



Tackling Shortcut Learning in Deep Neural Networks: An Iterative Approach with Interpretable Models

BOSTON UNIVERSITY

BATMAN LAB



Shantanu Ghosh¹, Ke Yu², Forough Arabshahi³, Kayhan Batmanghelich¹

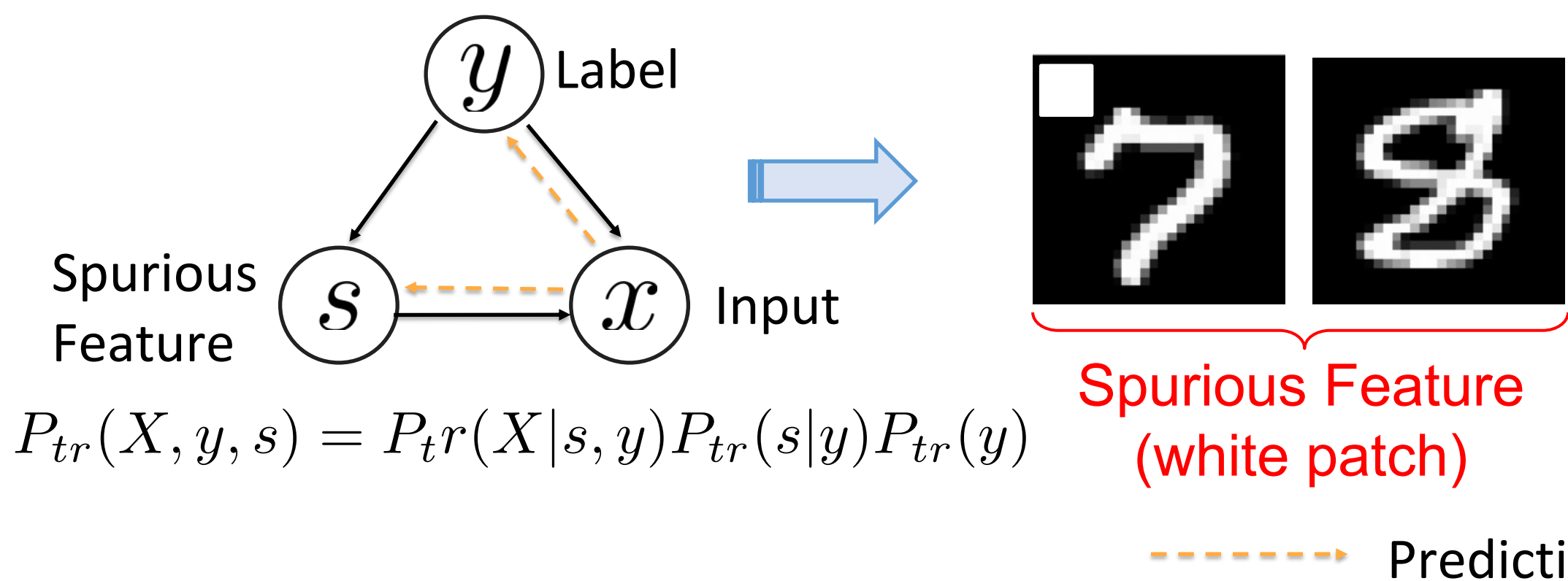
¹Dept. Of Electrical and Computer Engineering, Boston University

