# Fighting Fire with Fire: Avoiding DNN Shortcuts through Priming

Chuan Wen, Jianing Qian, Jierui Lin, Jiaye Teng, Dinesh Jayaraman, Yang Gao

Tsinghua University · UNIVERSITY of PENNSYLVANIA · TEXAS The University of Texas at Austin

## Problem: Shortcuts in DNNs

DNNs often struggle to disambiguate between competing hypotheses for a target concept, and end up learning "shortcuts".
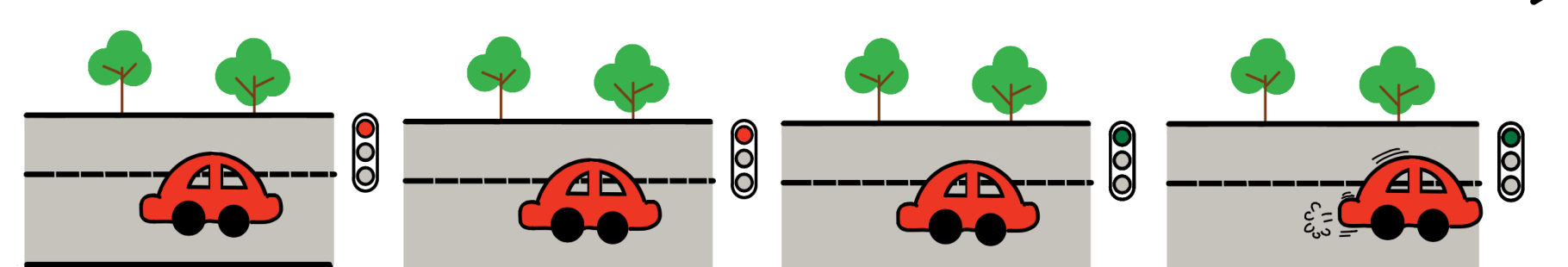
### Example 1: Image Classification

in-distribution | out-of-distribution

prediction: cow ✓ | prediction: cat ✗

DNNs cheat by relying on backgrounds.

### Example 2: Behavioral cloning

label $a_1 = 0$   $a_2 = 0$   $a_3 = 1$   $a_4 = 1$

copy → copy → copy →

**Shortcut Solution** $\hat{a}_1 = 0$ ✓   $\hat{a}_2 = 0$ ✓   $\hat{a}_3 = 0$ ✗   $\hat{a}_4 = 1$ ✓

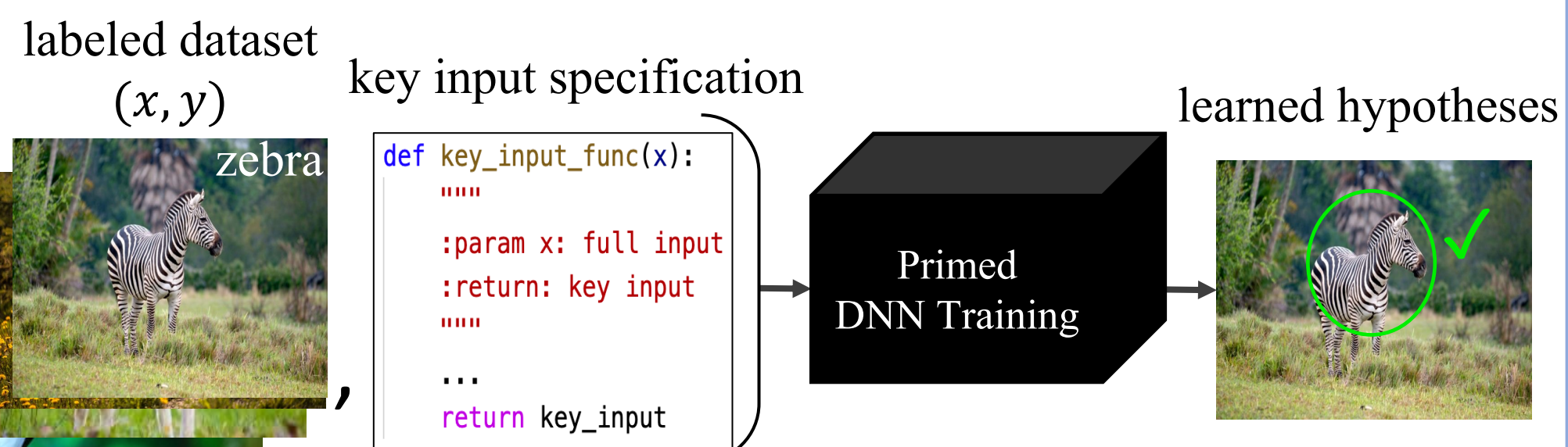DNNs cheat by copying from the previous action during training.

