



# Bridging the Gap: From Post Hoc Explanations to Inherently Interpretable Models for Medical Imaging



Shantanu Ghosh<sup>1</sup>, Ke Yu<sup>2</sup>, Forough Arabshahi<sup>3</sup>, Kayhan Batmanghelich<sup>1</sup>

<sup>1</sup>Dept. Of Electrical and Computer Engineering, Boston University

<sup>2</sup>Intelligent Systems Program (ISP), University of Pittsburgh, <sup>3</sup>Meta AI



**TLDR:** Extracting a mixture of interpretable models from a BlackBox to provide instance specific concept-based explanations using First-order logic (FOL).

## Post hoc explanation

### Pros

- Does not alter the Black box.

### Cons

- Inconsistent explanations.
- No recourse.

## Interpretable by design

### Pros

- Support concept intervention.

### Cons

- Harder to train.
- Sub par performance.

How to blur this gap?

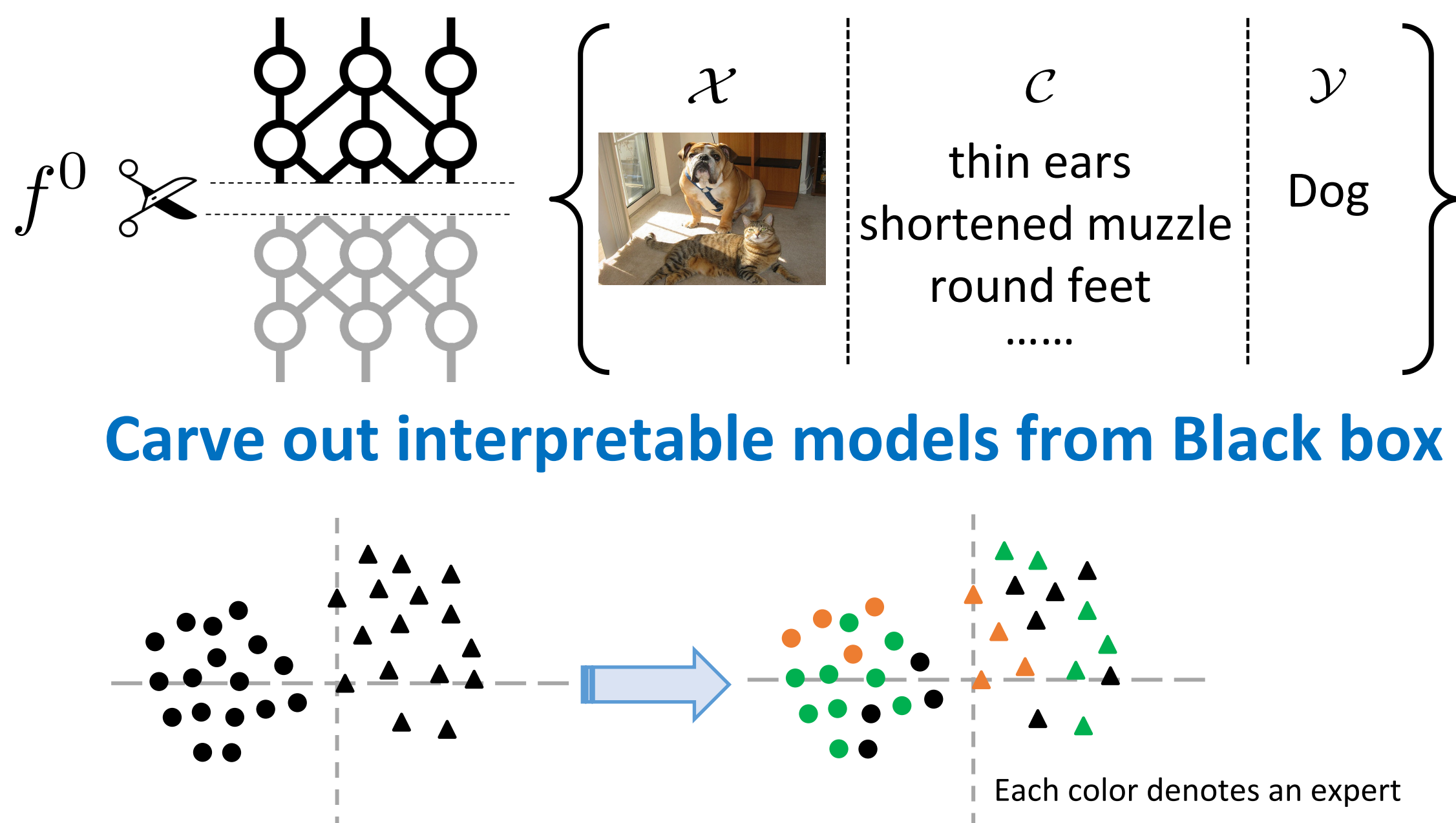
## Desirable properties

- Does compromise the performance.
- Can be intervened to fix the misclassification

