# Fighting Fire with Fire: Avoiding DNN Shortcuts through Priming

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

# Shortcuts in DNNs

- DNNs tends to take the shortcut solutions rather than the intended ones.
  - Shortcuts: simple; work well in training distribution; fail in out-of-distribution region.
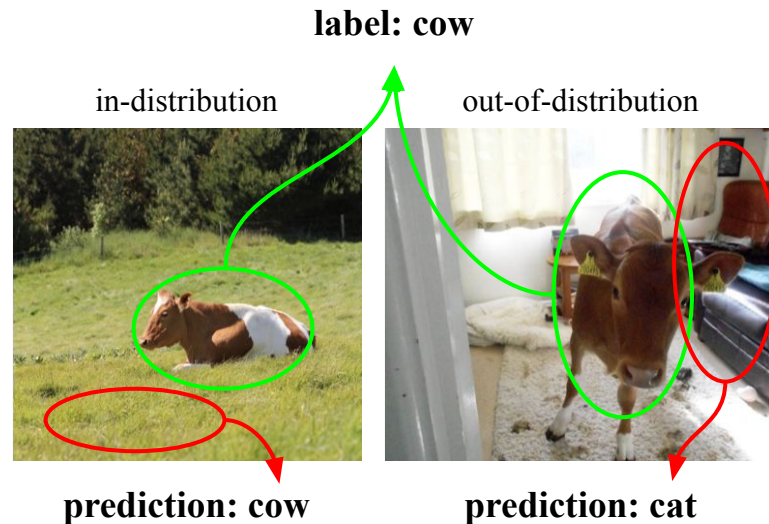  - Learn the suprious correlation and hurt the generalization performance.

labeled dataset (x,y)

DNN Training

learned hypothesis

# Examples: Image Classification

DNNs tend to classify the images according to the **contextual** features rather than the discrimitive **contents**.

Content: foreground object

Context: background scene

Shortcut: context → label



label: cow

in-distribution          out-of-distribution
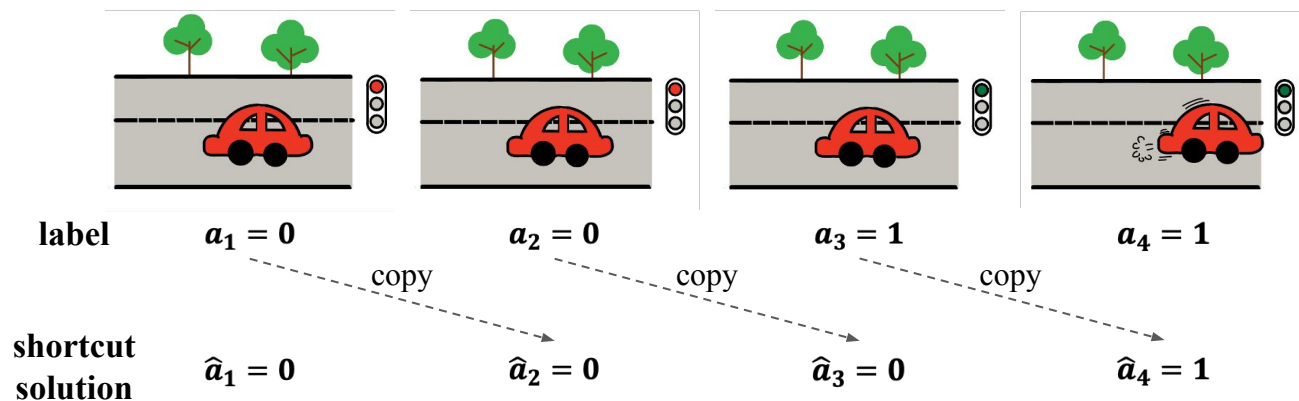
prediction: cow          prediction: cat

Beery, Sara, Grant Van Horn, and Pietro Perona. "Recognition in terra incognita." ECCV 2018.
Wang, Tan, et al. "Causal attention for unbiased visual recognition." ICCV 2021.

# Examples: Imitation Learning

DNNs are prone to simply copy the previous action rather than learn the complex decision policy from the observations.