## Motivation:

This is a zebra!

The training data does not fully specify the task!

Humans rely on additional knowledge e.g. task-relevant inputs

## Our Idea:

Use auxiliary knowledge to "prime" DNNs away from shortcuts.

labeled dataset $(x, y)$   key input specification   learned hypotheses

```
def key_input_func(x):
    """
    :param x: full input
    :return: key input
    """
    ...
    return key_input
```
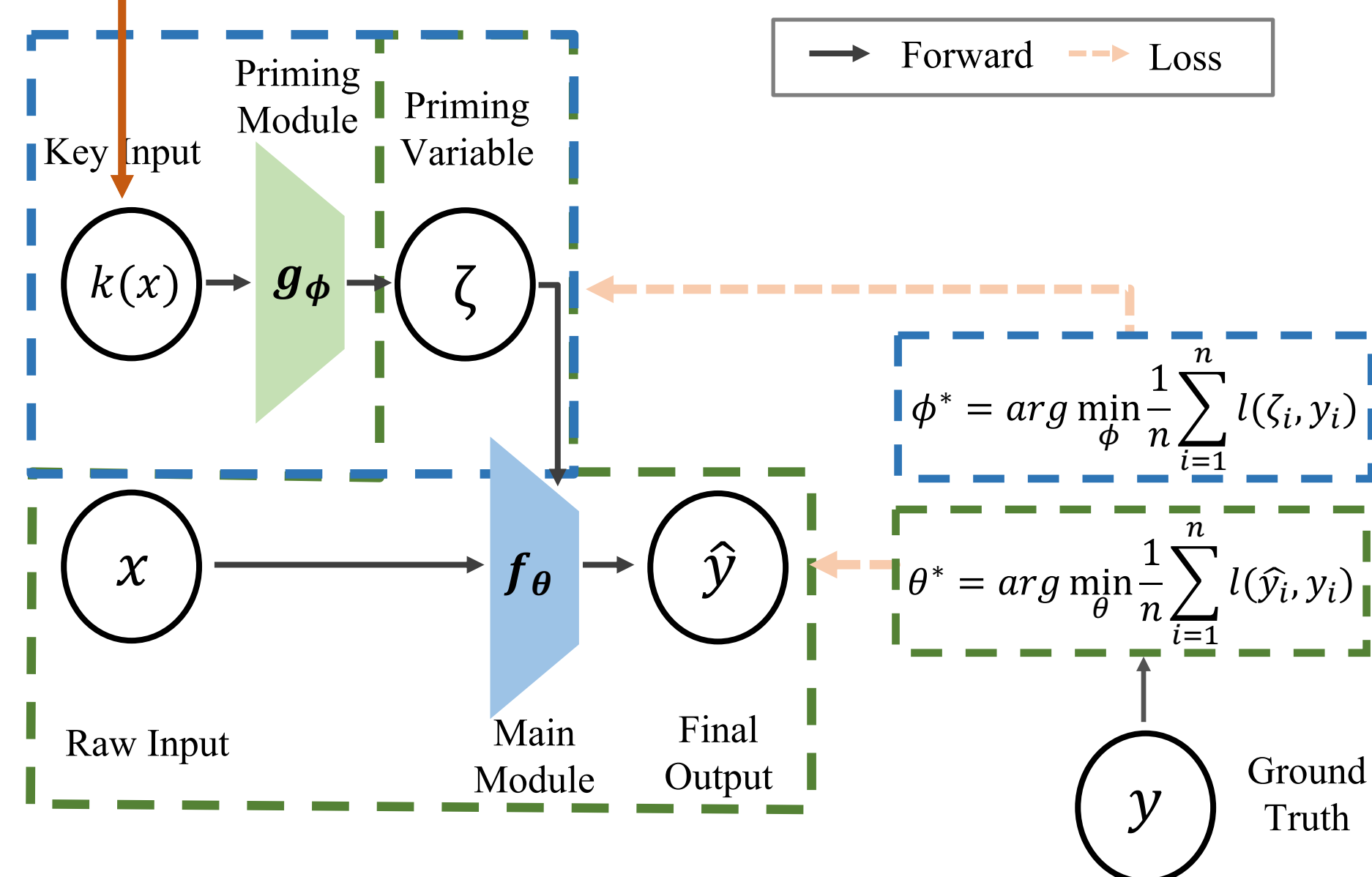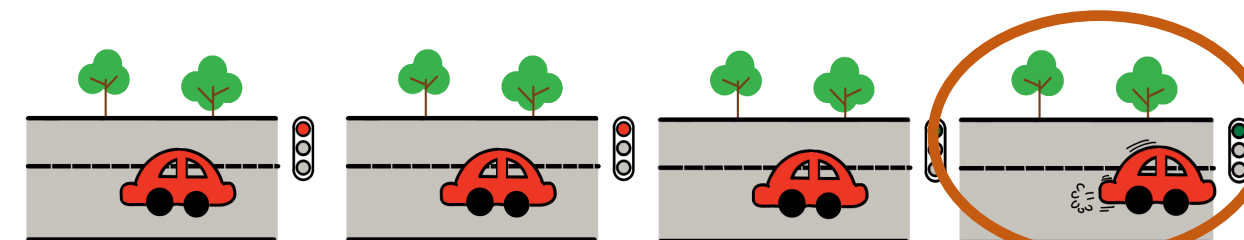
Primed DNN Training