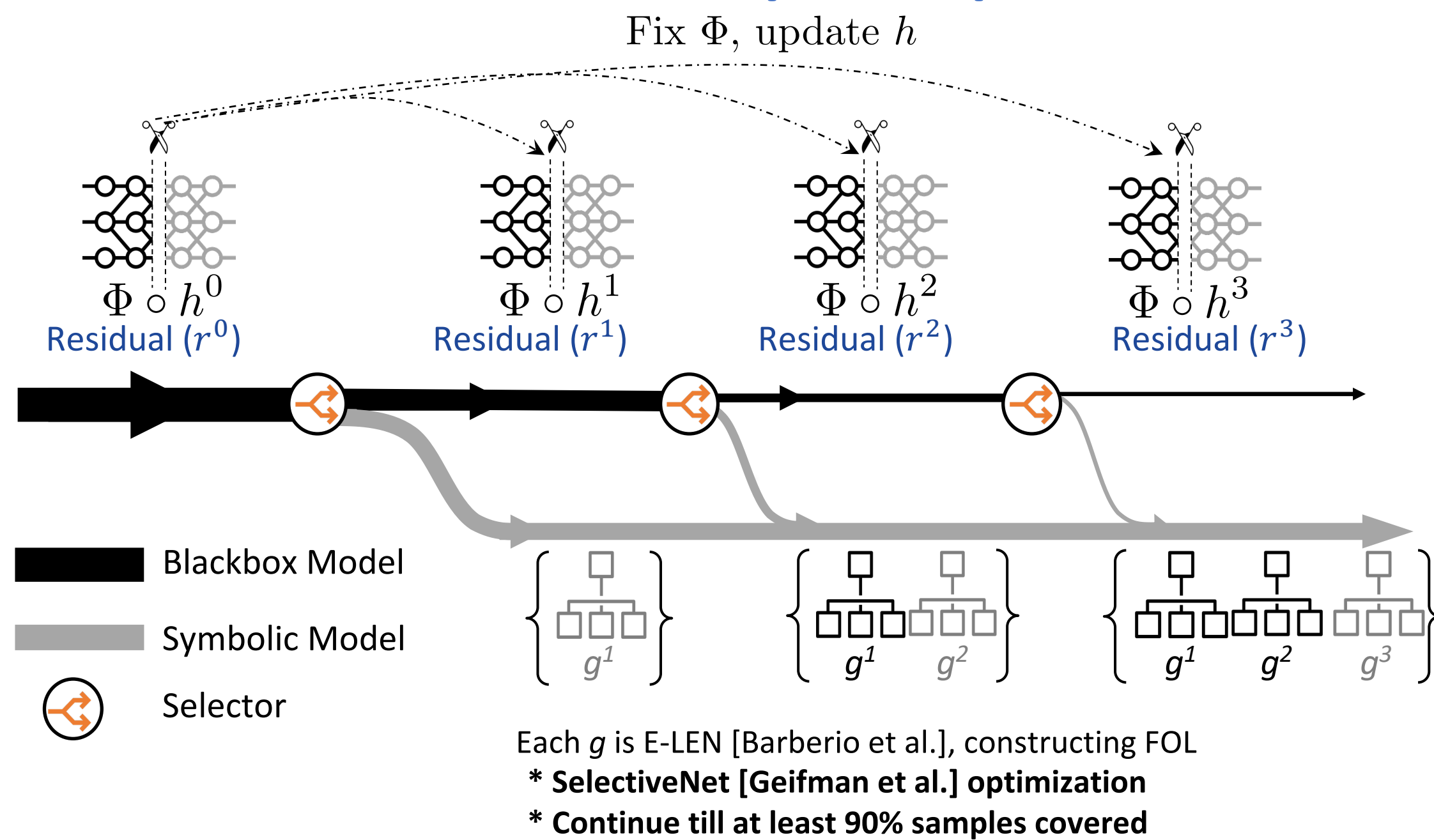
## Design choices

- Carve interpretable models from Blackbox.
- Concept based
- First order logic for concept interaction

