# Concept-based Explanations for Out-Of-Distribution Detectors

**Jihye Choi**, Jayaram Raghuram, Ryan Feng, Jiefeng Chen,

Somesh Jha, Atul Prakash
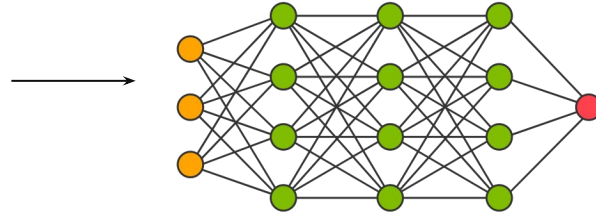
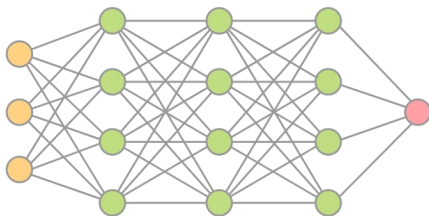# Standard Machine Learning (ML) Models



ML model in self-driving car

Training

"Car"

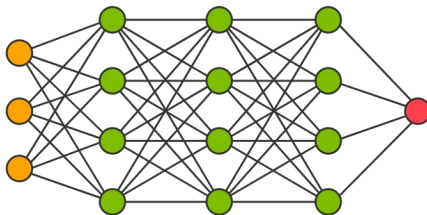# Standard Machine Learning (ML) Models

ML model in self-driving car

Training  →  → "Car"
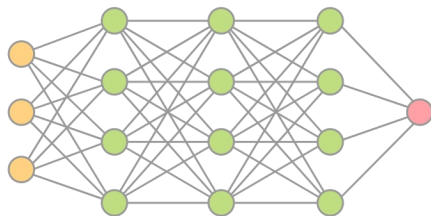
Testing  →  → "Car"

In-distribution
(ID)

# Standard Machine Learning (ML) Models
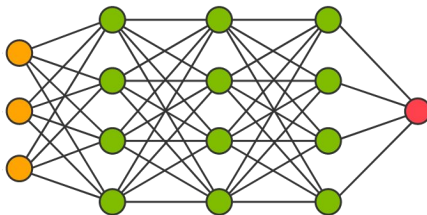
ML model in self-driving car

Training → → "Car"

Testing

Out-of-distribution
(OOD)
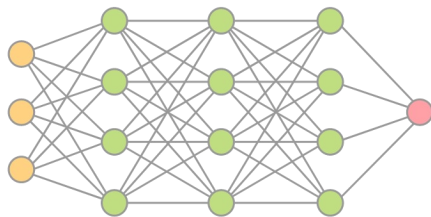
→ → "Car"
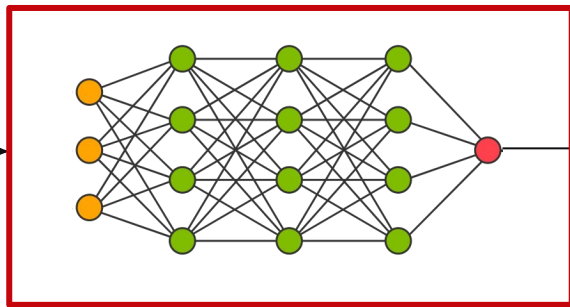
Overconfident prediction
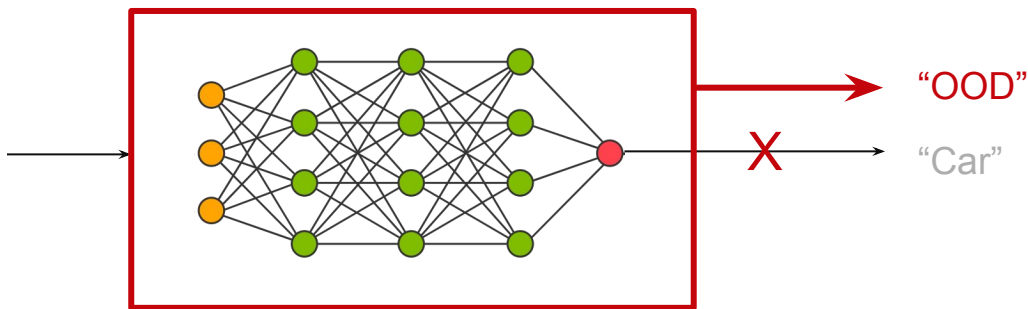for unseen object "Buffalo"

# Out-Of-Distribution (OOD) Detection

# Understanding OOD Detection

Why a given OOD detector decides an input to be OOD?