## PrimeNet: Using Key Input Knowledge

Providing key input appropriately creates a new desirable shortcut that "primes" the main module towards correct solutions.
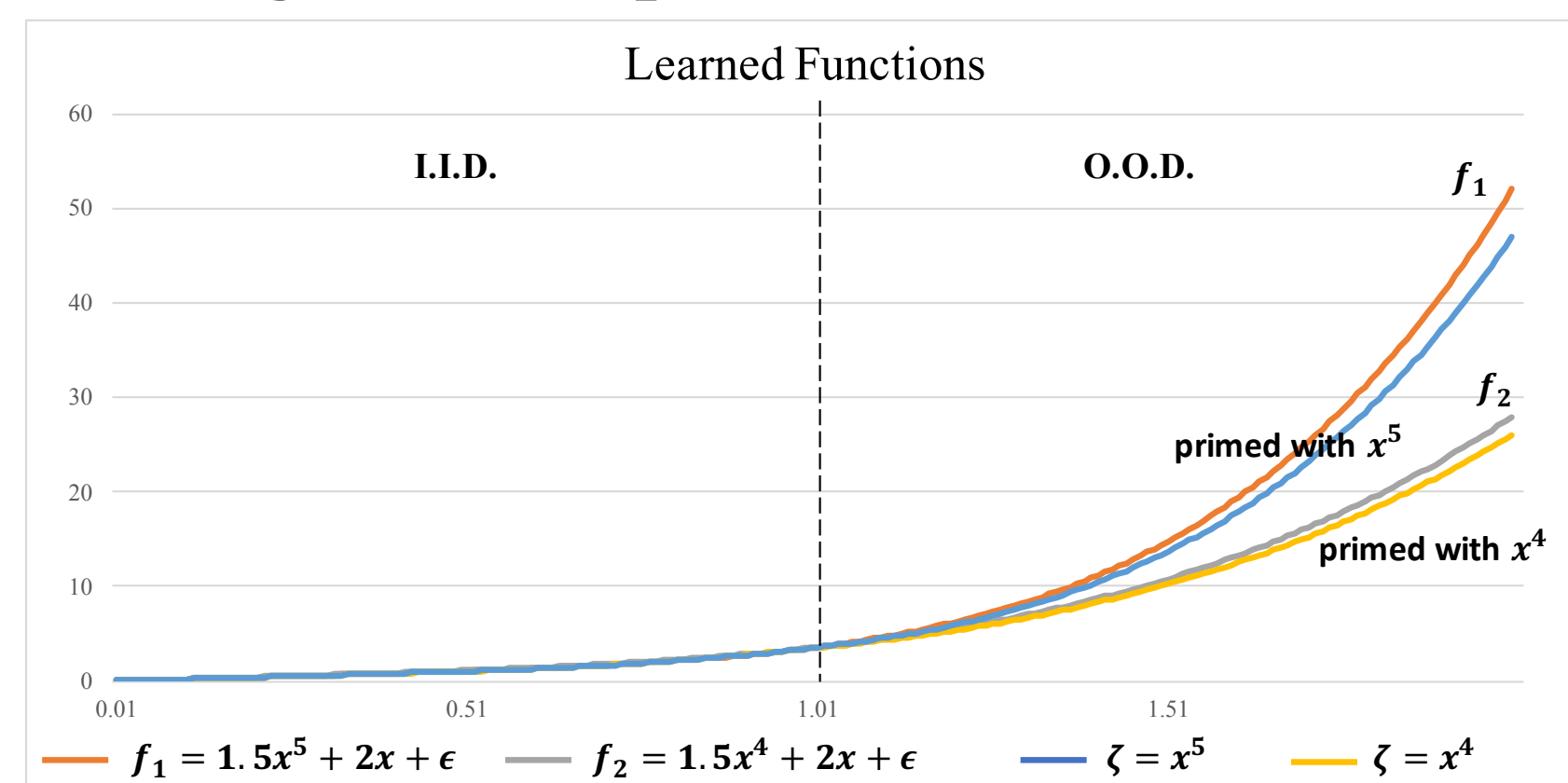
Approximate key inputs are easy to provide

**Classification:** saliency          **Behavioral Cloning:** last frame

→ Forward   --→ Loss

Key Input   Priming Module   Priming Variable

$k(x)$   $g_\phi$   $\zeta$

$x$   $f_\theta$   $\hat{y}$

Raw Input   Main Module   Final Output

$y$   Ground Truth

$$\phi^* = arg\min_\phi \frac{1}{n}\sum_{i=1}^{n} l(\zeta_i, y_i)$$

$$\theta^* = arg\min_\theta \frac{1}{n}\sum_{i=1}^{n} l(\hat{y}_i, y_i)$$

## Experiment Results & Analysis

### Toy 1-D regression experiment:

Learned Functions

I.I.D. | O.O.D.

$f_1$
primed with $x^5$
$f_2$
primed with $x^4$

— $f_1 = 1.5x^5 + 2x + \epsilon$   — $f_2 = 1.5x^4 + 2x + \epsilon$   — $\zeta = x^5$   — $\zeta = x^4$

Conclusion: priming variable $\zeta$ can guide DNN training towards the solution desired by key input.

