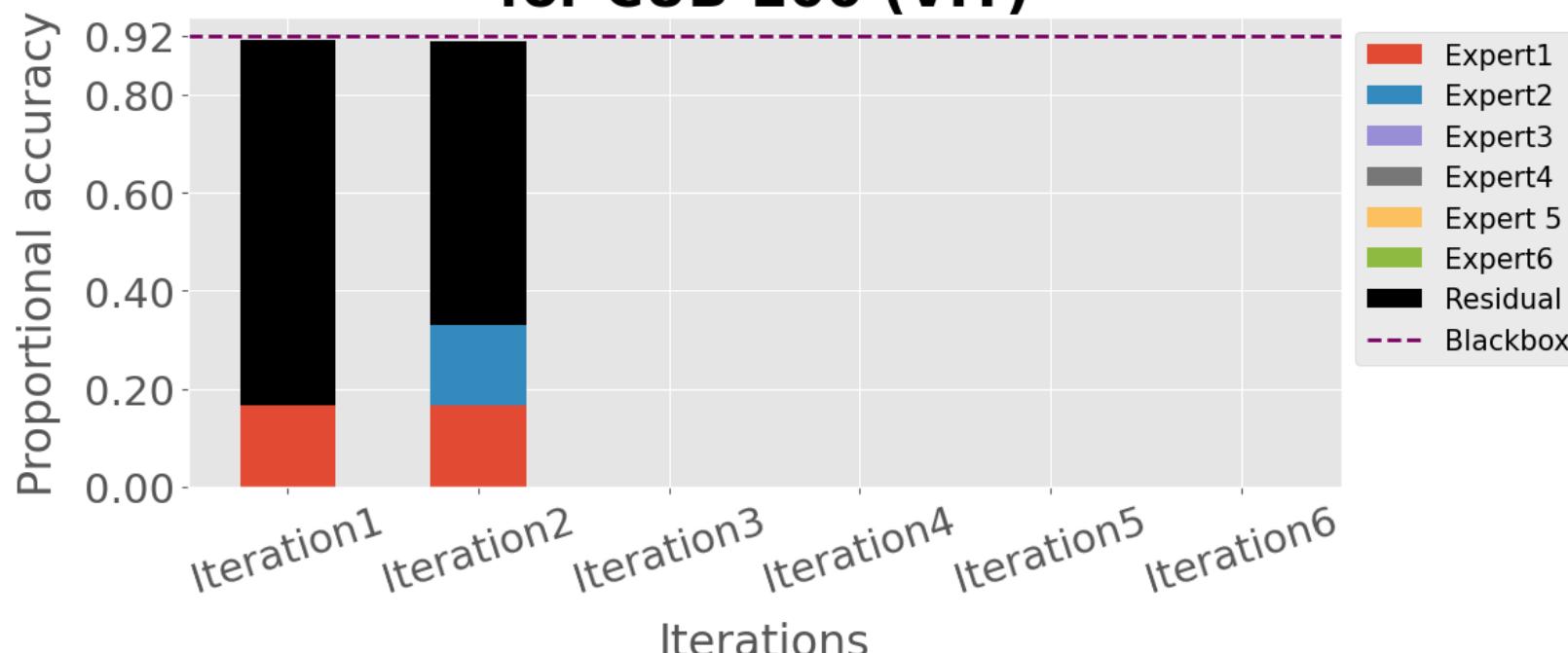


Residuals gradually performs poorly

Coverages of various experts for CUB-200 (ViT)



Performance across different iterations for CUB-200 (ViT)

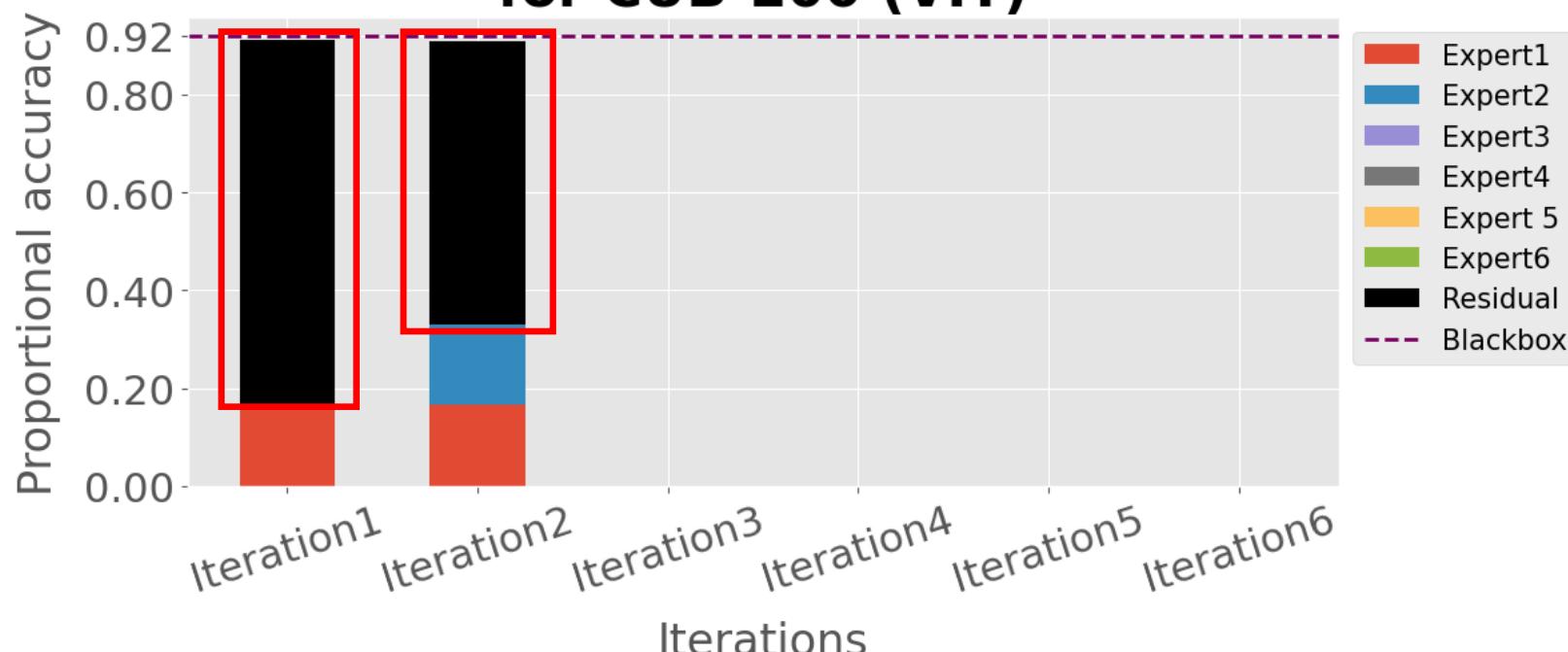


Residuals gradually performs poorly

Coverages of various experts for CUB-200 (ViT)



Performance across different iterations for CUB-200 (ViT)

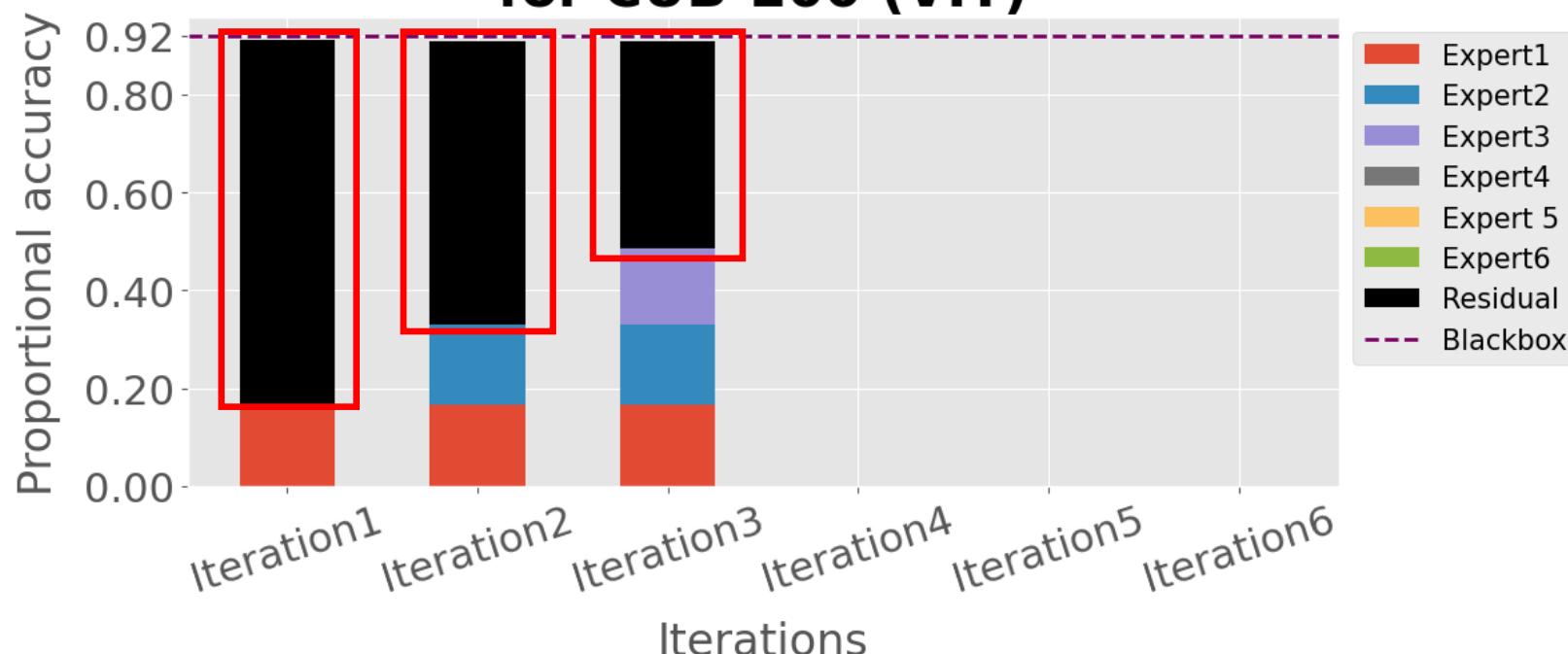


Residuals gradually performs poorly

Coverages of various experts for CUB-200 (ViT)