Shortcut: previous action → current action



Codevilla, Felipe, et al. "Exploring the limitations of behavior cloning for autonomous driving." ICCV 2019.
Wen, Chuan, et al. "Fighting copycat agents in behavioral cloning from observation histories." NeurIPS 2020.
Wen, Chuan, et al. "Keyframe-Focused Visual Imitation Learning." ICML 2021.

# How Do Human Avoid Shortcuts?

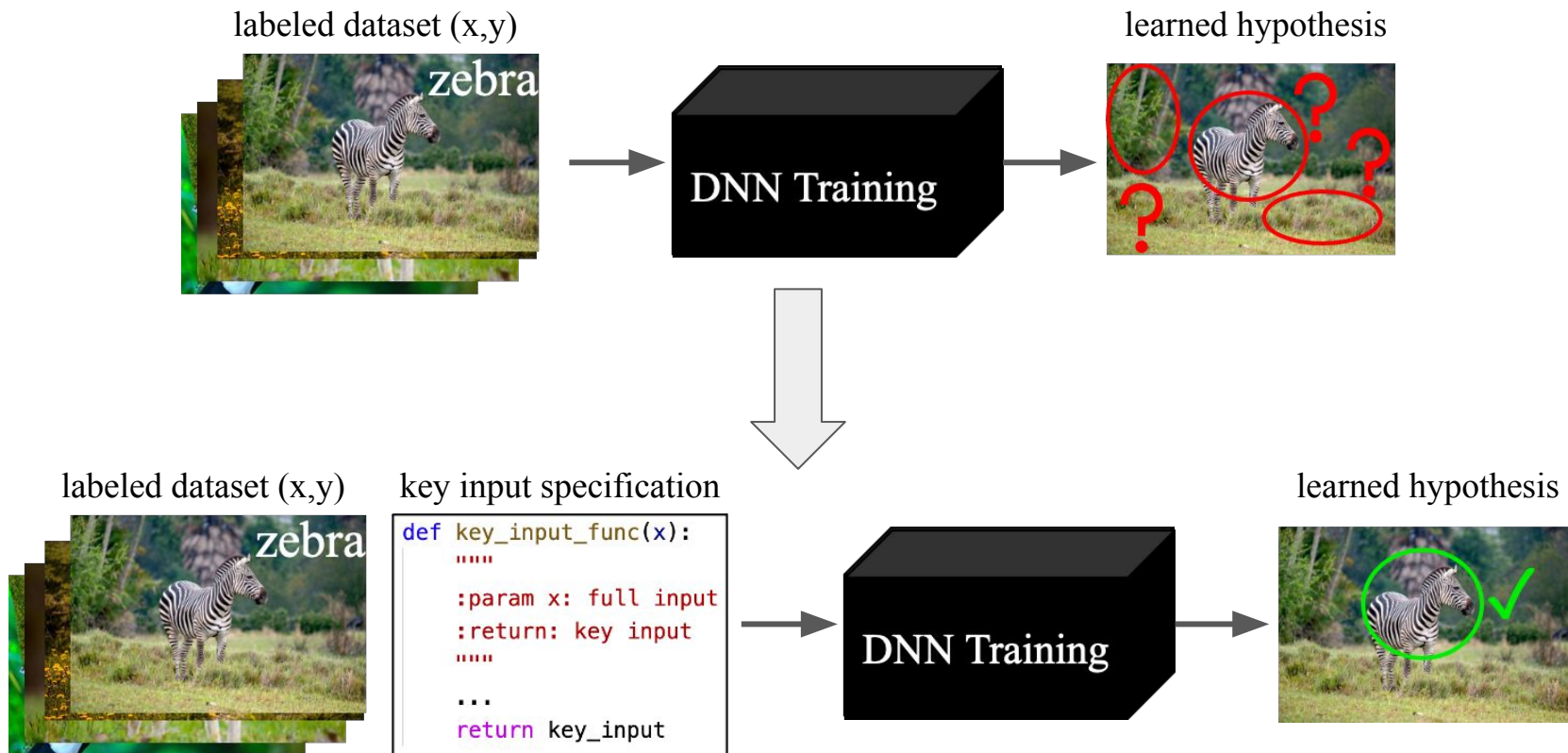Two critical components in Human learning process:

- (1) labels

- (2) domain knowledge about which part of the
  input signal is key to the task

However, in supervised learning, (2) is missing.



