**Auroc scores of the background concepts**

Auroc

Biased Model
Robust Model

Ocean  Lake  Bamboo  Forest
Concepts

**Accuracy scores the biased Blackbox**

Accuracy

Land
Water

LandBird    WaterBird

(a)

**Accuracy scores MoIE from the robust Blackbox**

Accuracy

Land
Water

LandBird    WaterBird

(b)

**Biased Blackbox**

Groundtruth: WaterBird

Prediction : LandBird

Explanation : LandBird ↔ WingShapeRoundedwings ∧ **Forest**

**Robust Blackbox**

Groundtruth: WaterBird

Prediction : WaterBird

Explanation : WaterBird ↔ BillLengthAboutTheSameAsHead
∧ ¬BillLengthShorterThanHead  ∧ ¬SizeSmall5_9in
∧ ¬ShapePerchingLike ∧ CrownColorWhite

(c)

**Accuracy scores of the background concepts**

Accuracy

Biased Model
Robust Model

Ocean   Lake   Bamboo   Forest
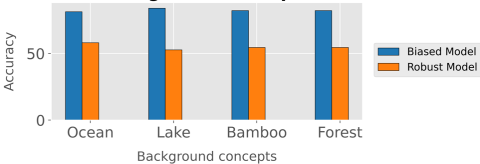Background concepts

(d)

(e)

**Accuracy scores the biased Blackbox**

Accuracy

Land
Water

LandBird    WaterBird
Background