Performance across different iterations for CUB-200 (ViT)

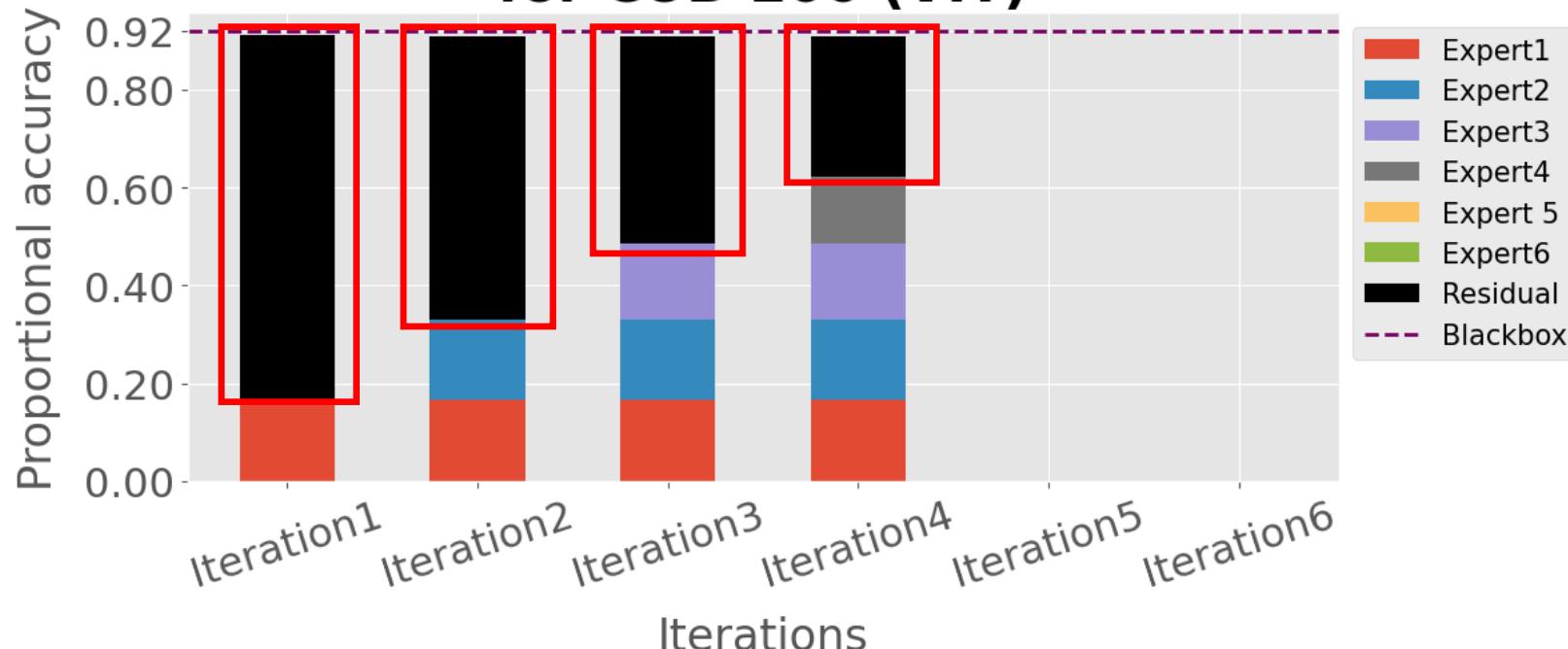


Residuals gradually performs poorly

Coverages of various experts for CUB-200 (ViT)



Performance across different iterations for CUB-200 (ViT)

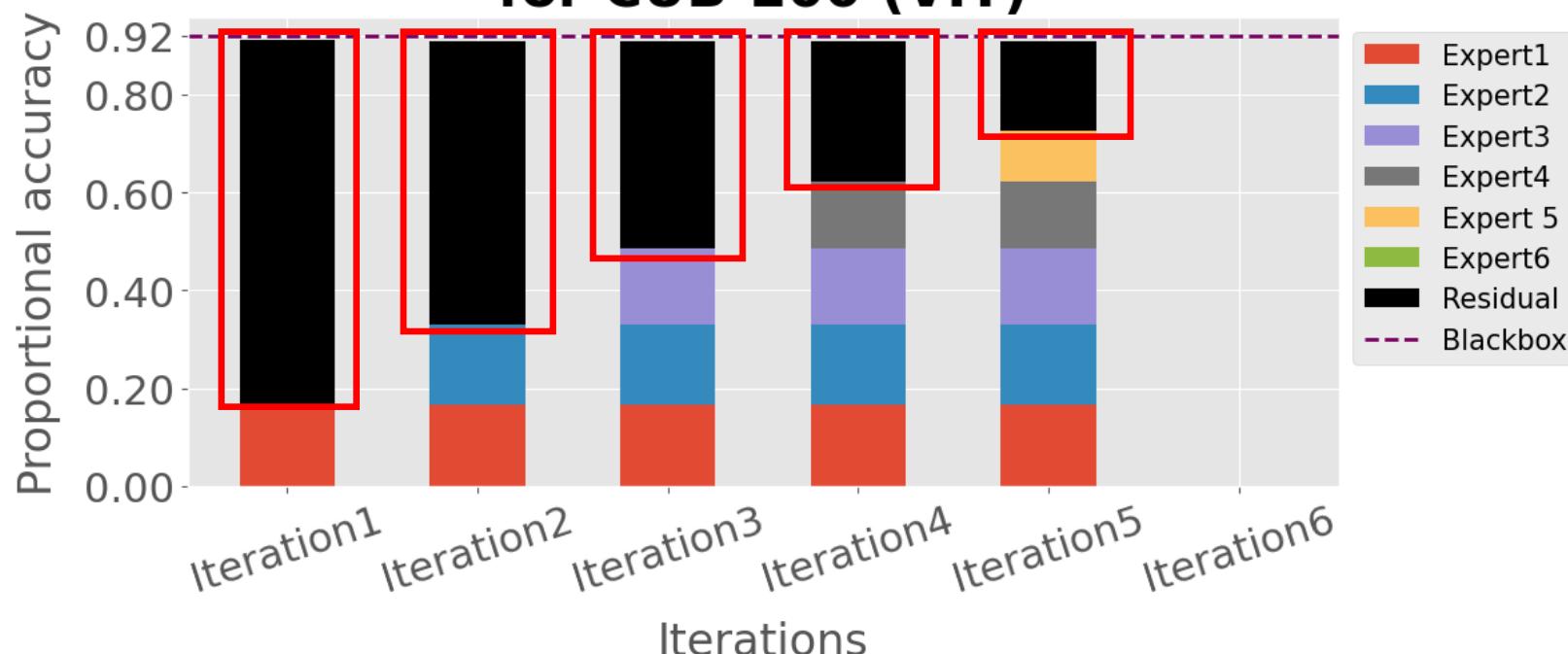


Residuals gradually performs poorly

Coverages of various experts for CUB-200 (ViT)



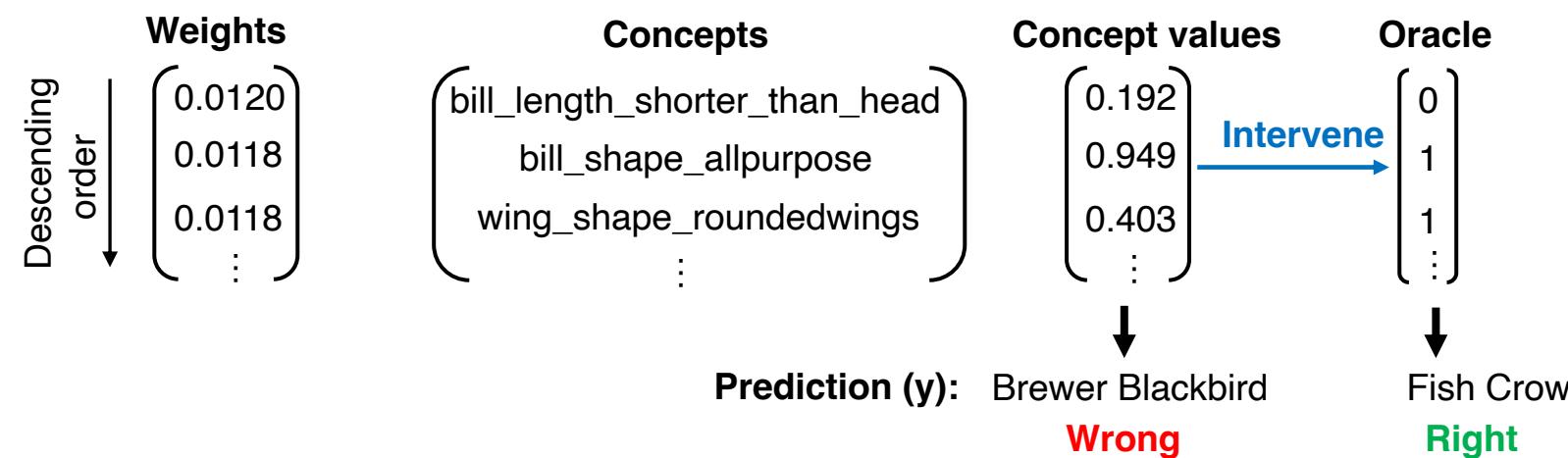
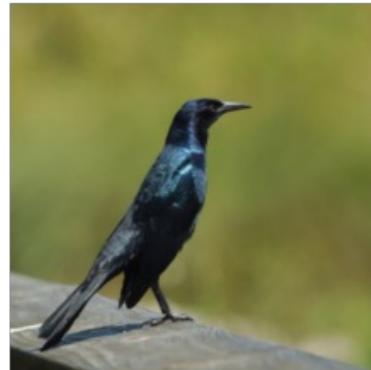
Performance across different iterations for CUB-200 (ViT)