**We propose to integrate such auxiliary knowledge into DNNs, to "prime" them away from shortcuts.**

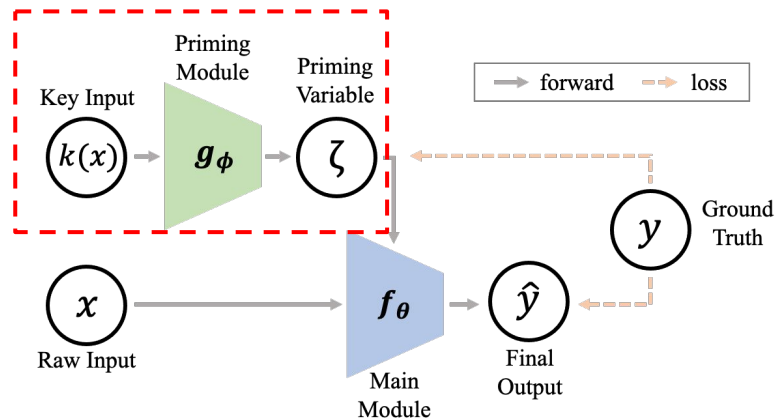# PrimeNet: Prime DNNs Away from Shortcuts

# PrimeNet: Prime DNNs Away from Shortcuts

- Architecture of PrimeNet:
  - Priming Module
    - Key Input: $k(x)$
    - Priming Variable: $\zeta = g_\phi(k(x))$

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

# PrimeNet: Prime DNNs Away from Shortcuts
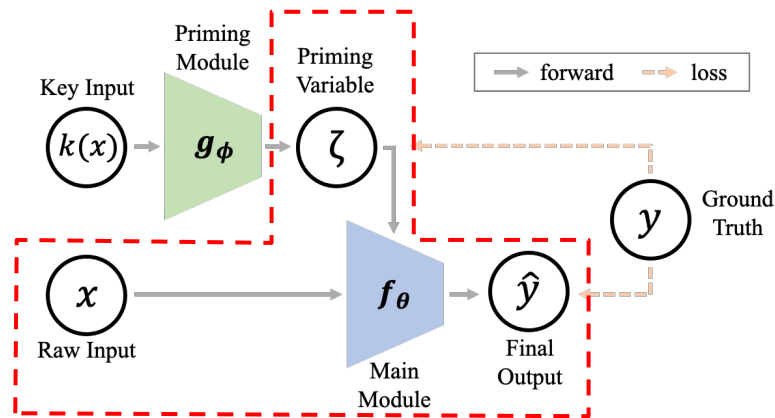
- Architecture of PrimeNet:
  - Priming Module
    - Key Input: $k(x)$
    - Priming Variable: $\zeta = g_\phi(k(x))$

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

  - Main Module

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

# PrimeNet: Prime DNNs Away from Shortcuts

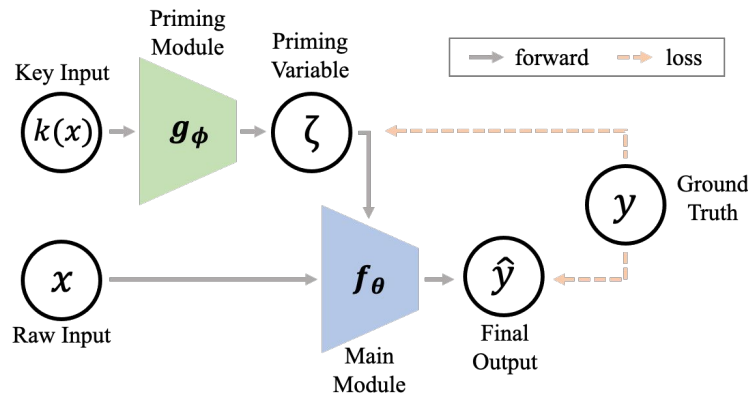- Architecture of PrimeNet:

  - Priming Module
    - Key Input: $k(x)$
    - Priming Variable: $\zeta = g_\phi(k(x))$