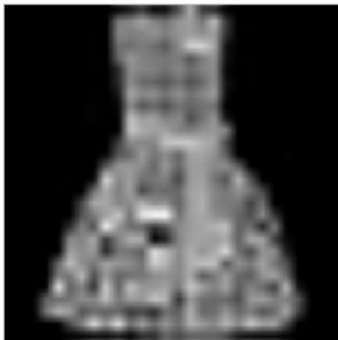Can we use existing ML explanation methods for classification to interpret OOD detection results?
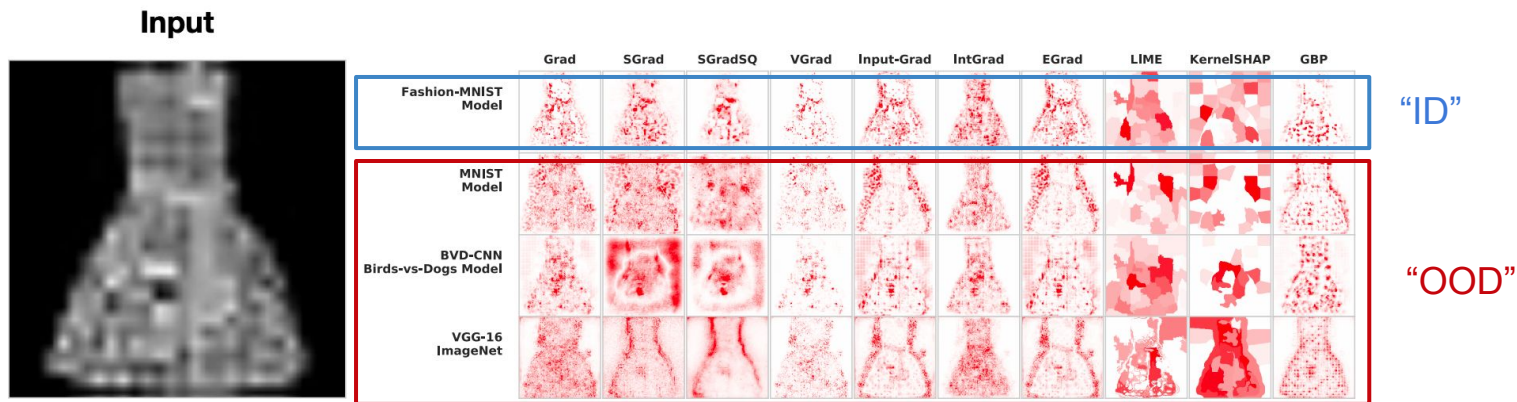


OOD Detector

Testing

"OOD"

X "Car"

# ML Explanations for Classification

[Type 1] Feature Attributions

**Input**



[Adebayo et al., NeurIPS'20]

# ML Explanations for Classification

[Type 1] Feature Attributions



[Adebayo et al., NeurIPS'20]