Test time interventions on hard samples

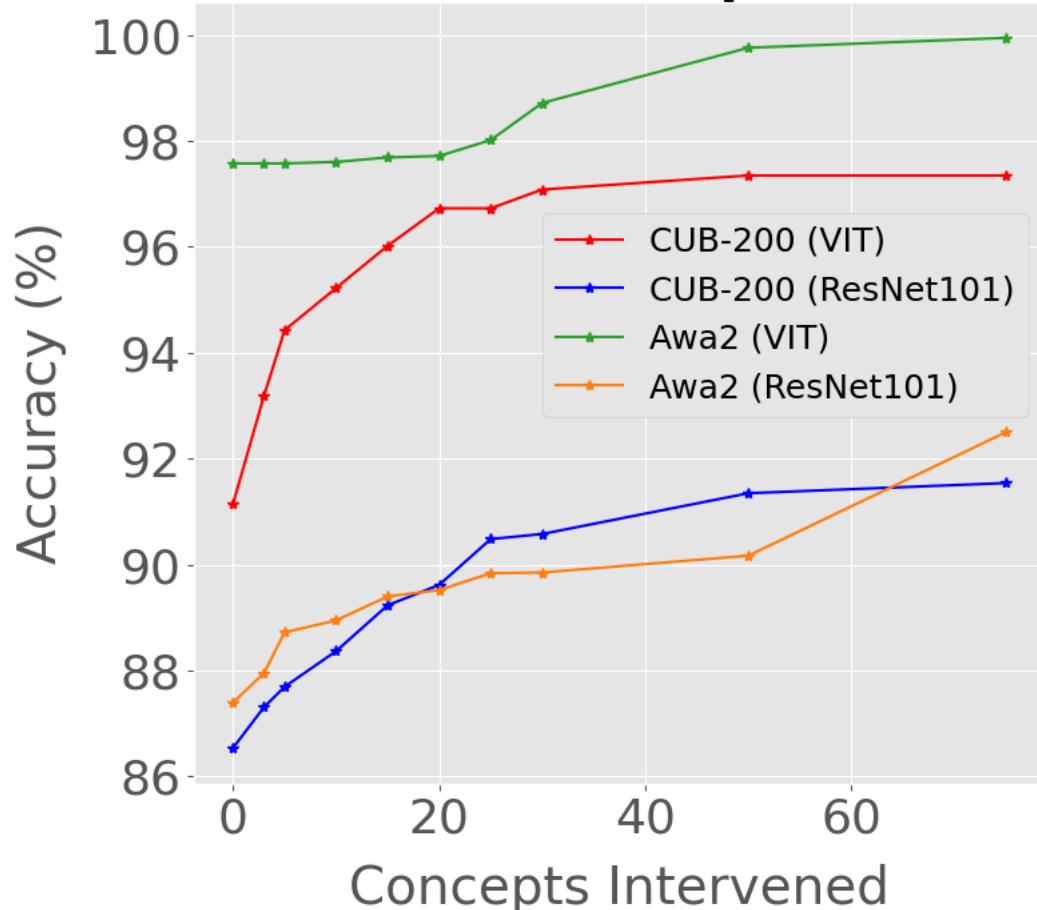
Chose the misclassified samples covered by the last 2 experts

Expert6



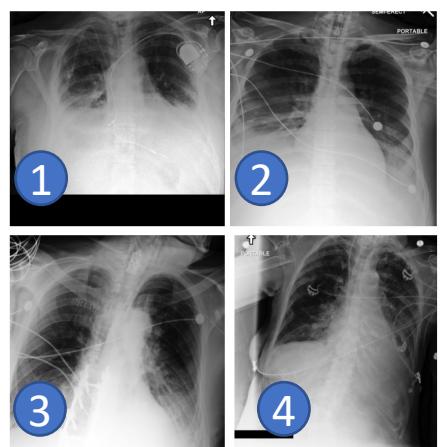
Test time interventions on hard samples

Test time interventions for the last two experts

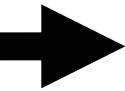


Intervening on top 75 concepts

- CUB-200 (ResNet101) – 5.49 % ↑
- CUB-200 (VIT) – 6.5 % ↑
- AWA2 (ResNet101) – 5.74 % ↑
- AWA2 (VIT) – 2.42 % ↑

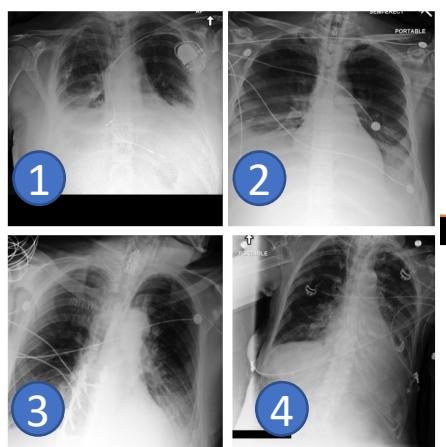


Examples on Chest X-ray



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

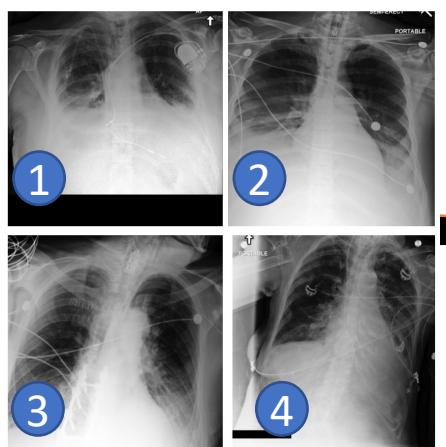
Expert 1

Effusion ↔
left_pleural
^ right_pleural
^ pleural_unspec

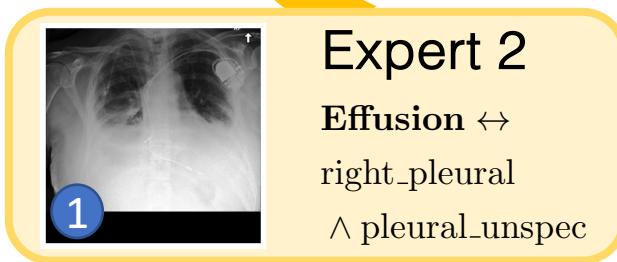
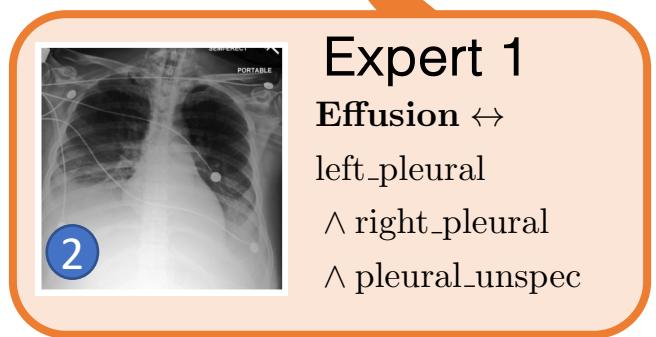


Pleural_unspec is “unspecified pleural effusion” referred to as “hydrothorax”.

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



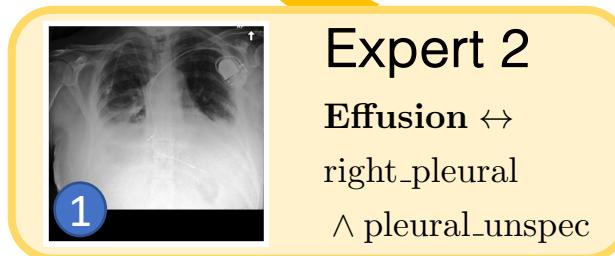
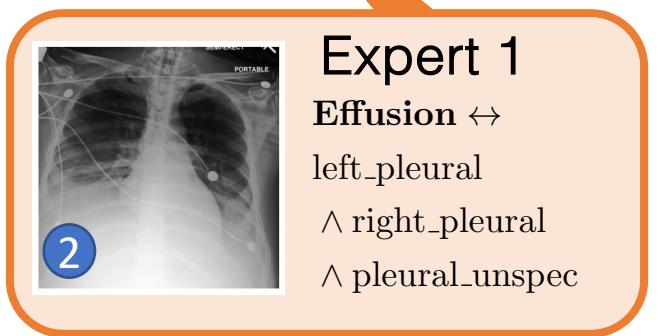
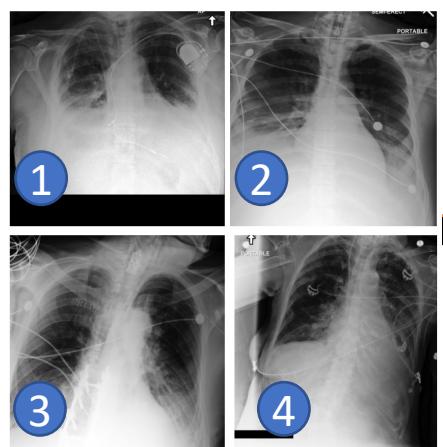
Examples on Chest X-ray



Pleural_unspec is “unspecified pleural effusion” referred to as “hydrothorax”.

