# Distilling Blackbox to Interpretable models for Efficient Transfer Learning

Shantanu Ghosh[1], Ke Yu[2], Kayhan Batmanghelich[1]

[1]Dept. Of Electrical and Computer Engineering, Boston University

[2]Intelligent Systems Program (ISP), University of Pittsburgh

**TLDR:** Extracting a mixture of interpretable models from a BlackBox to provide concept-based explanations for efficient transfer learning.

## Motivation

- Neural Networks fail to generalize due to scanner types, disease subtypes, patient subpopulation.
- Fine-tuning a Blackbox to a new domain can solve.
- This is data and computationally expensive.
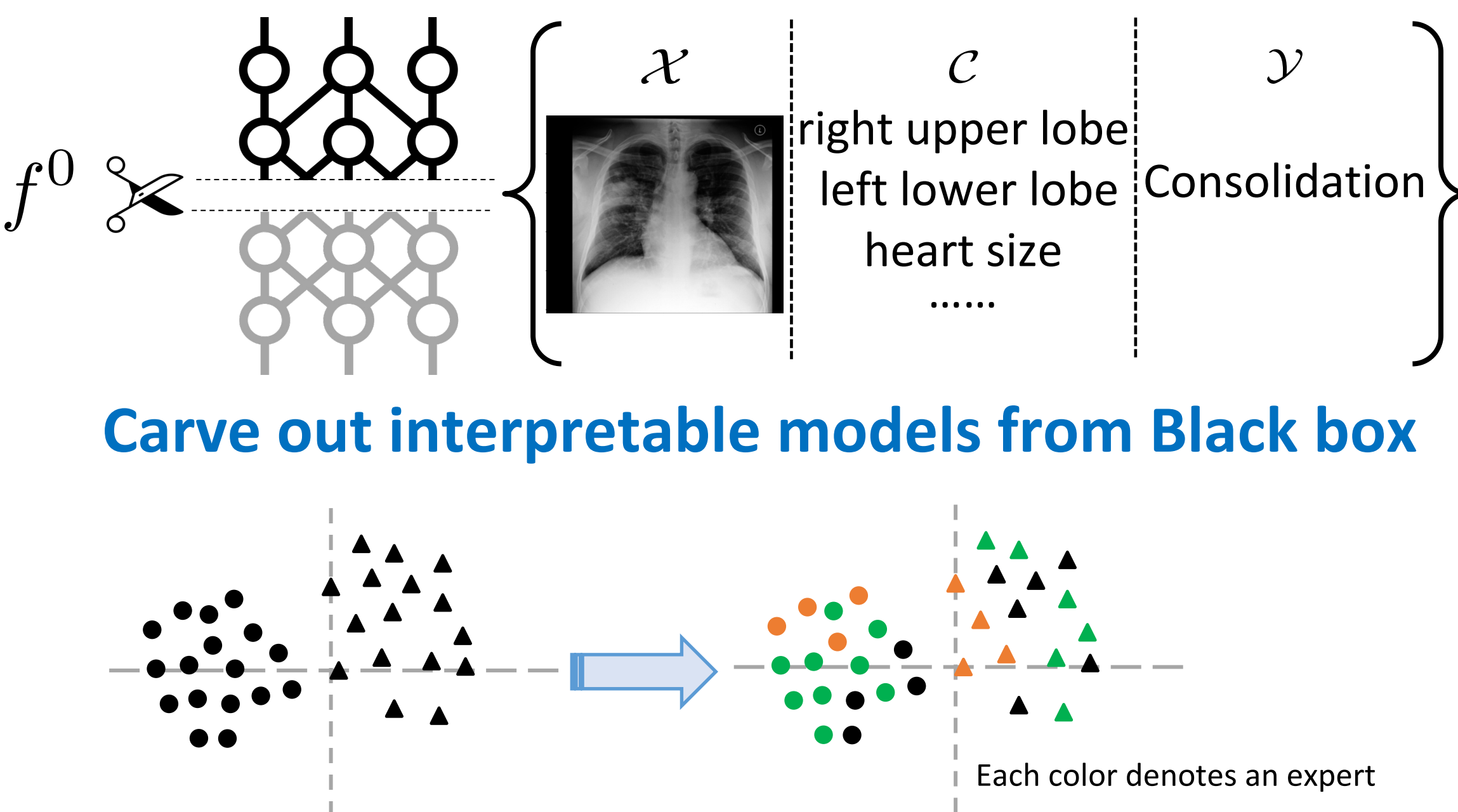- Whole process is not interpretable.