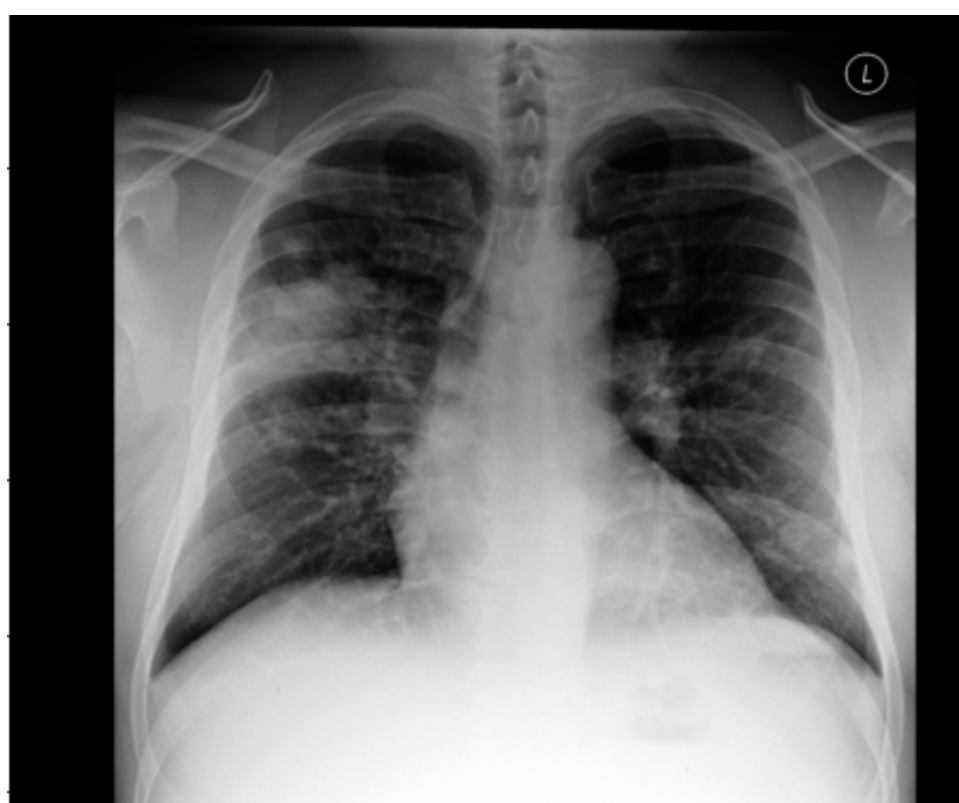
## Assumption



## Route Interpret Repeat



## Extract concepts from MIMIC-CXR using Radgraph NLP pipeline



Ke Yu et al.