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

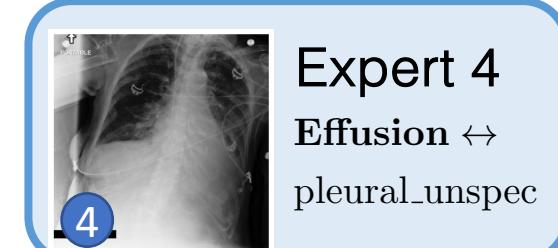
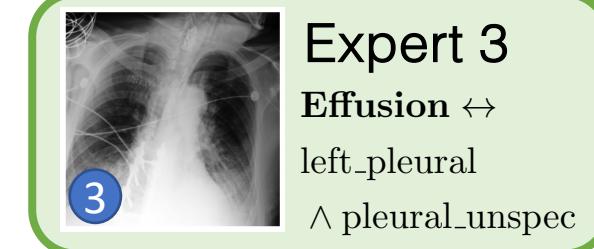
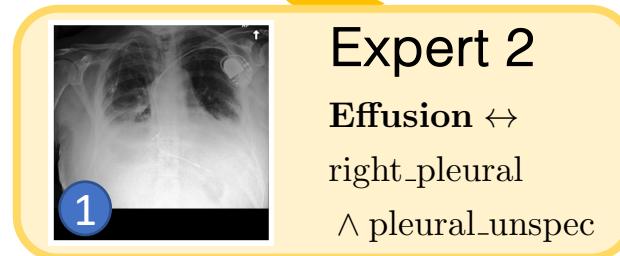
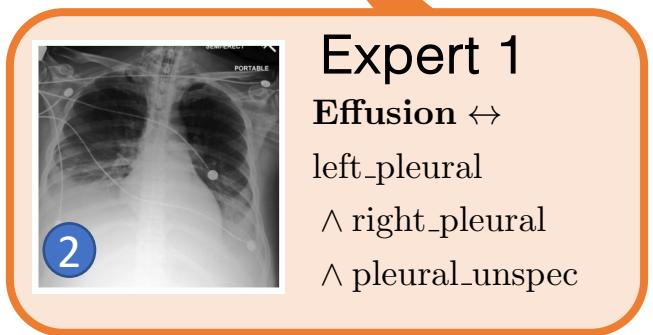
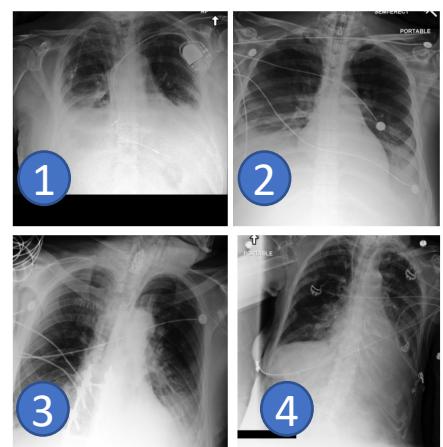
Examples on Chest X-ray



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

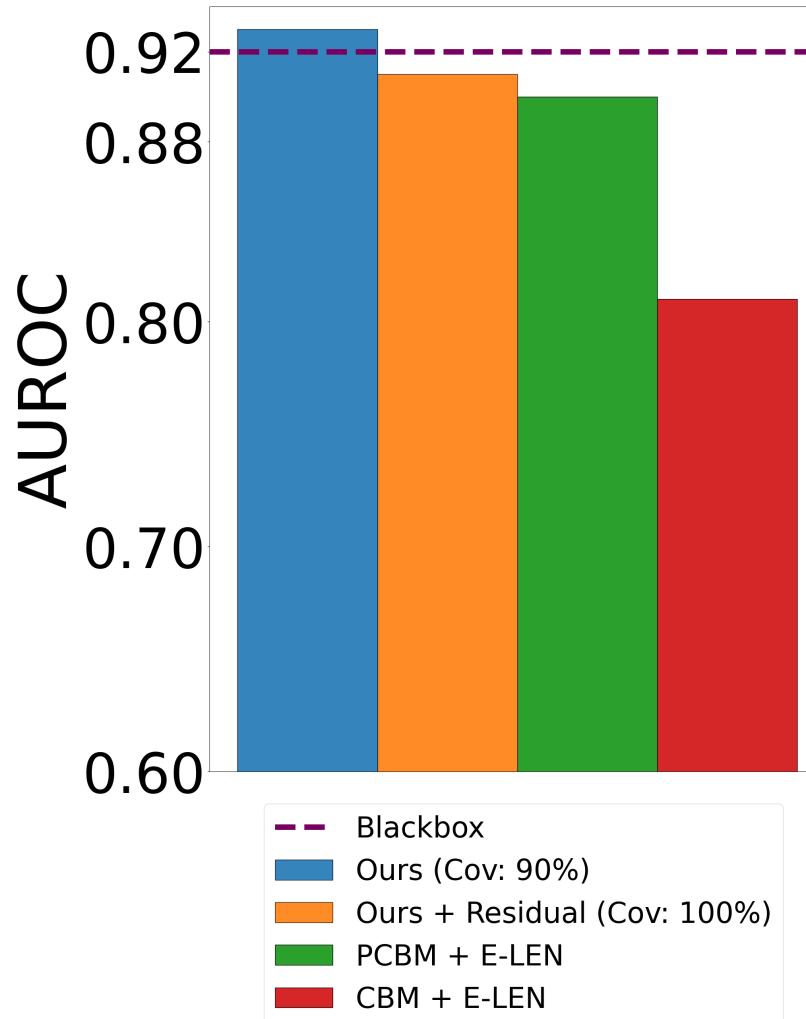


Pleural_unspec is “unspecified pleural effusion” referred to as “hydrothorax”.