## Image Classification on NICO:

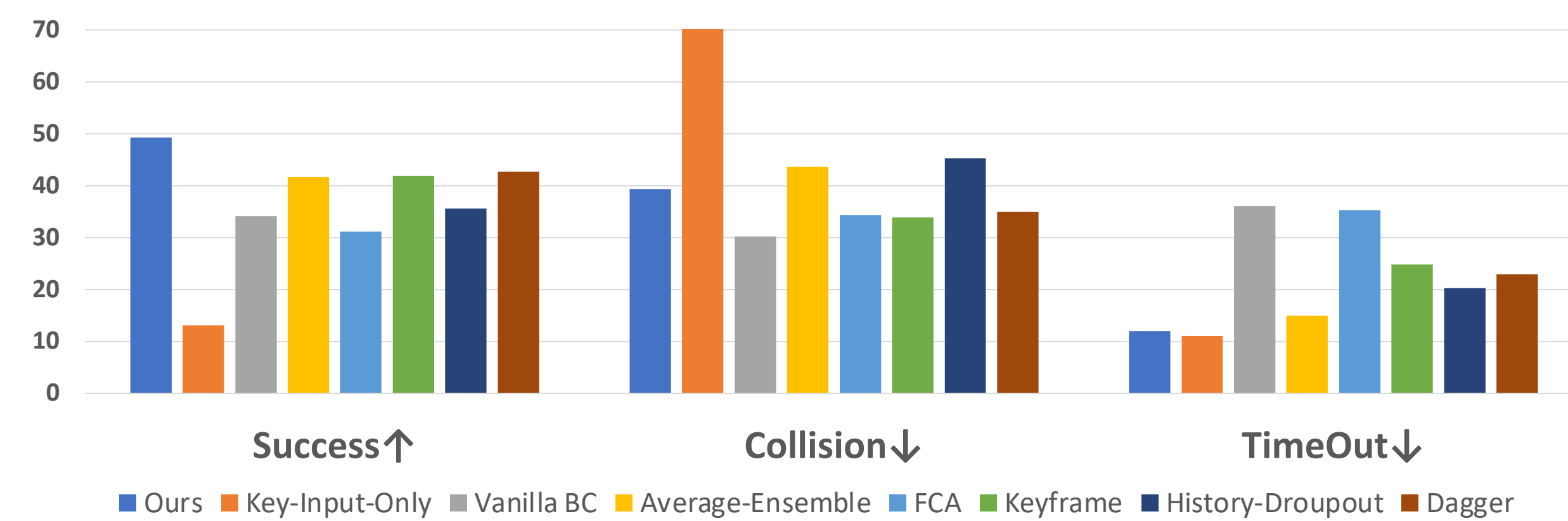| METHOD | IN-DOMAIN TEST | OOD TEST |
|---|---|---|
| VANILLA RESNET18 | 66.11 | 42.61 |
| KEY-INPUT-ONLY | 62.78 | 47.54 |
| AVERAGE-ENSEMBLE | 63.33 | 47.69 |
| RUBI (CADENE ET AL., 2019) | - | 44.37 |
| REBIAS (BAHNG ET AL., 2020) | - | 45.23 |
| CUTOUT (DEVRIES & TAYLOR, 2017) | - | 43.77 |
| MIXUP (ZHANG ET AL., 2017) | 62.78 | 41.46 |
| IRM (ARJOVSKY ET AL., 2019) | - | 41.46 |
| STABLENET (ZHANG ET AL., 2021B) | 63.33 | 43.62 |
| CAAM (WANG ET AL., 2021B) | 70.00 | 46.62 |
| PRIMENET (OURS) | 71.11 | **49.00** |