²Intelligent Systems Program (ISP), University of Pittsburgh, ³Meta AI

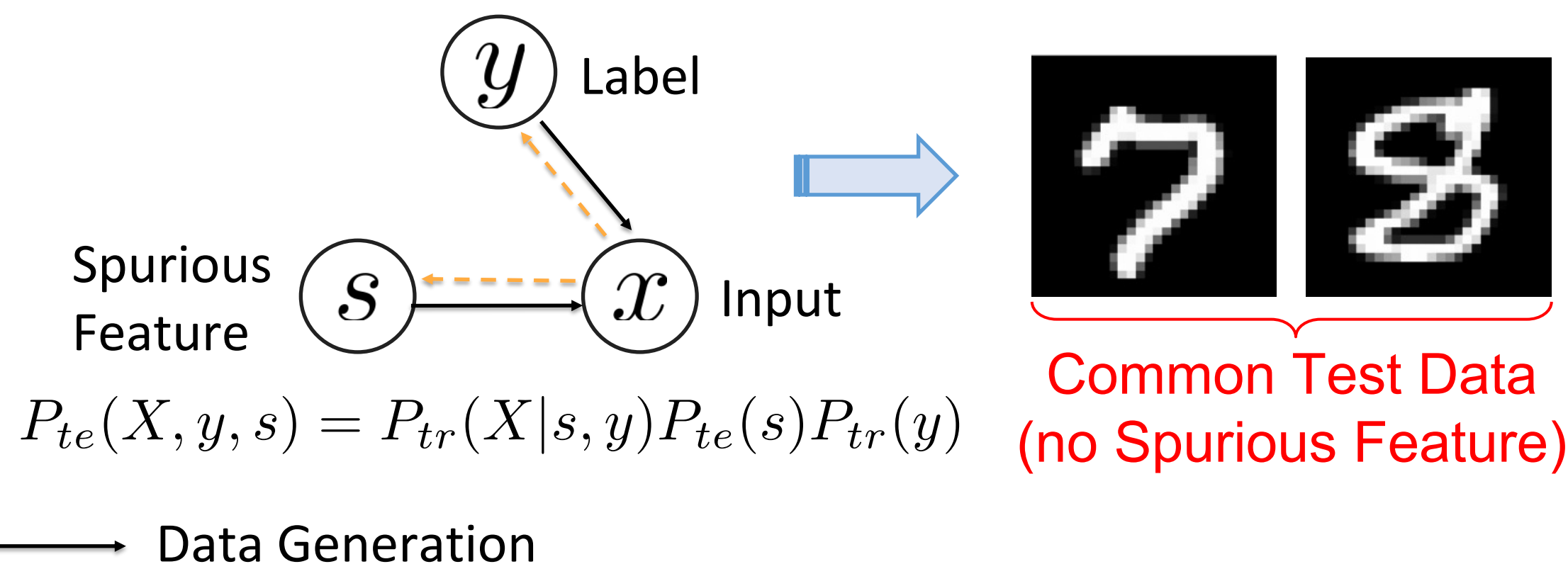
TLDR: We use concept based interpretable method to mitigate the problem of shortcut learning.

What is shortcut?

