

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.

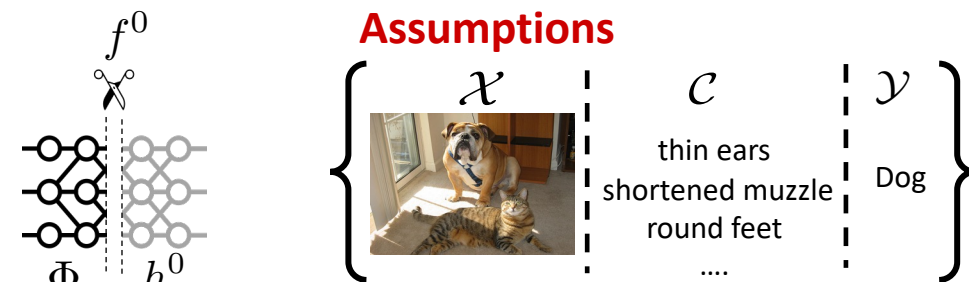
How to blur this gap?

Desirable properties

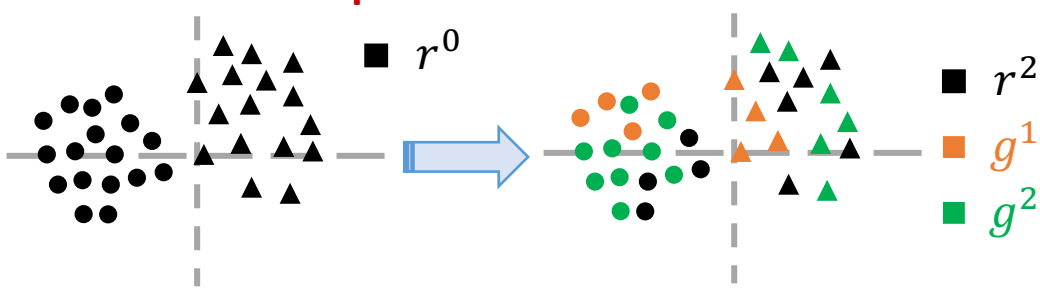
- Does compromise the performance.
- Can be intervened to fix the misclassification

Design choices

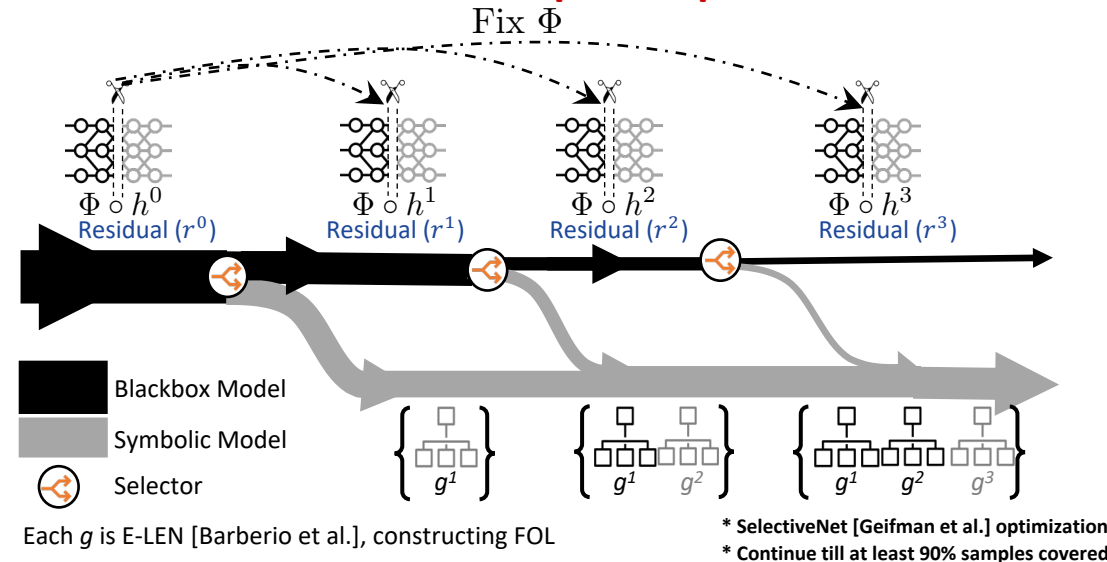
- Carve interpretable models from Blackbox.
- Concept based
- First order logic for concept interaction



Carve out interpretable models from Black box

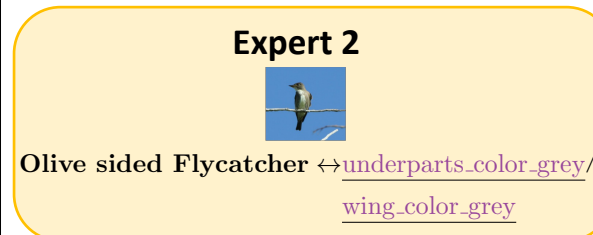


Route Interpret Repeat

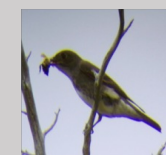


Capturing heterogenous explanations

* Extracted from ViT-based BlackBox

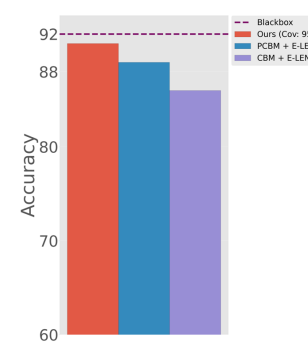


Final residual (Unexplained)

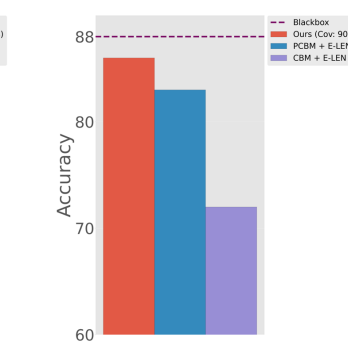


Not compromising performance

CUB-200 with ViT

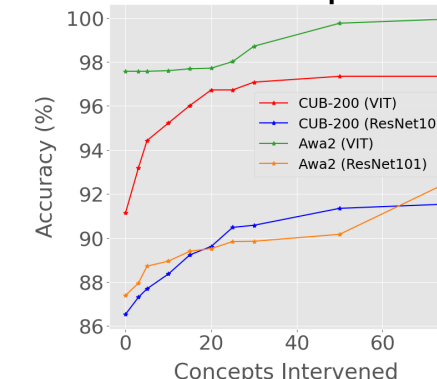


CUB-200 with ResNet101



Test time intervention on harder samples

Test time interventions for the last two experts



Also in our paper,

- + we compare with more datasets and baselines
- + we achieve higher concept completeness scores
- + we achieve higher accuracy during test time intervention
- + we eliminate shortcut learning problem (SCIS w)
- + we efficiently transfer the experts to new domain (IMLH w)
- + ViT-based experts compose less concepts than CNN-based