# ML Explanations for Classification

[Type 1] Feature Attributions
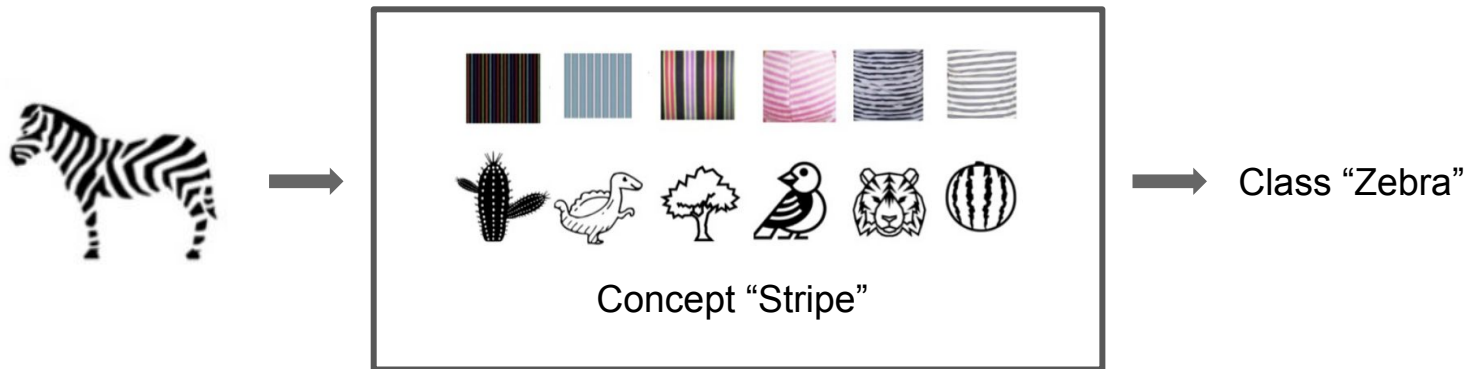


[Adebayo et al., NeurIPS'20]

# ML Explanations for Classification

[Type 1] Feature Attributions



Pixel-level activations might not be the most intuitive form of explanations for humans
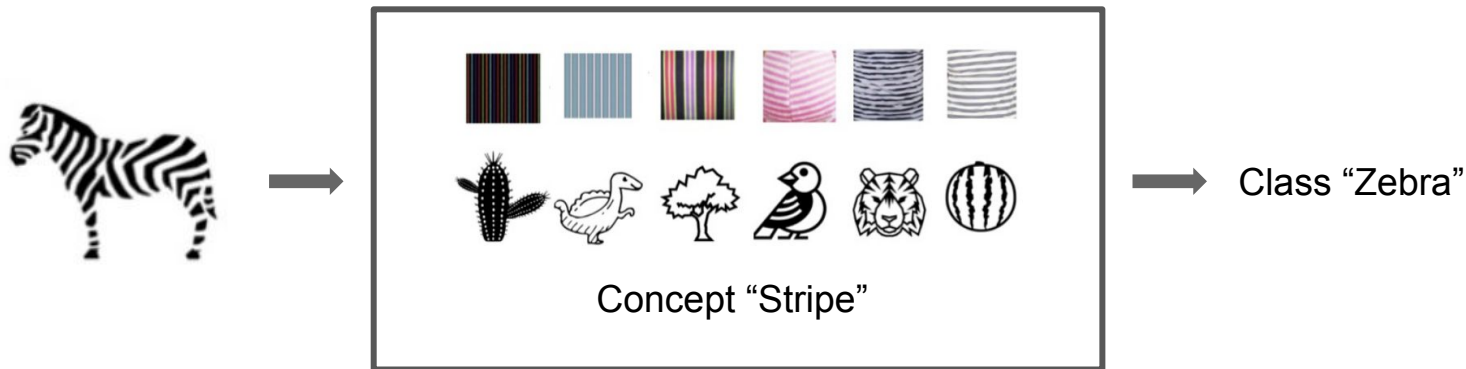
[Adebayo et al., NeurIPS'20]

# ML Explanations for Classification

[Type 2] Concept-based Explanations



Concept "Stripe"

Class "Zebra"

[Kim et al., ICML'18]

# ML Explanations for Classification

[Type 2] Concept-based Explanations



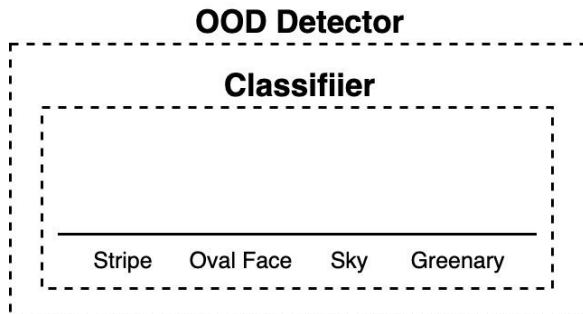Concept "Stripe"