### Report:

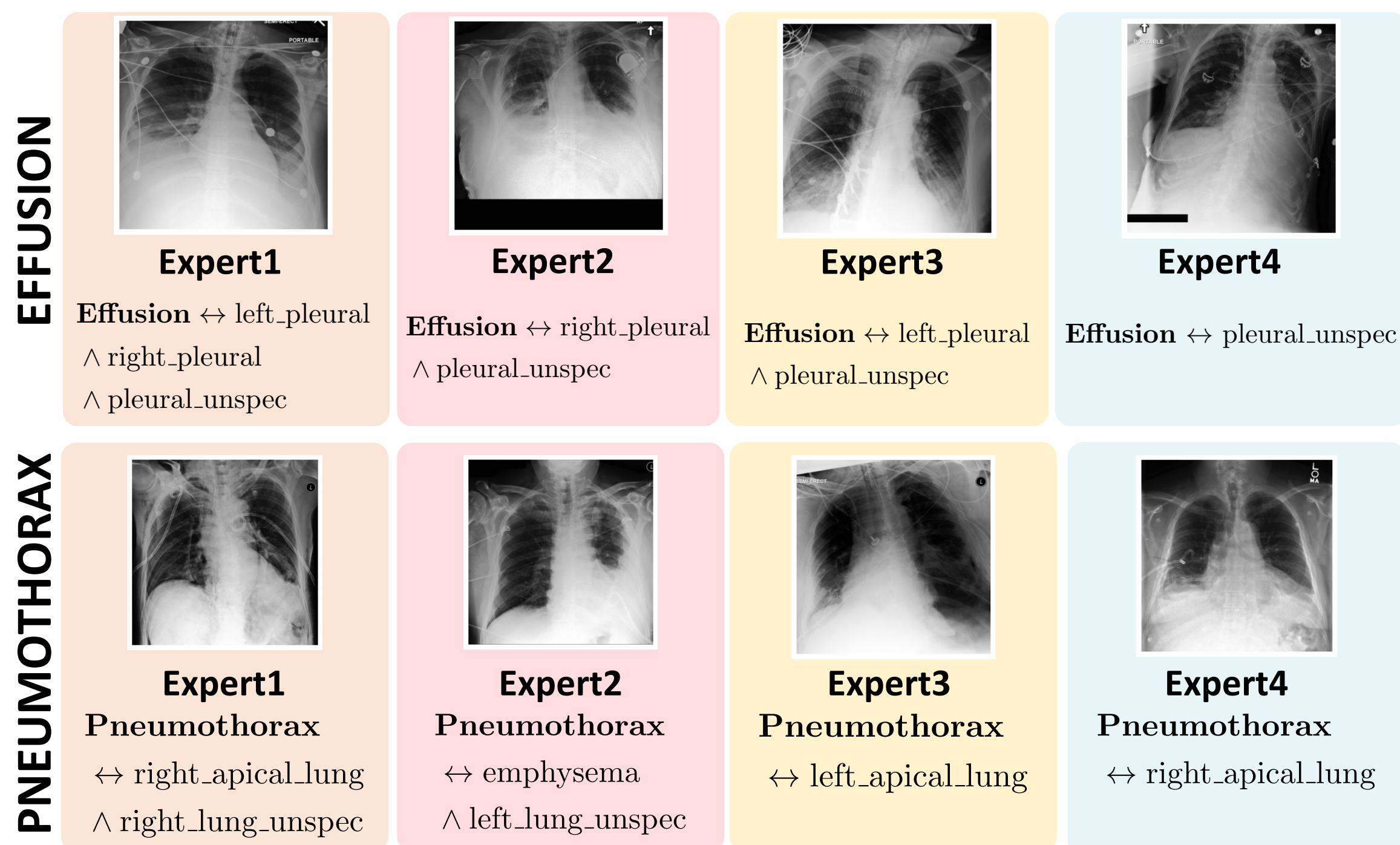
Right upper lobe consolidation

with adjacent.