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

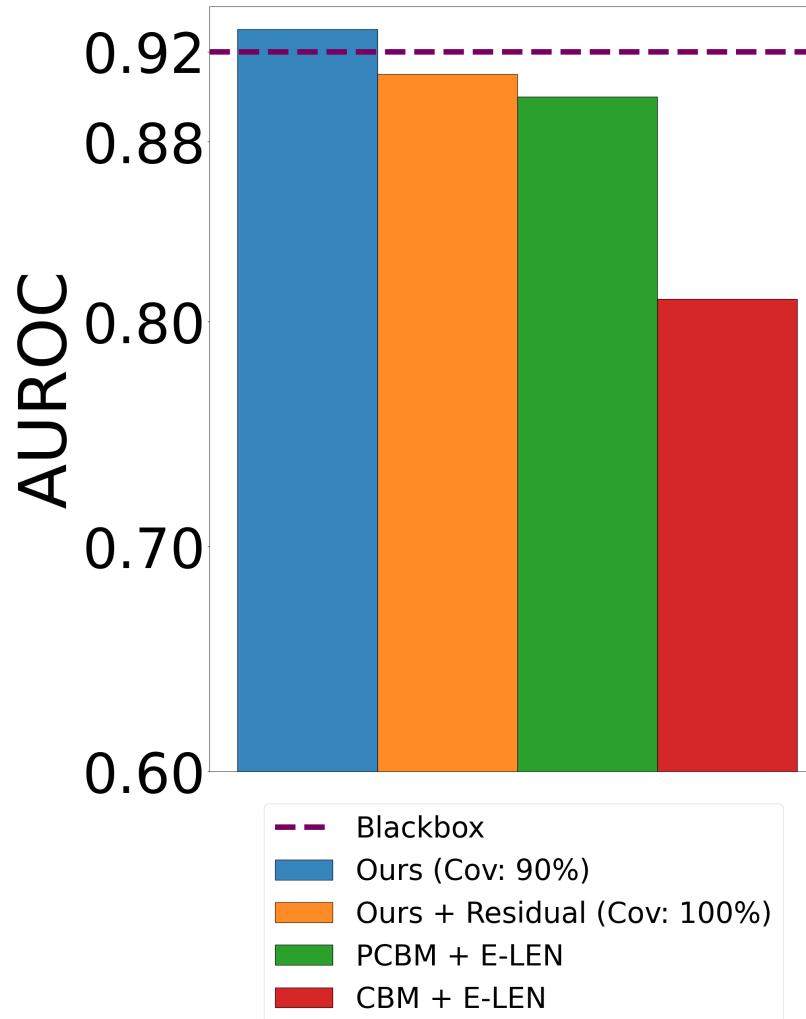
Comparing Performance

EFFUSION

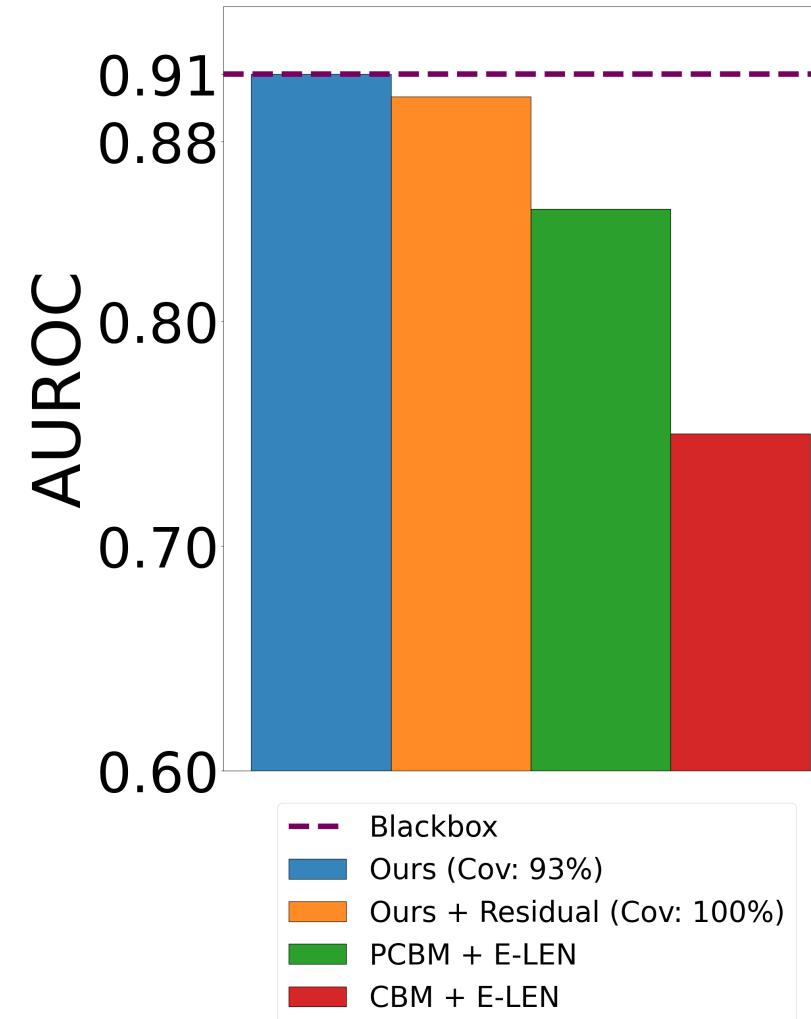


Comparing Performance

EFFUSION

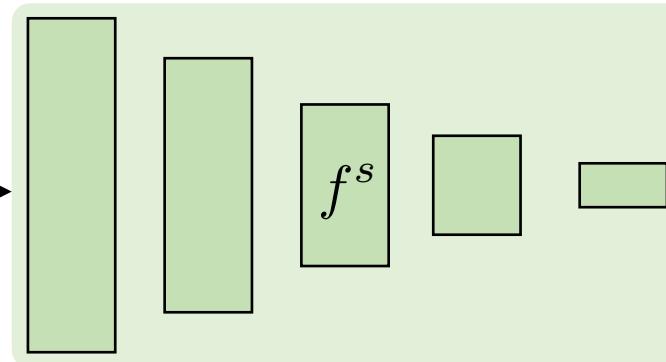


PNEUMOTHORAX



Data-Efficient Fine-tuning

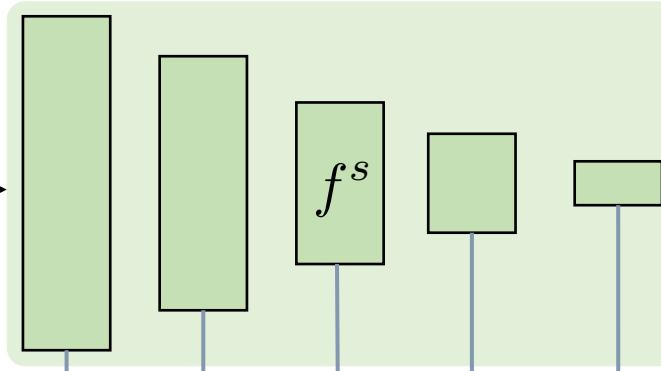
MIMIC-CXR



Pneumothorax

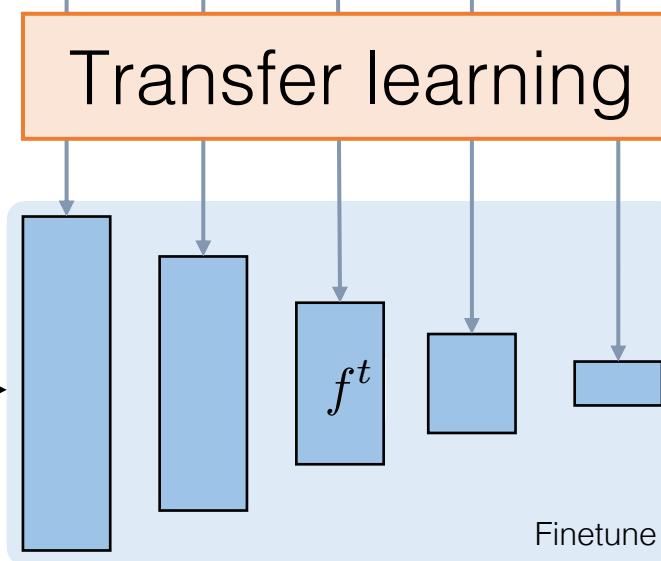
Data-Efficient Fine-tuning

MIMIC-CXR



Pneumothorax

Stanford-CXR



Pneumothorax

Data-Efficient Fine-tuning

MIMIC-CXR



Pneumothorax

Data and Computationally inefficient

Stanfo



Pneumothorax

The clinical rules are “invariant”

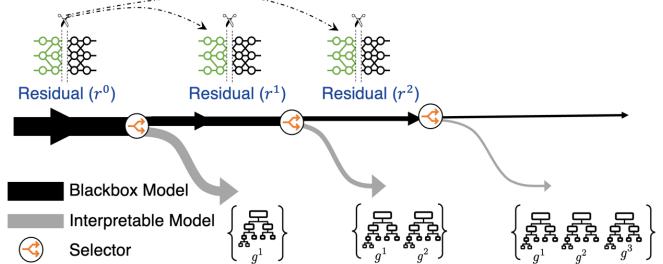


Finetune

Fine-tune to a New Domain

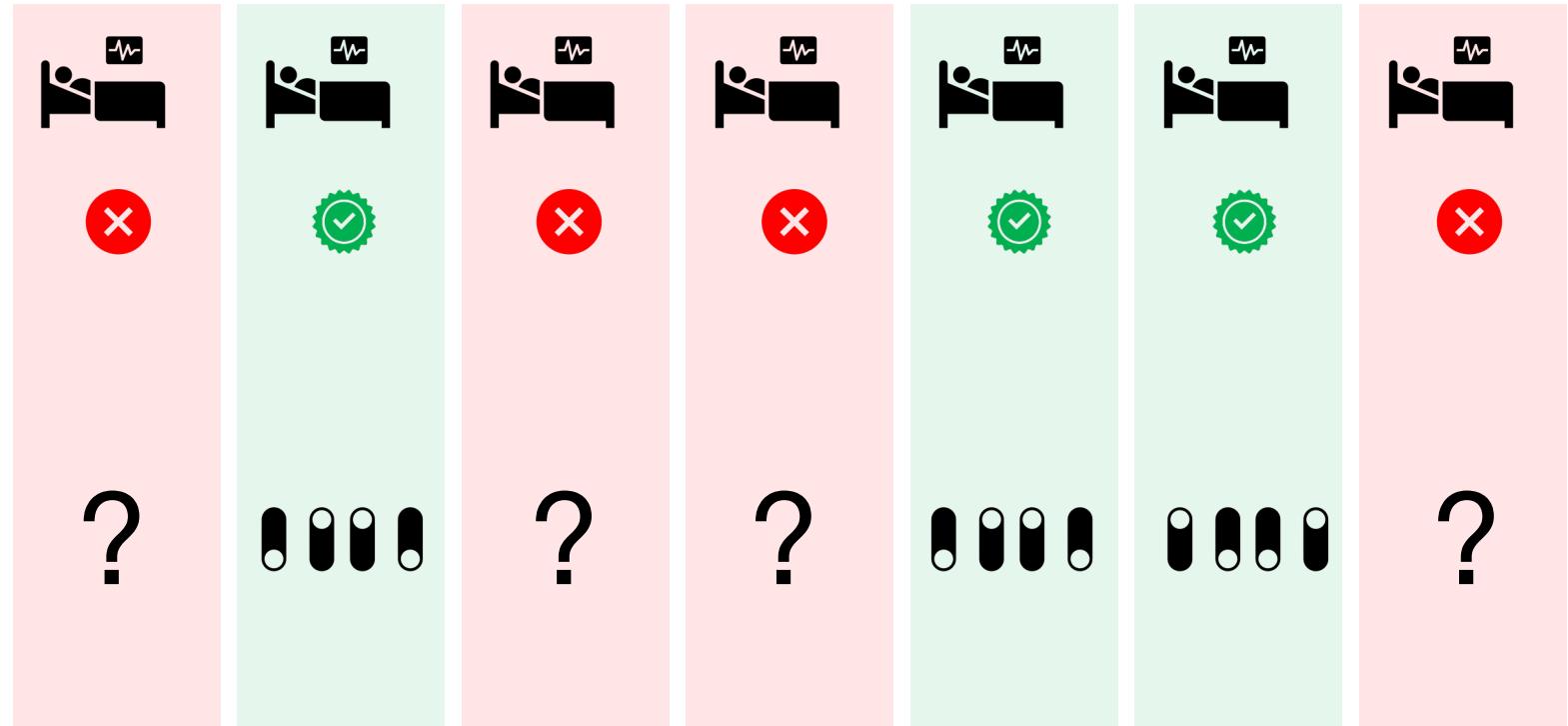
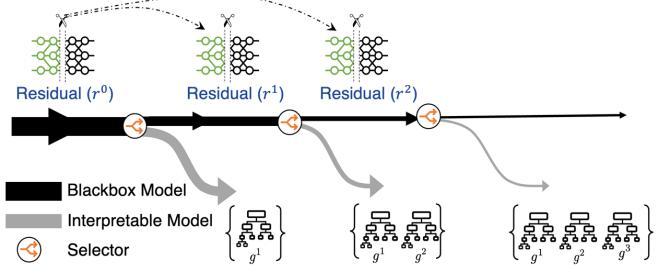
1