## Approach by radiologist

- Search for patterns for anatomical changes to read abnormality.
- Apply generalizable logical rules for disease diagnosis.
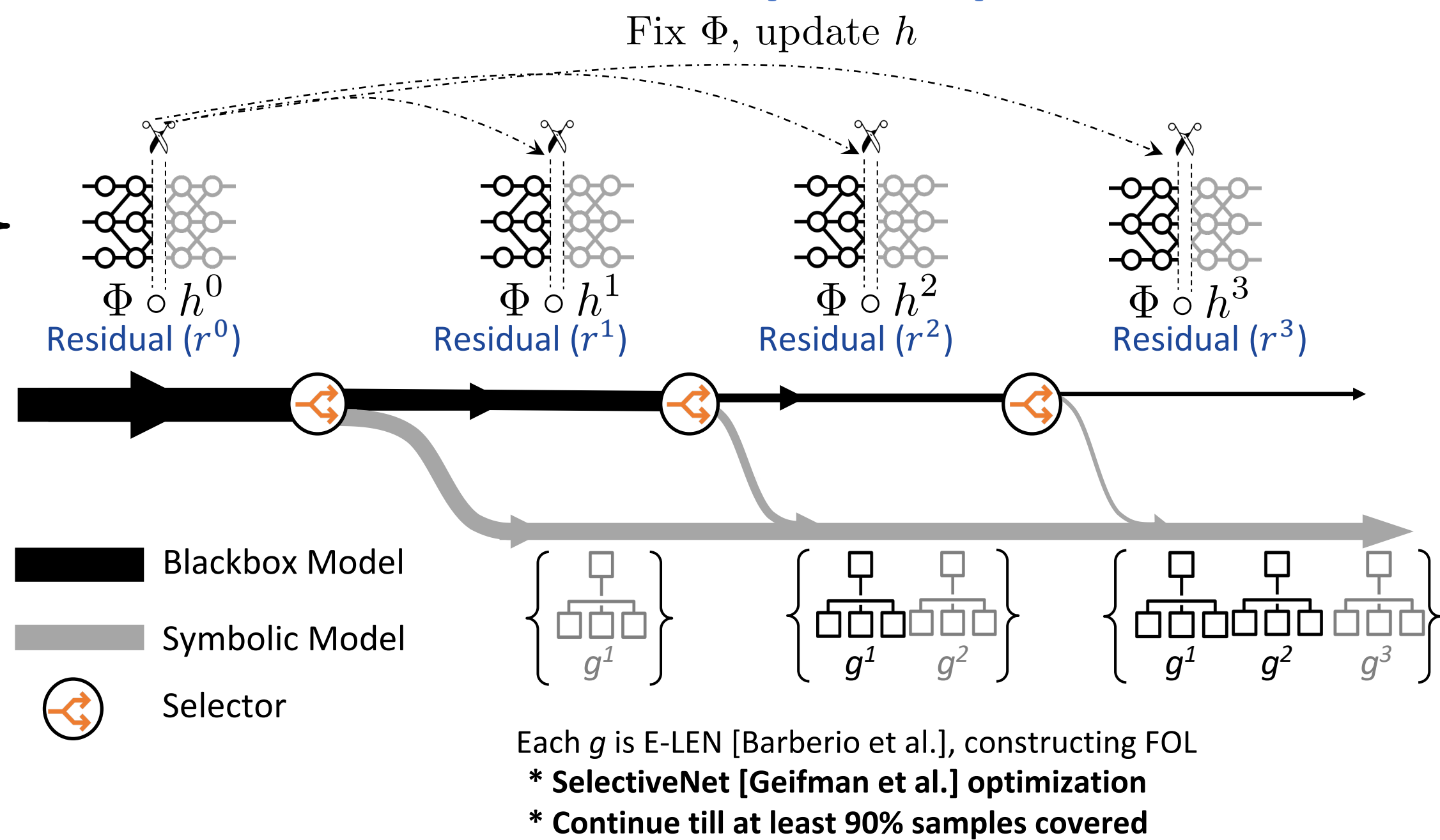- Whole process is interpretable.