(A) Training Data w/ Shortcut

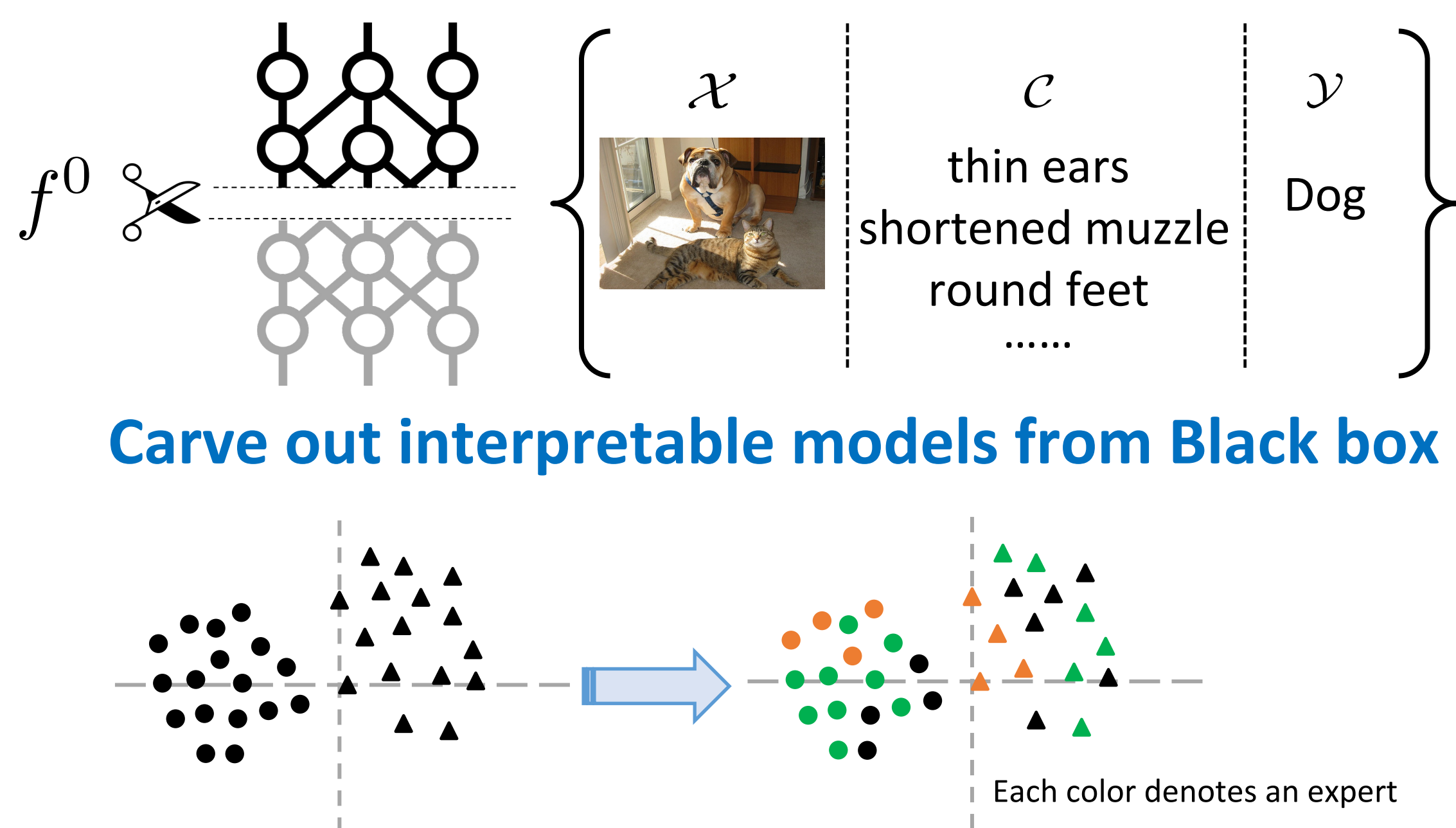


(B) Testing Data w/o Shortcut