Apply source model



Fine-tune to a New Domain

1 Apply source model

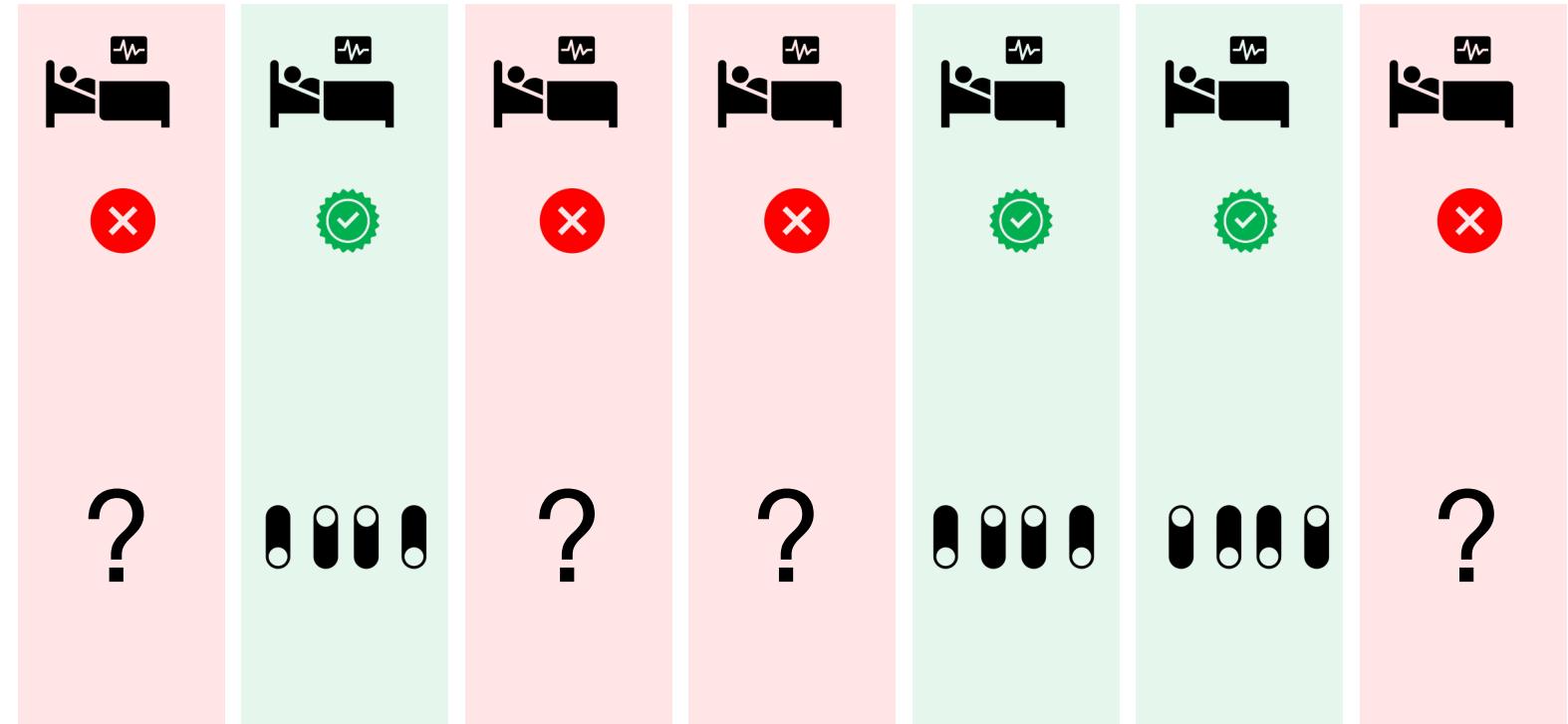
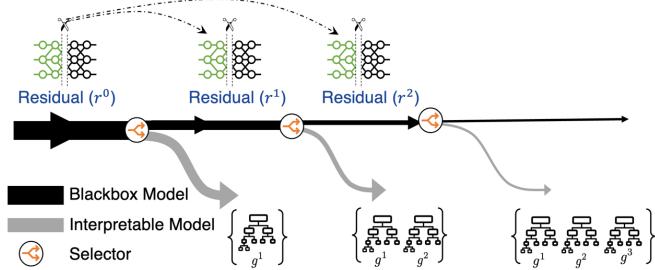