## Design choices

- Carve a mixture of interpretable models from Blackbox.
- Built on domain-invariant anatomical concepts.
- Transfer the interpretable models to an unseen domain without any concept annotation.
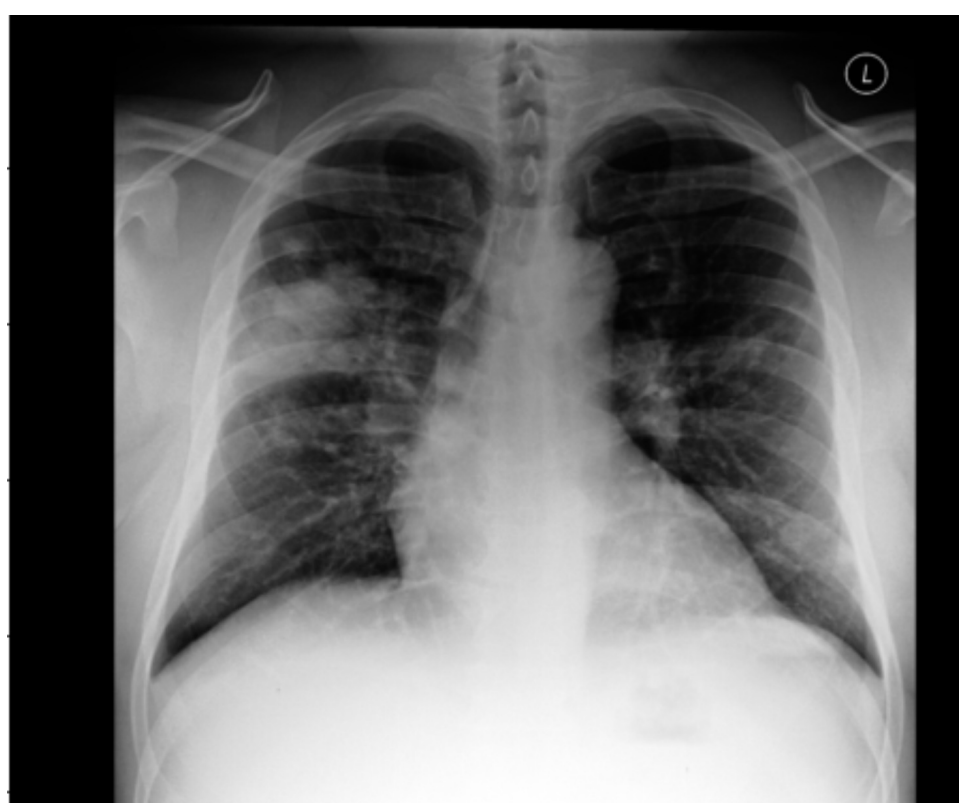
## Assumption



$f^0$

$\mathcal{X}$ — $\mathcal{C}$ right upper lobe, left lower lobe, heart size ...... — $\mathcal{Y}$ Consolidation

### Carve out interpretable models from Black box



Each color denotes an expert

## Route Interpret Repeat

Fix $\Phi$, update $h$



$\Phi \circ h^0$ Residual $(r^0)$    $\Phi \circ h^1$ Residual $(r^1)$    $\Phi \circ h^2$ Residual $(r^2)$    $\Phi \circ h^3$ Residual $(r^3)$

- Blackbox Model
- Symbolic Model
- Selector

$g^1$    $g^1$ $g^2$    $g^1$ $g^2$ $g^3$

Each $g$ is E-LEN [Barberio et al.], constructing FOL
* SelectiveNet [Geifman et al.] optimization
* Continue till at least 90% samples covered