Class "Zebra"

The use of concept-based explanations for OOD detectors remains unexplored
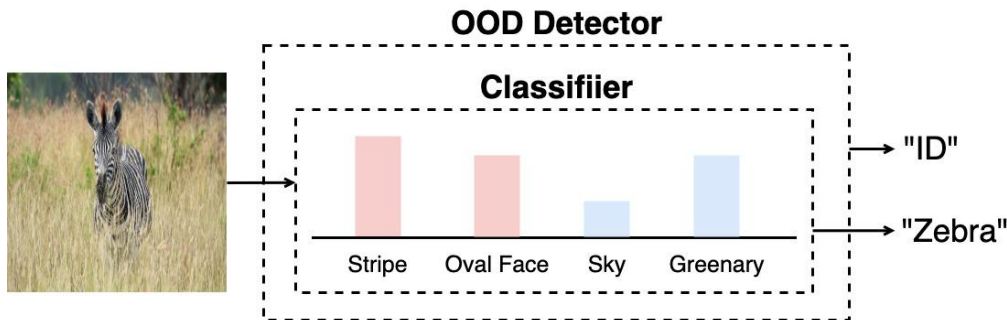
[Kim et al., ICML'18]

# Concept-based Explanation for OOD Detection

Our work: the first method to understand the decisions of an OOD detector in terms of *high-level concepts*

GIven DNN classifier and OOD detector, find a set of concepts that sufficiently explain their behaviors.
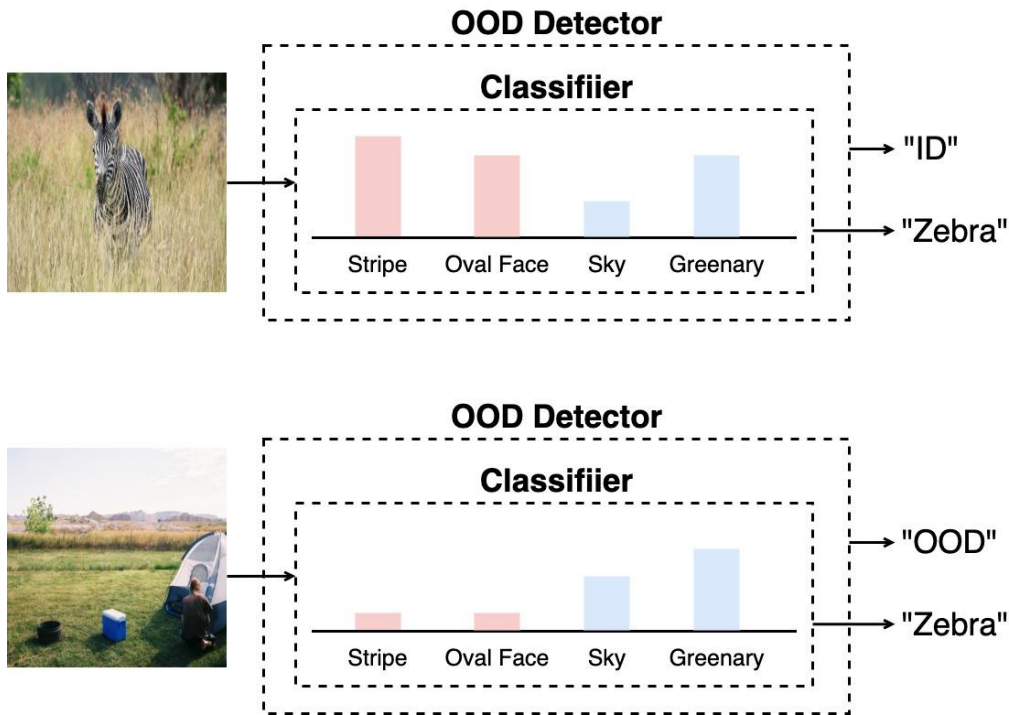
# Concept-based Explanation for OOD Detection

Our work: the first method to understand the decisions of an OOD detector in terms of *high-level concepts*

Observe normal concept activations patterns given ID inputs.