2 Use concepts from matching patients

C

Fine-tune to a New Domain

1 Apply source model



2 Use concepts from matching patients

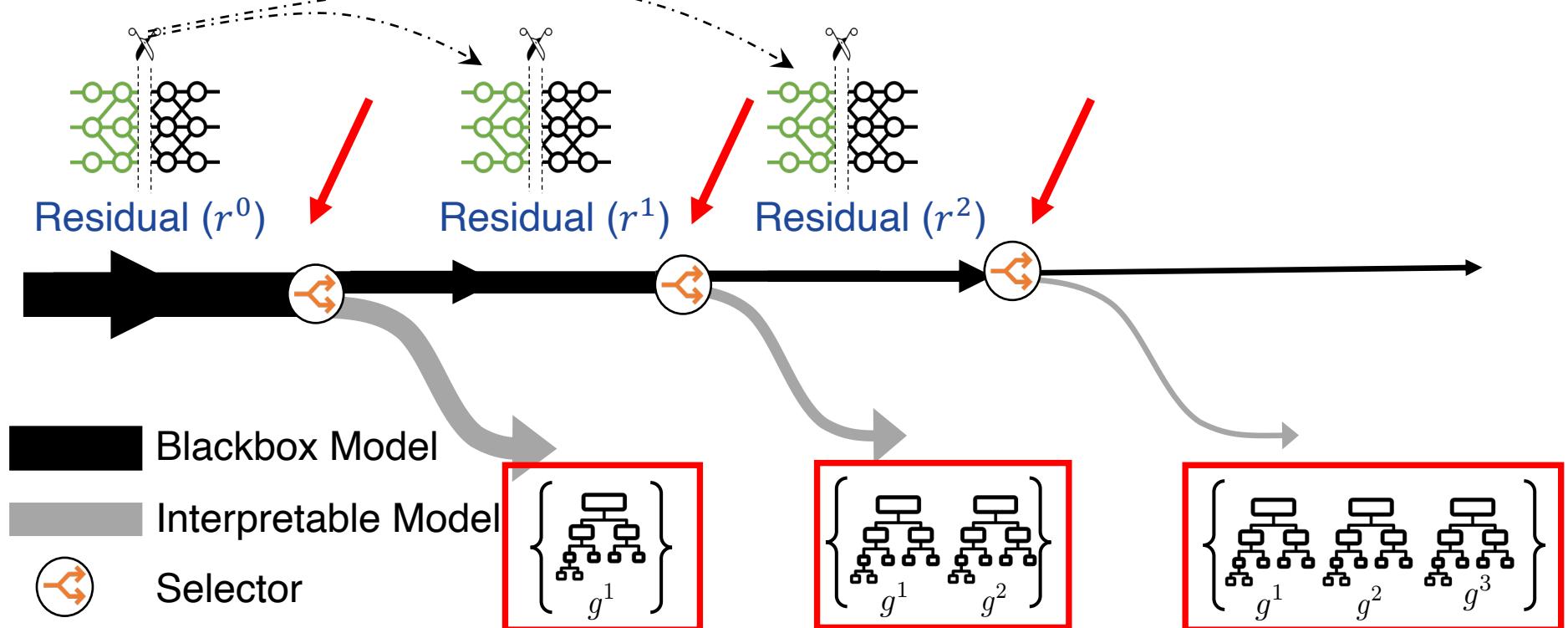
C

3 Propagate the concepts and update the concept extractor

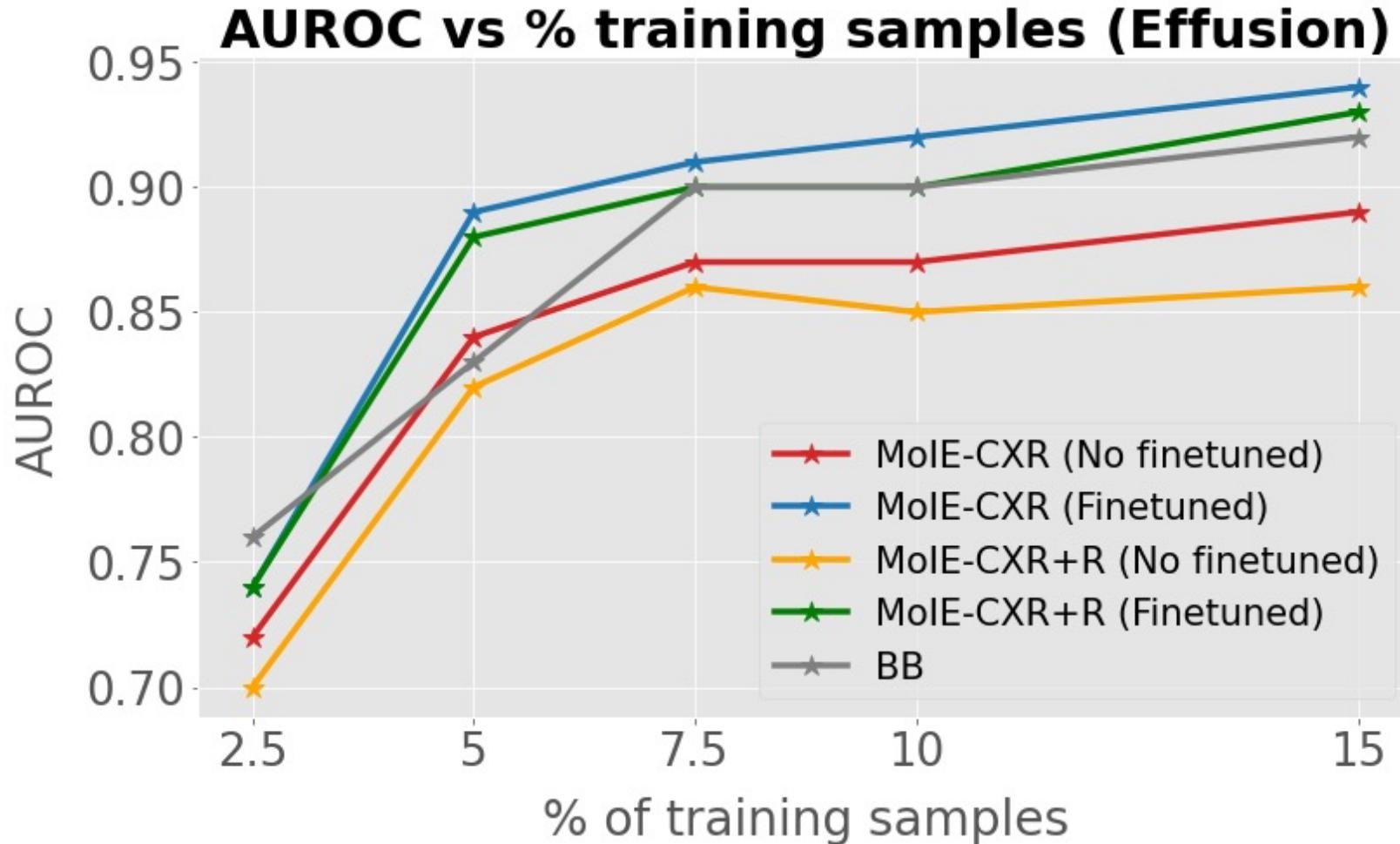
Fine-tune to a New Domain

4

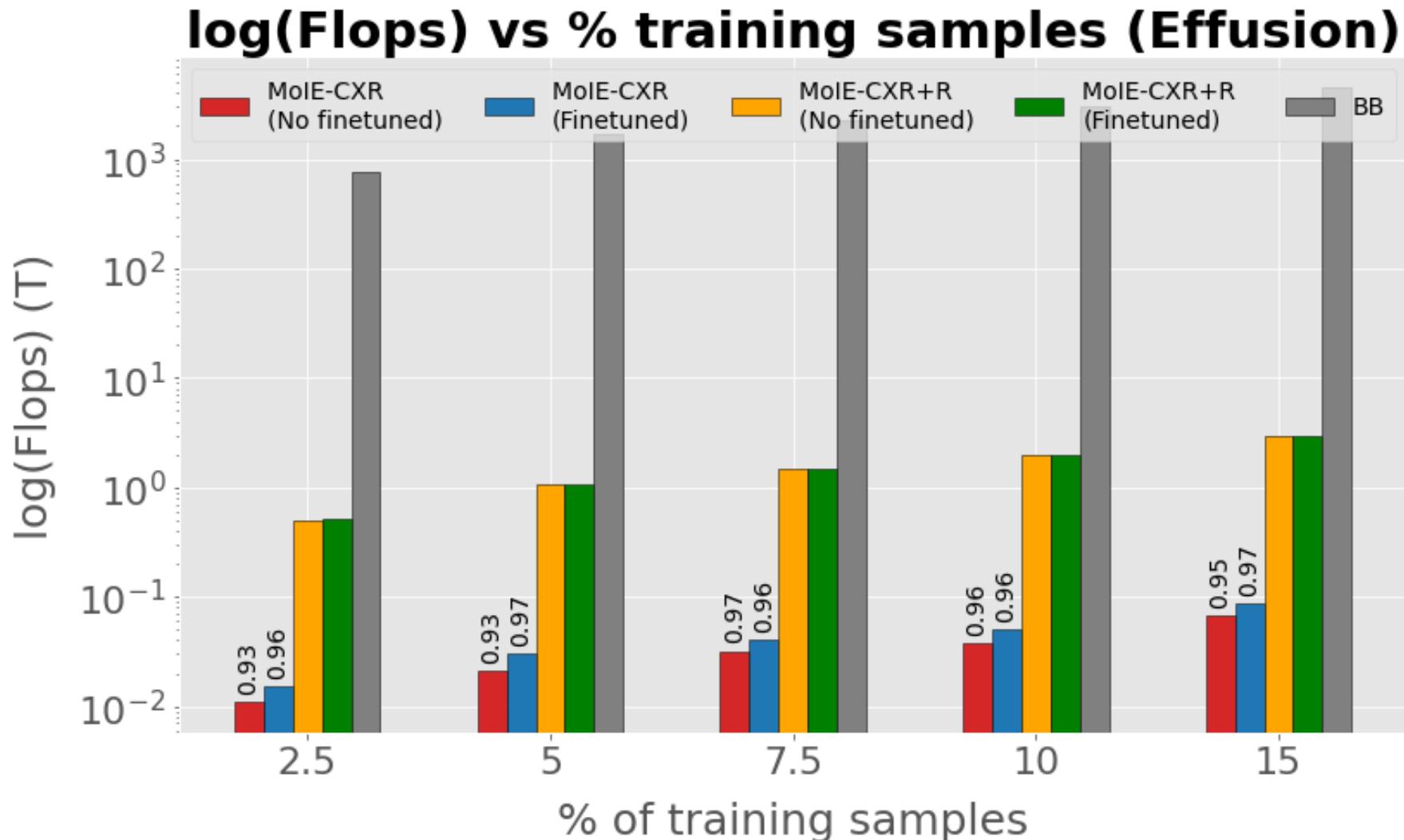
Update the selectors and interpretable models for 5 epochs



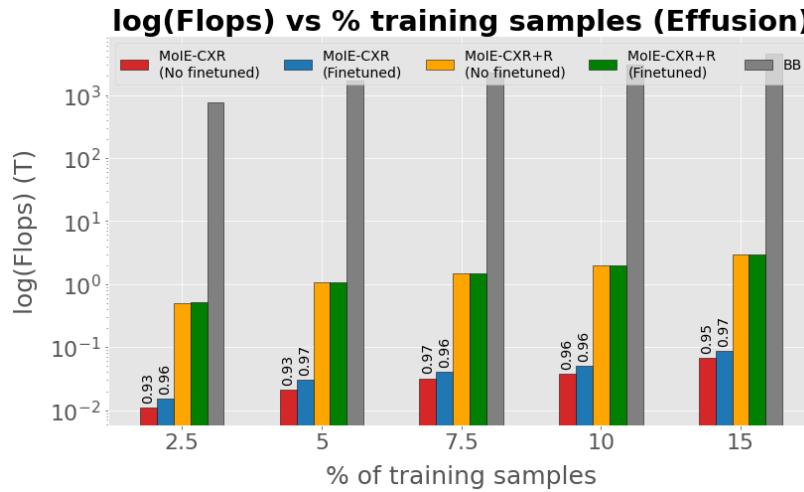
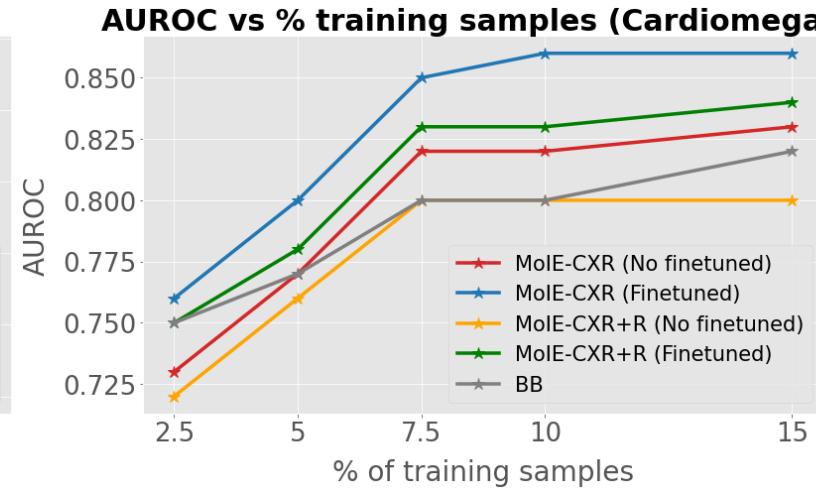
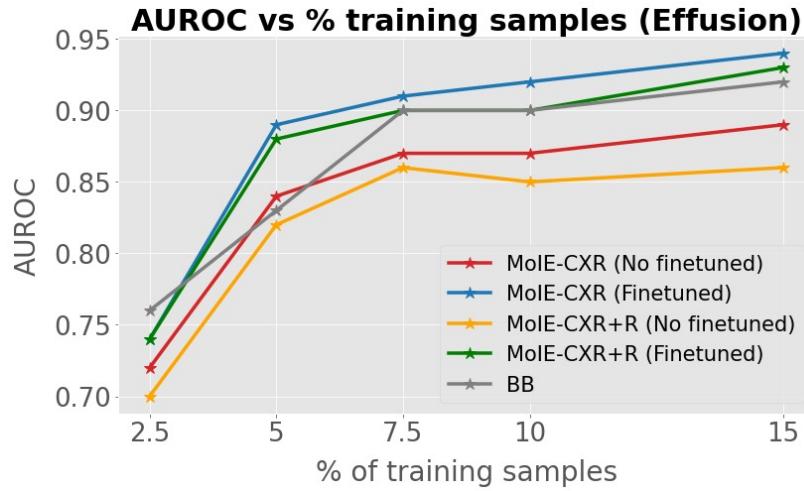
Transferring to Stanford-CXR



