

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?

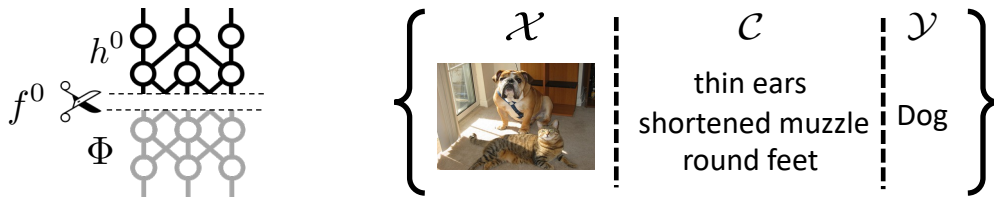
Desirable properties

- Does not compromise the performance.
- Can be intervened to fix the misclassification

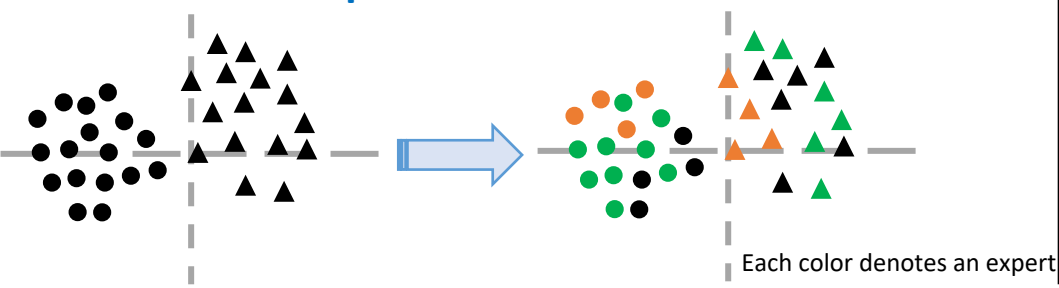
Design choices

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

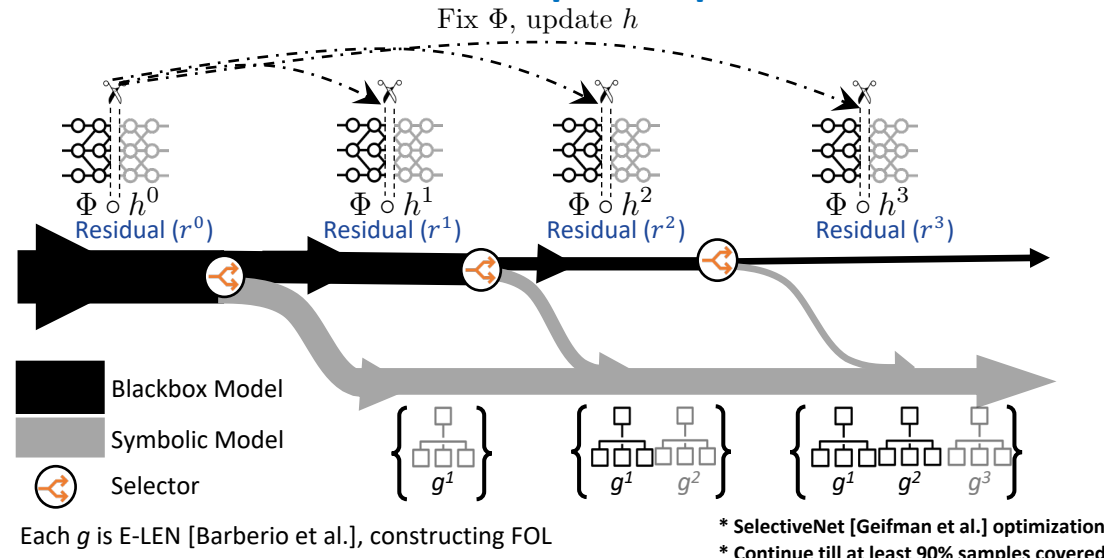
Assumptions



Carve out interpretable models from Black box

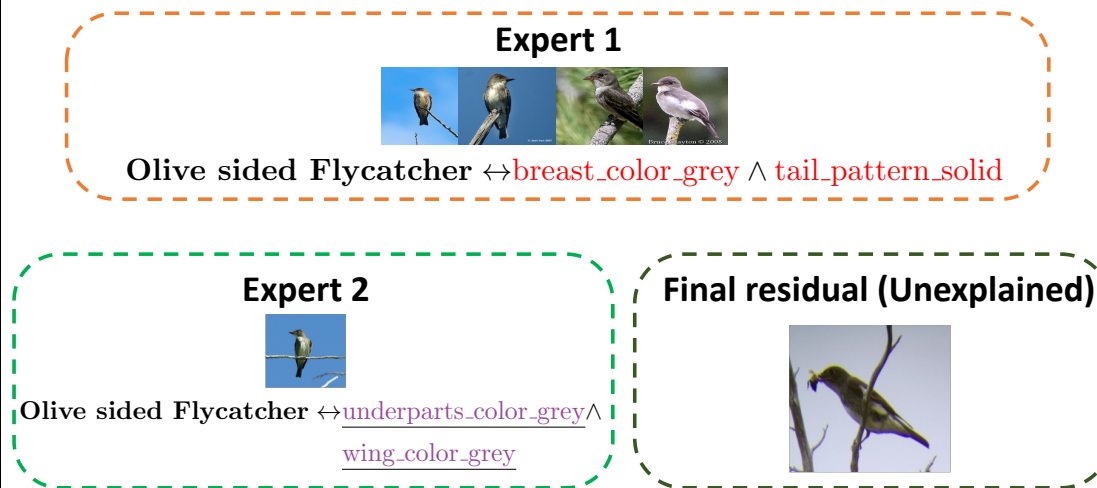


Route Interpret Repeat

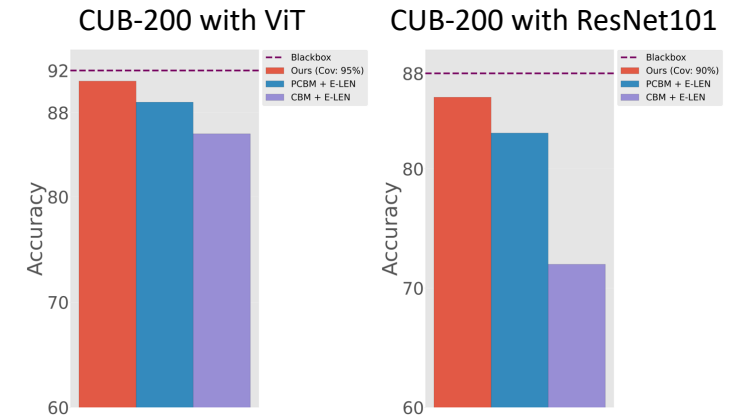