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

  - Main Module

    $$\theta^* = \arg\min_{\theta} \frac{1}{n} \sum_{i=1}^{n} l(f_\theta(x_i, \zeta_i), y_i)$$
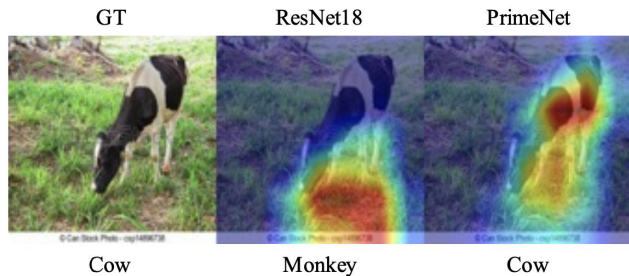
- Defination of Key Input:
  - Image Classification: image patch crop from unsupervised saliency detection.
  - Imitaiton Learning: the most recent frame.
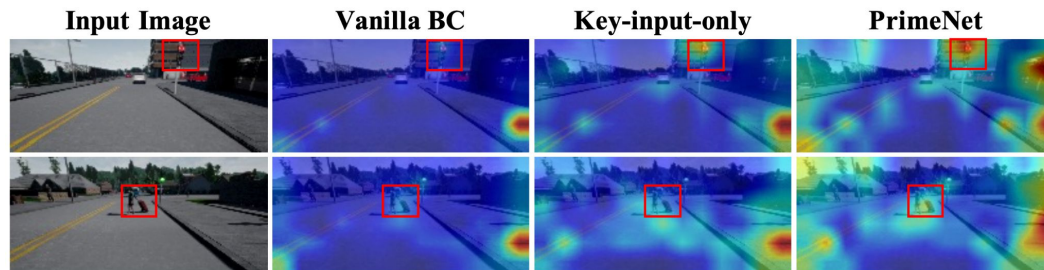
# Experimental Results

## Image Classification: 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** |



| GT | ResNet18 | PrimeNet |
|---|---|---|
| Cow | Monkey | Cow |

## Imitation Learning: CARLA, MuJoCo

| METHOD | CARLA RESULTS | | MUJOCO REWARDS | | |
| | %SUCCESS | #TIMEOUT | HOPPER | ANT | HALFCHEETAH |
|---|---|---|---|---|---|
| VANILLA BC | 34.1 ± 7.5 | 36.1 ± 14.5 | 628 ± 99 | 2922 ± 1266 | 639 ± 121 |
| KEY-INPUT-ONLY | 13.1 ± 1.8 | **11.1 ± 2.9** | 589 ± 94 | 4198 ± 433 | 489 ± 77 |
| AVERAGE-ENSEMBLE | 41.7 ± 3.1 | 15.0 ± 0.8 | 504 ± 47 | 4659 ± 396 | 729 ± 50 |
| **PRIMENET (OURS)** | **49.3 ± 3.6** | 12.0 ± 1.9 | **1124 ± 135** | **4798 ± 304** | **1448 ± 74** |
| FCA (WEN ET AL., 2020) | 31.2 ± 5.2 | 35.3 ± 9.6 | 831 ± 108 | 3727 ± 926 | 1148 ± 81 |
| KEYFRAME (WEN ET AL., 2021) | 41.9 ± 6.2 | 24.8 ± 7.9 | 696 ± 28 | 2930 ± 1321 | 1062 ± 127 |
| HISTORY-DROPOUT (BANSAL ET AL., 2019) | 35.6 ± 3.5 | 20.3 ± 5.6 | 539 ± 33 | 4069 ± 517 | 1215 ± 70 |
| DAGGER (ROSS ET AL., 2011) | 42.7 ± 5.7 | 23.0 ± 7.1 | 2383 ± 294 | 4097 ± 418 | 1842 ± 10 |



Input Image    Vanilla BC    Key-input-only    PrimeNet

# Experimental Results

- Visualization



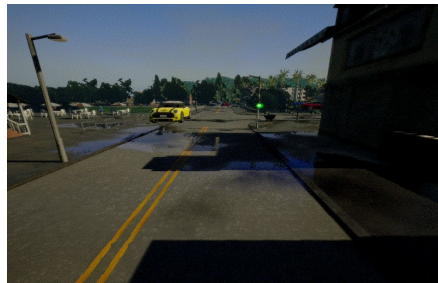| | Shortcut Policy | PrimeNet | |
|---|---|---|---|
| copy the previous action: acceleration | | | Successfully stop |
| copy the previous action: keeping still | | | Successfully start up |