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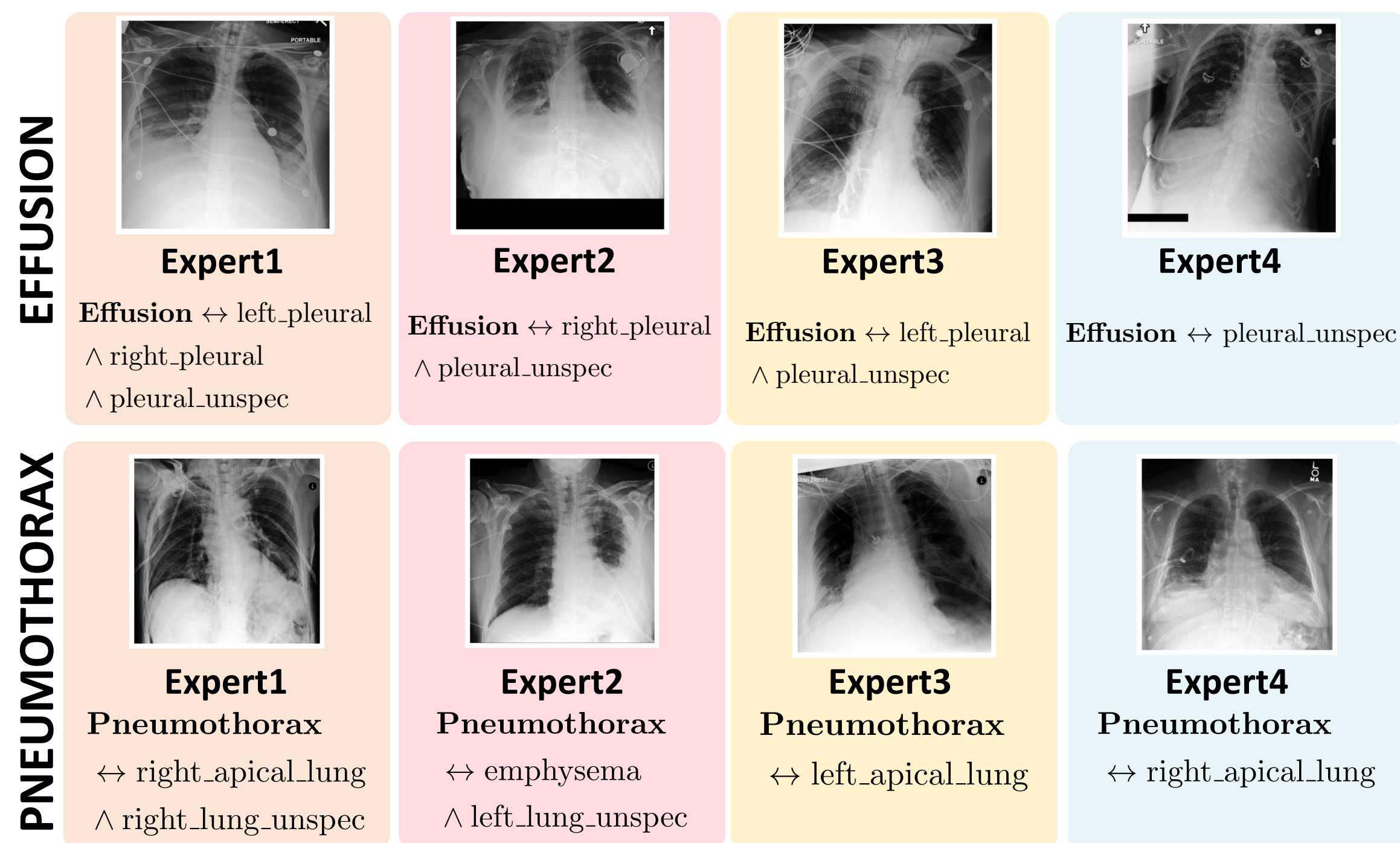