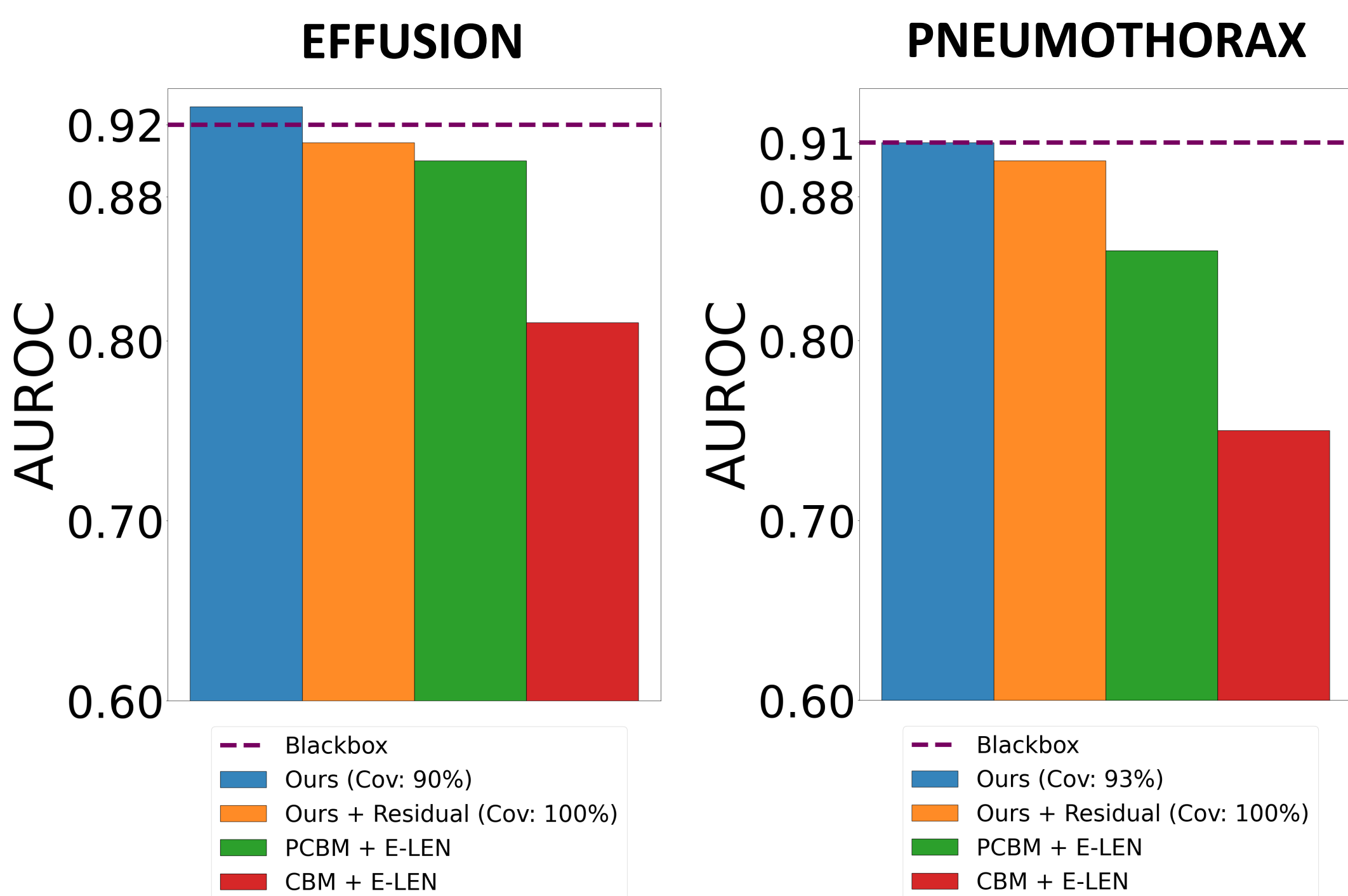
While this may be infectious in

nature, a CT scan is recommended for further clarification.