In-domain context

OOD context

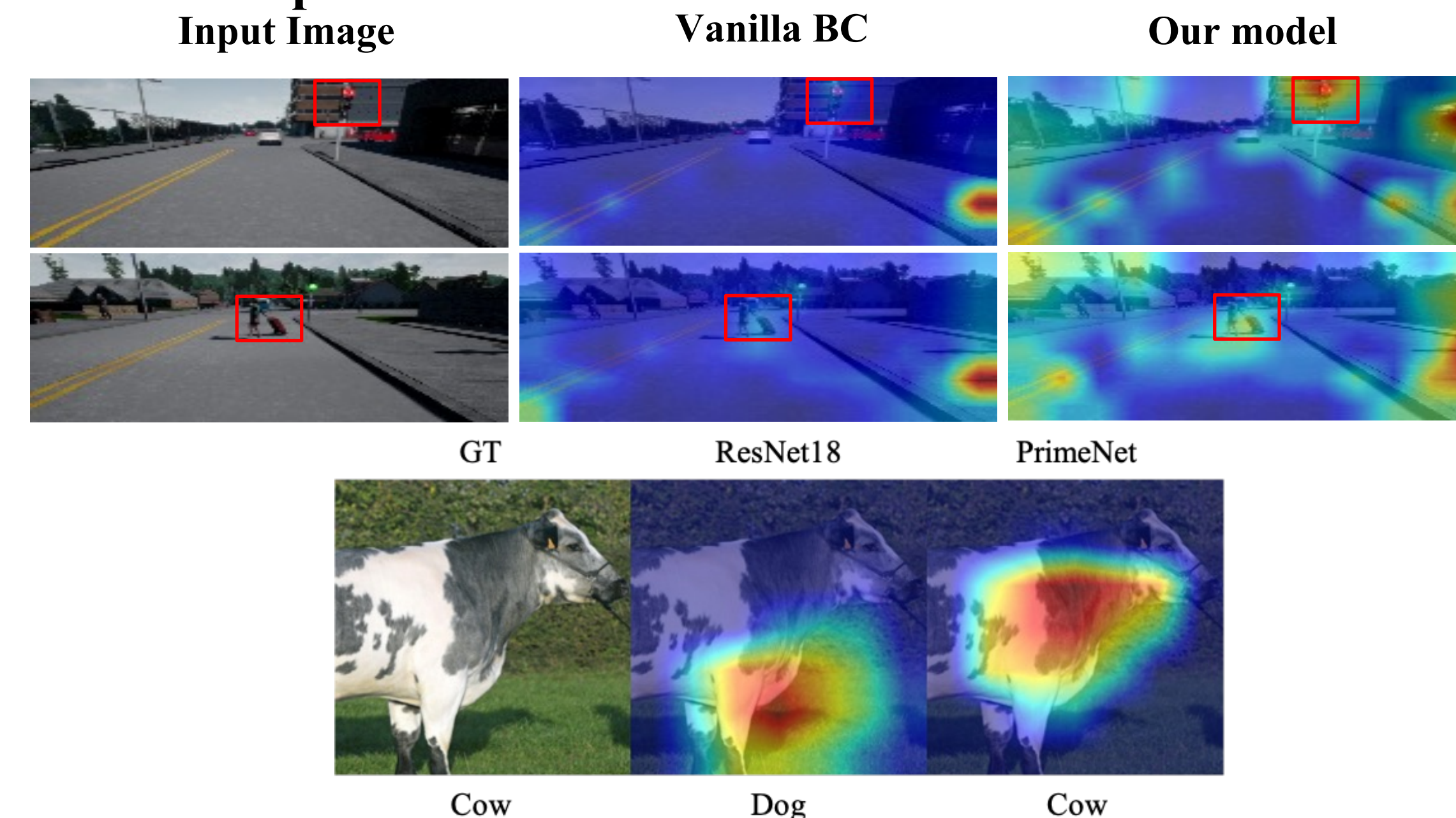PrimeNet generalizes best, without harming in-domain performance.

## Imitation Learning on CARLA Driving:

**CARLA NoCrash Benchmark Results**

Success↑ | Collision↓ | TimeOut↓

Ours | Key-Input-Only | Vanilla BC | Average-Ensemble | FCA | Keyframe | History-Dropout | Dagger

PrimeNet performs best, correctly responding to new environment cues such as traffic lights, neighboring cars, pedestrians.

## Activation Maps:

Input Image | Vanilla BC | Our model

GT | ResNet18 | PrimeNet

Cow | Dog | Cow

Our model attends to the appropriate visual cues in the scene.