Capturing heterogenous explanations

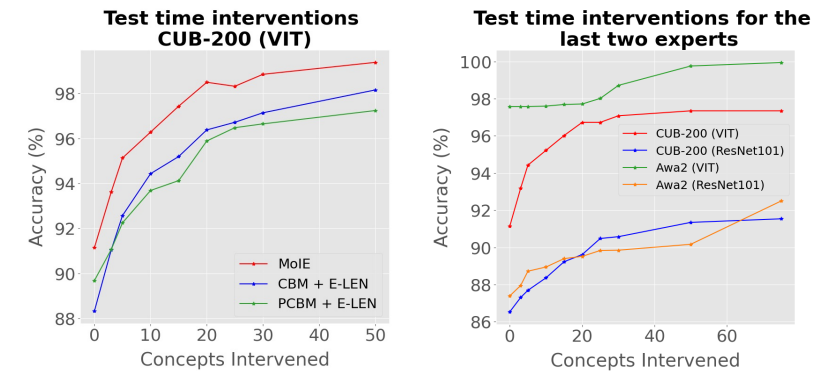
* Extracted from ViT-based BlackBox



Not compromising performance



Test time interventions



- Also in our paper,
- + we experiment with a diverse set of datasets and architectures
 - + we achieve higher concept completeness scores
 - + ViT-based experts compose less concepts than CNN-based
 - + we eliminate shortcut learning problem (SCIS w)
 - + we efficiently transfer the experts to new domain (IMLH w)