Ke Yu et al., MICCAI, 2022
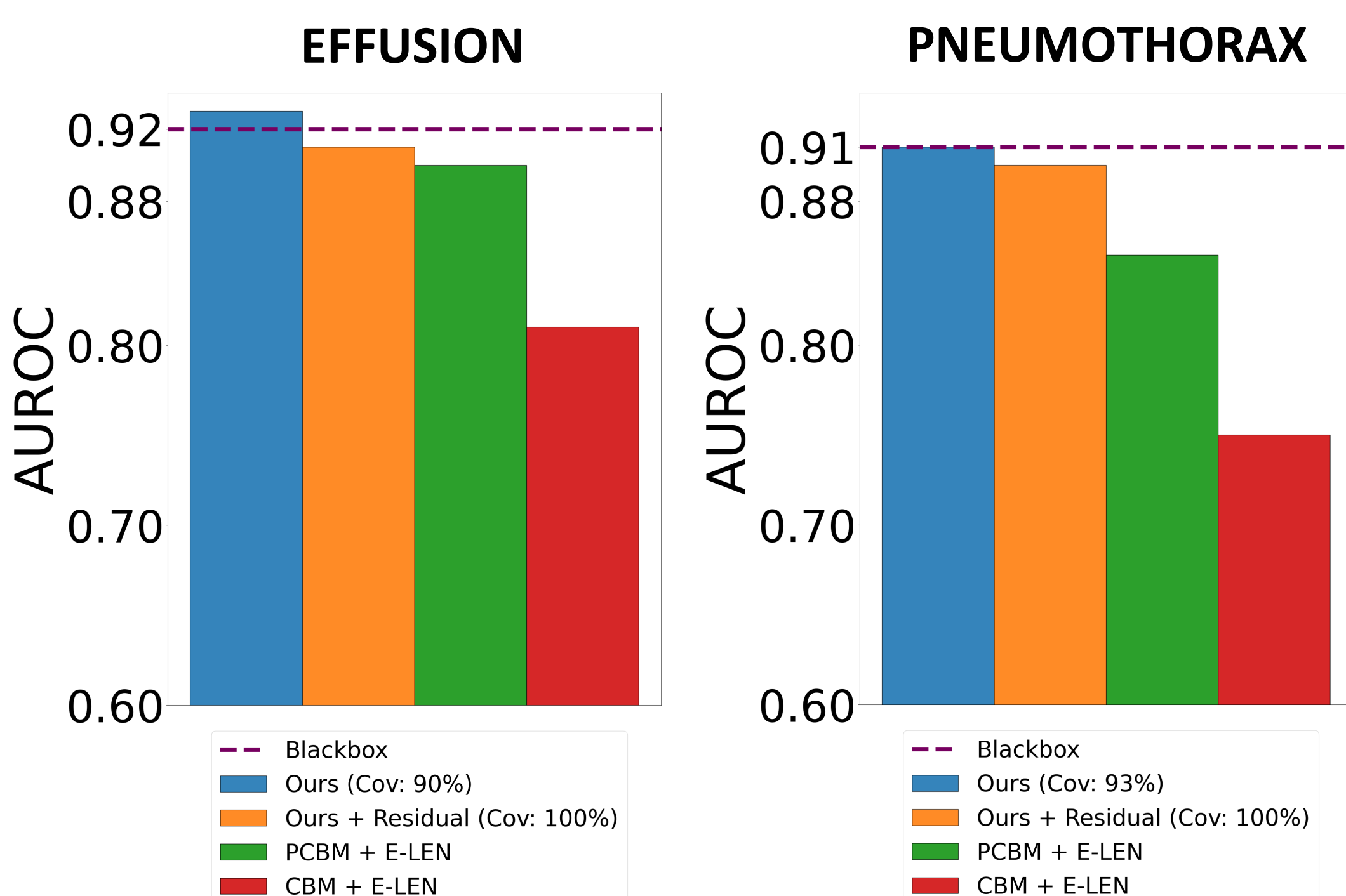
**Report:**

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



**EFFUSION**

**Expert1**
**Effusion** $\leftrightarrow$ left_pleural $\wedge$ right_pleural $\wedge$ pleural_unspec

**Expert2**
**Effusion** $\leftrightarrow$ right_pleural $\wedge$ pleural_unspec

**Expert3**
**Effusion** $\leftrightarrow$ left_pleural $\wedge$ pleural_unspec

**Expert4**
**Effusion** $\leftrightarrow$ pleural_unspec

**PNEUMOTHORAX**

**Expert1**
Pneumothorax $\leftrightarrow$ right_apical_lung $\wedge$ right_lung_unspec

**Expert2**
Pneumothorax $\leftrightarrow$ emphysema $\wedge$ left_lung_unspec

**Expert3**
Pneumothorax $\leftrightarrow$ left_apical_lung

**Expert4**
Pneumothorax $\leftrightarrow$ right_apical_lung

## Not compromising the accuracy in MIMIC-CXR



EFFUSION    PNEUMOTHORAX

- Blackbox
- Ours (Cov: 90%) / Ours (Cov: 93%)
- Ours + Residual (Cov: 100%)
- PCBM + E-LEN
- CBM + E-LEN

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



AUROC vs % training samples (Effusion)    AUROC vs % training samples (Cardiomegaly)

log(Flops) vs % training samples (Effusion)    log(Flops) vs % training samples (Cardiomegaly)

- MoIE-CXR (No finetuned)
- MoIE-CXR (Finetuned)
- MoIE-CXR+R (No finetuned)
- MoIE-CXR+R (Finetuned)
- BB