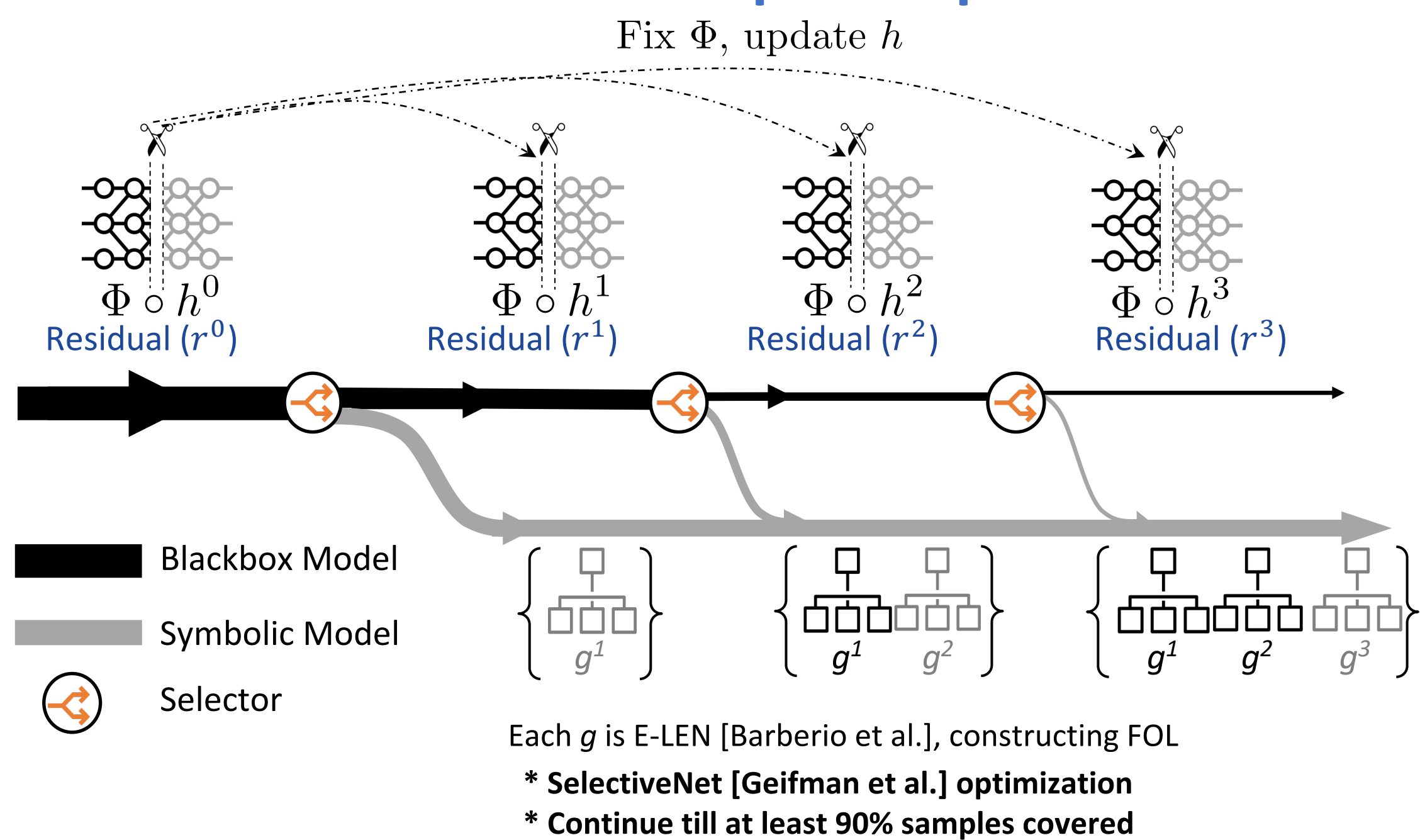


GAP: Existing shortcut elimination methods are not interpretable.