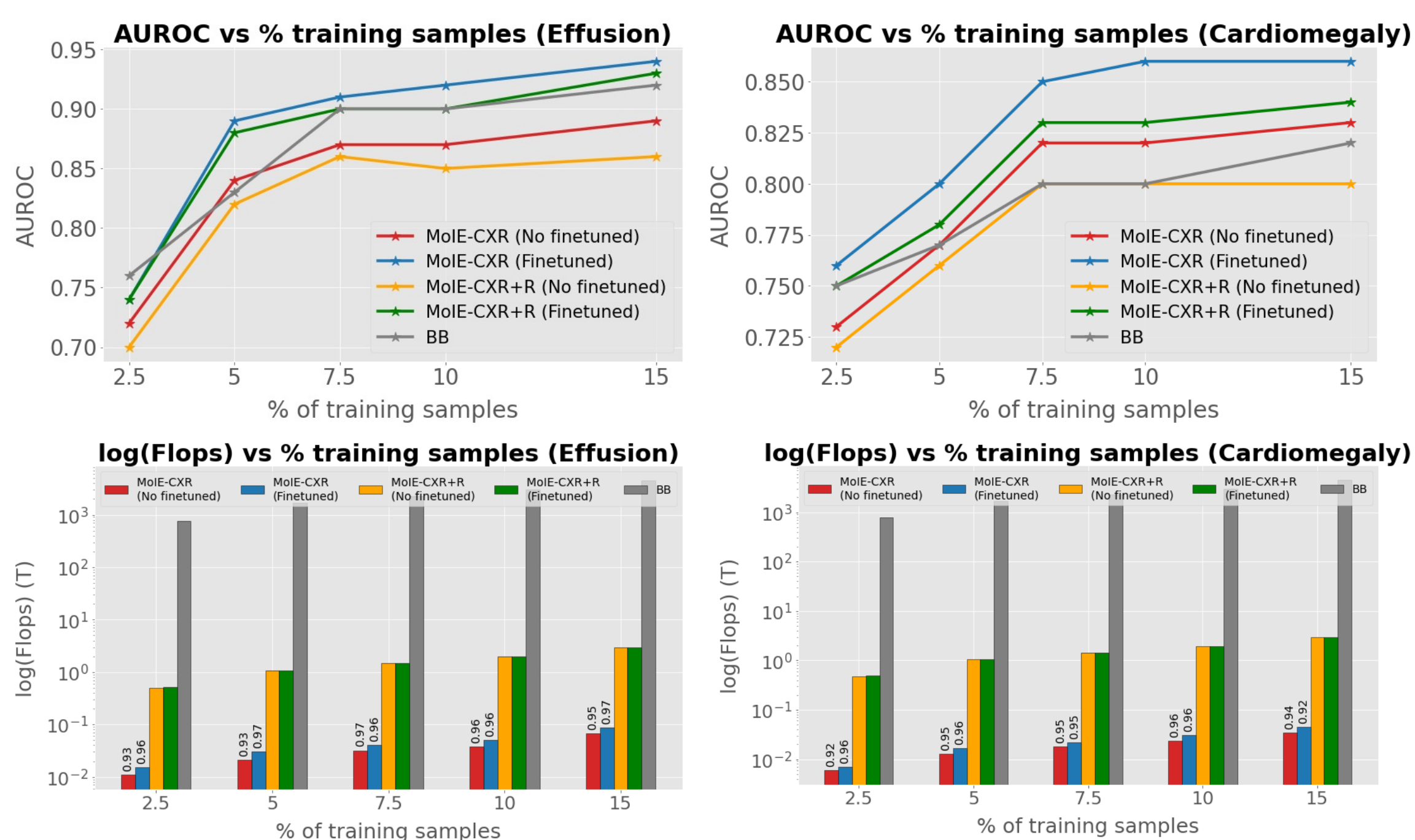
## Diversity in local explanations



## Not compromising the accuracy in MIMIC-CXR



## Transferring the first 3 experts of MIMIC-CXR to Stanford-CXR



We will also be in MICCAI 2023 to present **Distilling BlackBox to Interpretable models for Efficient Transfer Learning** (Early accept, top ~ 14%)