# Concept-based Explanation for OOD Detection

Our work: the first method to understand the decisions of an OOD detector in terms of *high-level concepts*

Given OOD inputs, we observe different concept activation patterns compared to that of ID inputs.
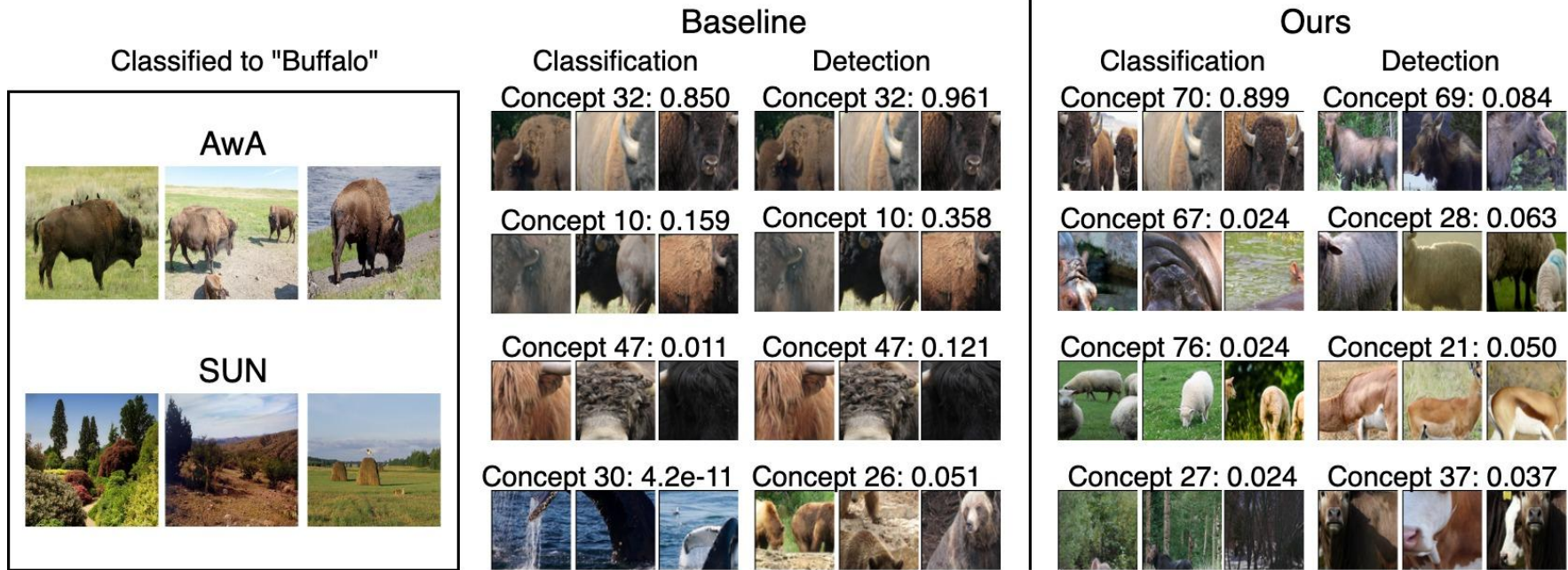
# Our Contributions

Given DNN classifier and OOD detector, we

1.  propose metrics to quantify the effectiveness of concept-based explanation for OOD detection:
    a.  *Detection Completeness*: are the concept scores sufficient statistics for class predictions and OOD detection?
    b.  *Concept Separability*: are ID and OOD inputs clearly distinctive in terms of concepts?

2.  introduce general concept learning framework that discovers a set of concepts that have good detection completeness and concept separability.

3.  by using the concepts learned by our framework, show how to identify prominent concepts that contribute to an OOD detector's decisions, and provide insights for popular OOD detectors.