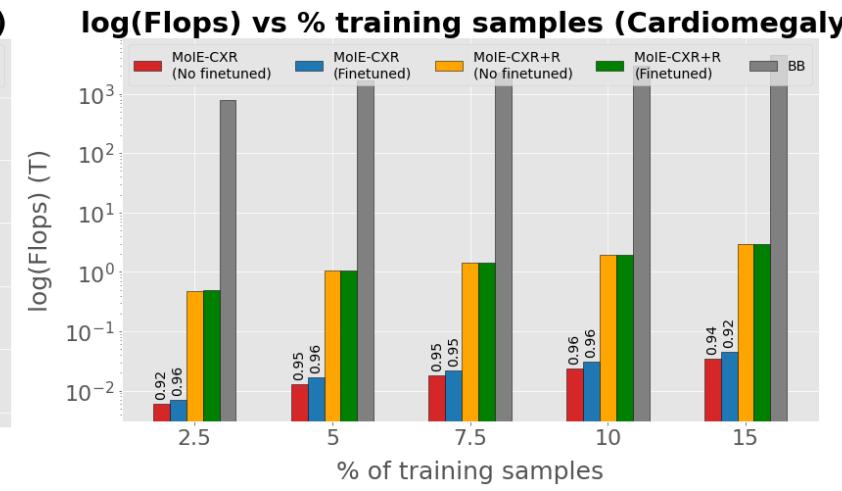
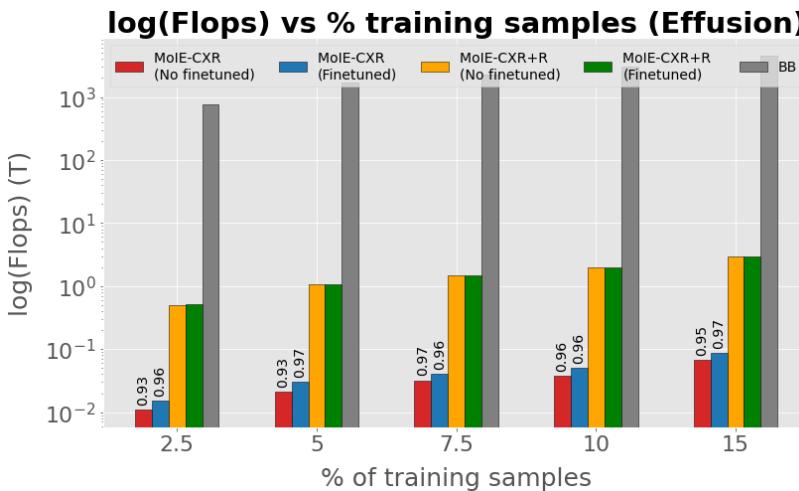
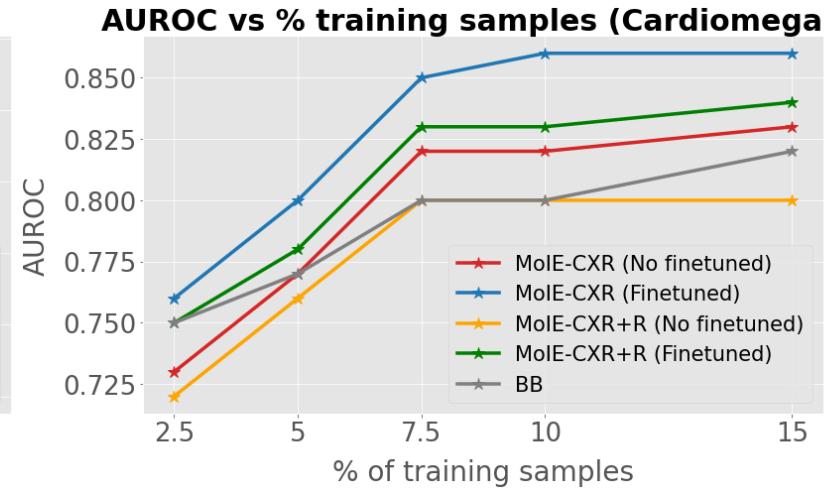
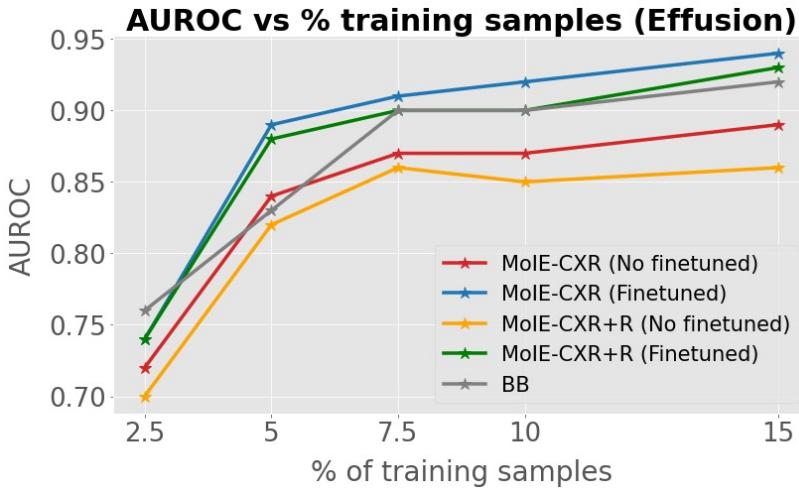
Transferring to Stanford-CXR



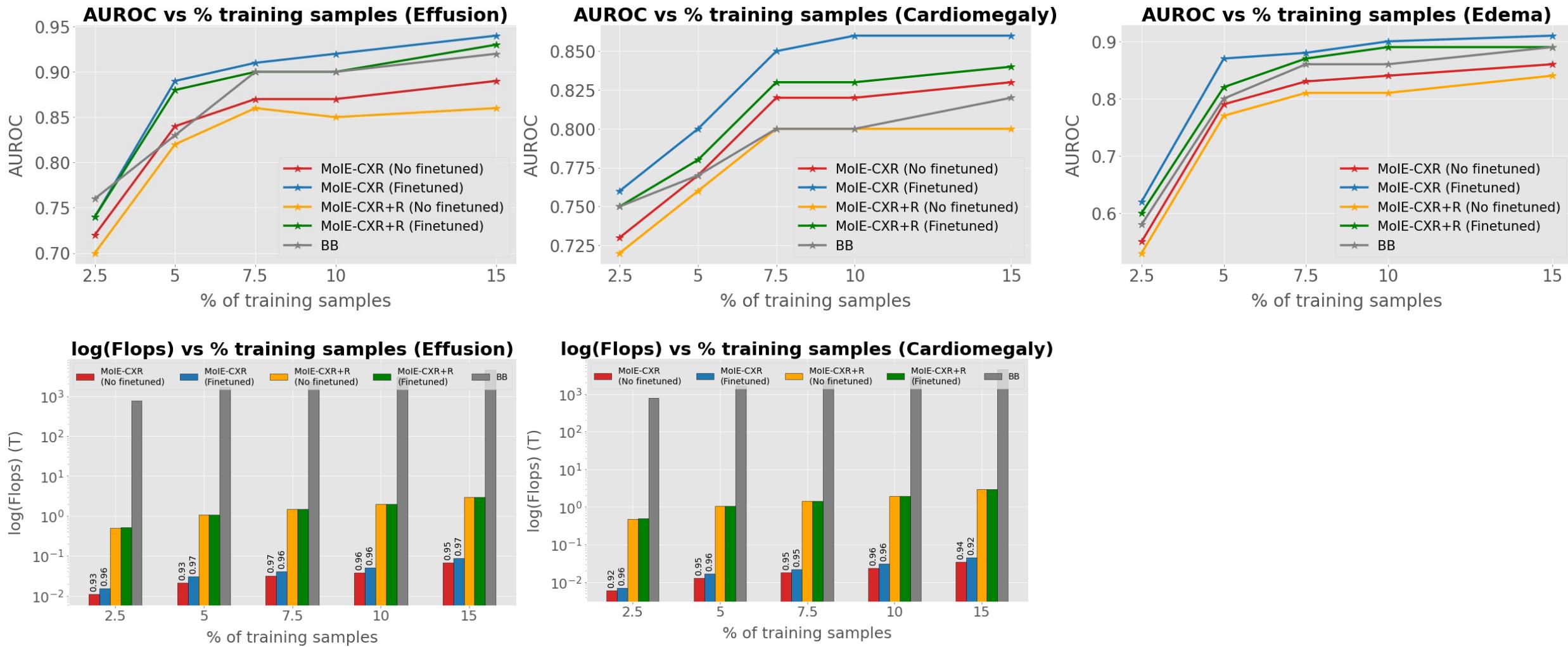
Transferring to Stanford-CXR



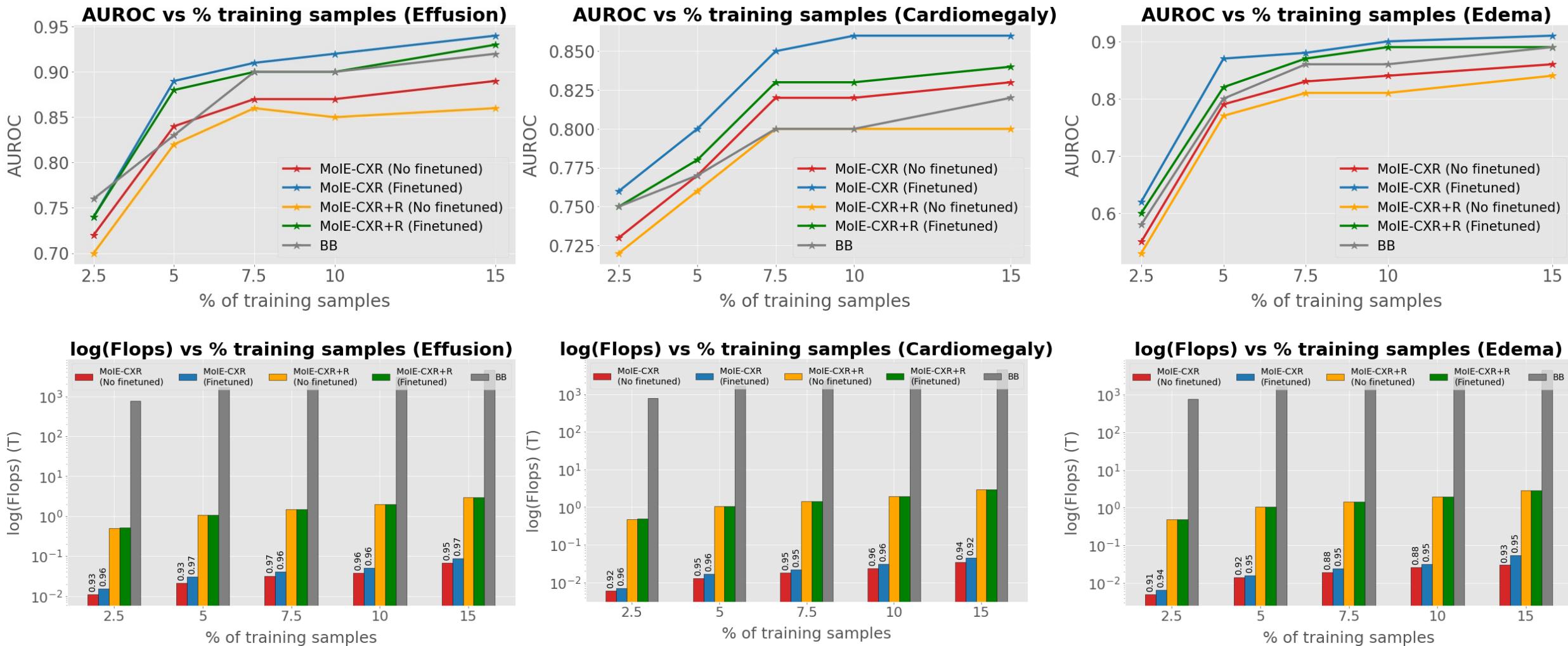
Transferring to Stanford-CXR



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**

Future work

- We aim to use language to detect and fix the blackbox.
- We are interested in Breast imaging, i.e 2D mammograms
- We developed the first vision language model for mammograms, i.e Mammo-CLIP (Early accept, top 11%, MICCAI 2024)
- We will present our work at MICCAI 2024

Date: Poster Session 5, Wednesday, October 9, 2024, 10:30 to 11:30

Poster number: W-AM-169

Project: <https://shantanu-ai.github.io/projects/MICCAI-2024-Mammo-CLIP/>



Shantanu Ghosh¹, Ke Yu², Kayhan Batmanghelich¹

¹BU ECE, ²Pitt ISP



Thank you