Assumption



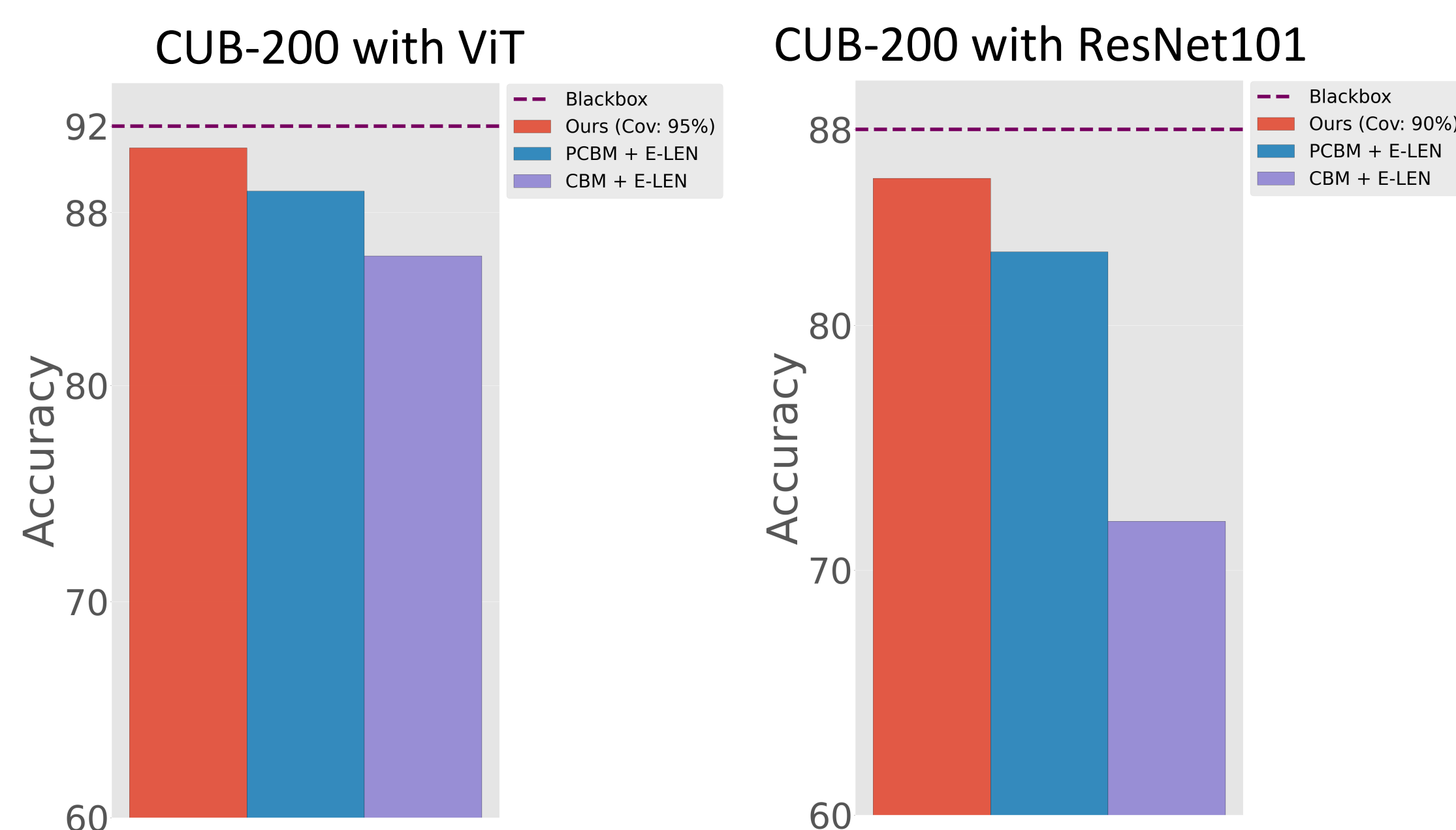
Route Interpret Repeat



Algorithm 1 Applying MoIE to eliminate shortcuts

- Input:** $\mathcal{D} = \{x_j, c_j, y_j\}_{j=1}^n$; biased BB $f^0 = h^0(\Phi(\cdot))$; The total iterations K ; Coverages τ_1, \dots, τ_K . Freeze Φ .
- Using (Yuksekgonul et al., 2022) learn the projection t to predict the concept value.
- Detection step.** Learn the experts in MoIE $\{g\}_{k=1}^K$ and extract the FOLs. The FOL contains shortcuts.
- Elimination step.** Consider the detected shortcut concept in the “Detection” step as metadata and finetune BB (f^0) with MDN (Lu et al., 2021) to remove the role of that shortcut.
- Retrain t with Φ of finetuned BB to get the concepts.
- Verification step.** Learn MoIE $\{g\}_{k=1}^K$ again from retrained t and recompute the FOLs. The final FOLs do not contain spurious concepts as they have been eliminated in the “Elimination step”.

Not compromising the accuracy



Comparison w/ other methods on Waterbirds dataset

| Method | Avg Acc. | Worst Acc. |
|---|---------------------------|---------------------------|
| ERM (Wah et al., 2011) | 97.0 ± 0.2% | 63.7 ± 1.9% |
| ERM+aug (Wah et al., 2011) | 87.4 ± 0.5% | 76.4 ± 2.0% |
| UW (Xian et al., 2018) | 96.3.0 ± 0.3% | 76.2 ± 1.4% |
| IRM (Arjovsky et al., 2020) | 87.5 ± 0.7% | 75.6 ± 3.1% |
| IB-IRM (Ahuja et al., 2022) | 88.5 ± 0.9% | 76.5 ± 1.2% |
| V-REx (Krueger et al., 2021) | 88.0 ± 1.4% | 73.6 ± 0.2% |
| CORAL (Sun & Saenko, 2016) | 90.3 ± 1.1% | 79.8 ± 1.8% |
| Fish (Shi et al., 2021) | 85.6 ± 0.4% | 64.0 ± 0.3% |
| GroupDRO (Sagawa et al., 2019) | 91.8 ± 0.3% | 90.6 ± 1.1% |
| JTT (Liu et al., 2021) | 93.3 ± 0.3% | 86.7 ± 1.5% |
| DM-ADA (Xu et al., 2020) | 76.4 ± 0.3% | 53.0 ± 1.3% |
| LISA (Yao et al., 2022) | 91.8 ± 0.3% | 88.5 ± 0.8% |
| BB w MDN (ours) | 95.01 ± 0.5% | 94.4 ± 0.5% |
| MoIE from BB w MDN (ours) (COVERAGE) | 91.0 ± 0.5% (0.91) | 93.7 ± 0.4% (0.87) |
| MoIE+R from BB w MDN (ours) | 90.2 ± 0.5% | 92.1 ± 0.4% |

Qualitative result on Waterbirds dataset

Biased Blackbox

Groundtruth: WaterBird
 Prediction : LandBird
 Explanation : LandBird ↔ WingShapeRoundedwings ∧ **Forest**

Robust Blackbox

Groundtruth: WaterBird
 Prediction : WaterBird
 Explanation : WaterBird ↔ BillLengthAboutTheSameAsHead
 ∧ ¬BillLengthShorterThanHead ∧ ¬SizeSmall5...9in
 ∧ ¬ShapePerchingLike ∧ CrownColorWhite