# Results

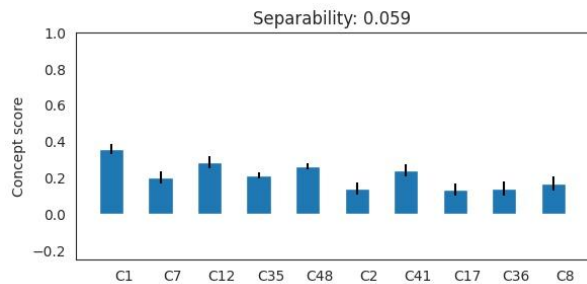Concept-based explanations for Energy detector



Energy detector: [Liu et al., NeurIPS'20]
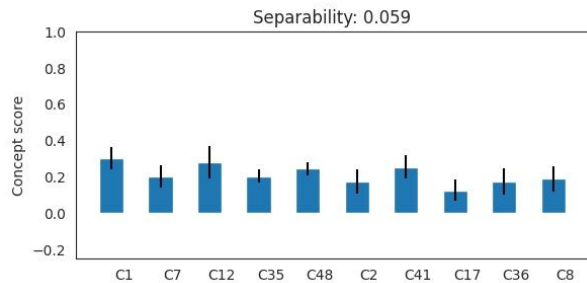Baseline: [Yeh et al., NeurIPS'20]

# Results

Concept score patterns between inputs detected as ID vs OOD
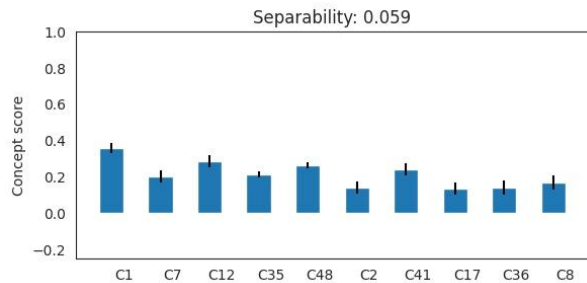


low separability, inputs detected as ID
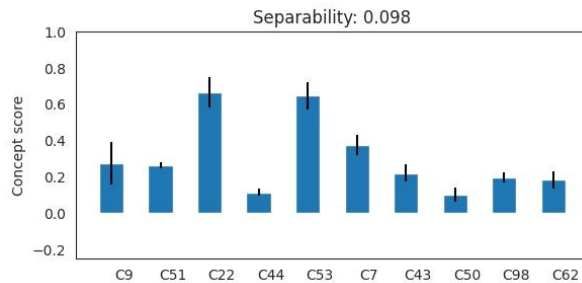


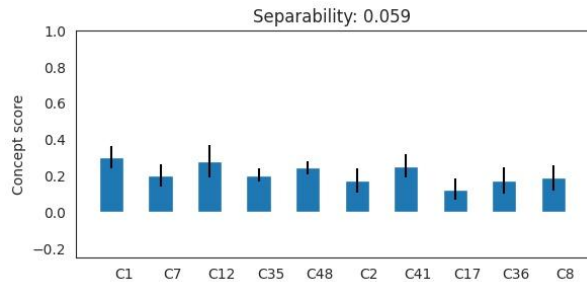low separability, inputs detected as ID

# Results

Concept score patterns between inputs detected as ID vs OOD



low separability, inputs detected as ID
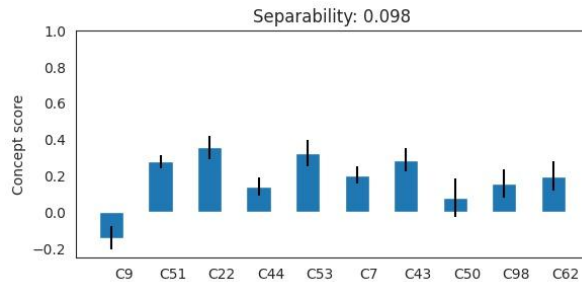
high separability, inputs detected as ID

low separability, inputs detected as ID

high separability, inputs detected as ID

# Thank you

For complete description of our method and full results, please check out our paper!