# BEARS Make Neuro-Symbolic Models Aware of their Reasoning Shortcuts
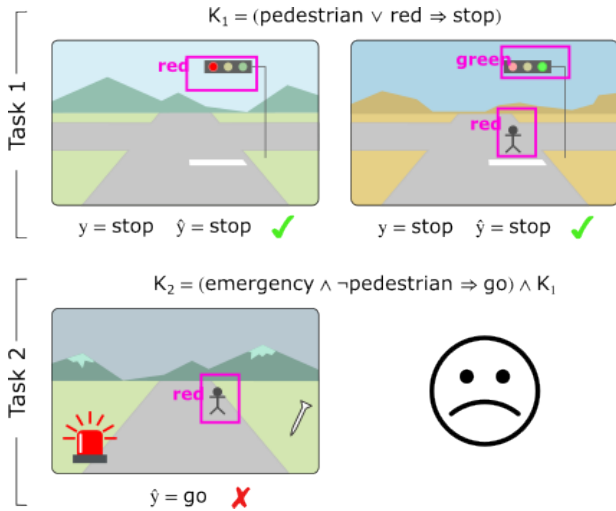
Emanuele Marconato [1,2]   **Samuele Bortolotti** [1]   Emile van Krieken [3]
Antonio Vergari [3]   Andrea Passerini [1]   Stefano Teso [1]

[1]University of Trento   [2]University of Pisa   [3]University of Edinburgh

## REASONING SHORTCUTS

NeSy predictors such as **DeepProbLog**[1], and **Logic Tensor Networks**[2], acquire concepts that comply with the knowledge.

*Are learned concepts interpretable and is the model trustworthy?* **Not always!**[3]



$K_1 = (\text{pedestrian} \vee \text{red} \Rightarrow \text{stop})$

Task 1

$y = \text{stop}$  $\hat{y} = \text{stop}$ ✓    $y = \text{stop}$  $\hat{y} = \text{stop}$ ✓

$K_2 = (\text{emergency} \wedge \neg\text{pedestrian} \Rightarrow \text{go}) \wedge K_1$

Task 2

$\hat{y} = \text{go}$ ✗

**Reasoning Shortcuts (RSs) like this might affect any NeSy predictor!**

### MITIGATION STRATEGIES

| STRATEGY | REQUIRES |
|---|---|
| **Multi-Task** | tasks |
| **Concept Sup.** | concepts |
| **Reconstruction** | (decoder) |
| **Disentanglement** | structure |

### DESIDERATA

- 🟩 **Concept calibration**
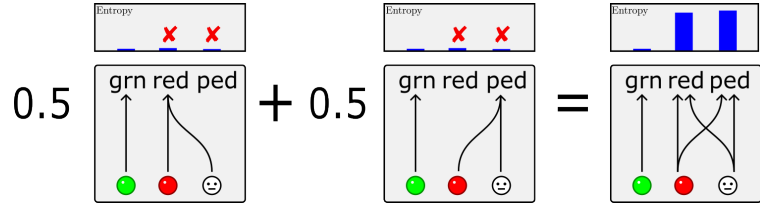- 🟫 **Performance**
- 🟥 **Cost effectiveness**

### 🐻 BEARS: BE AWARE OF REASONING SHORTCUTS! 🐻

**Effective mitigation** strategies for RSs, like concept supervision, are often **impractical**. If the model learns a RS what concepts can we trust?

**Over-confident** solutions are dangerous: impossible to be aware of wrong concepts!

$\{go, stop\}$   $K = (\text{pedestrian} \vee \text{red} \Rightarrow \text{stop})$

K

grn, ped, ..., red

NN

x

Entropy   ✗ ✗

grn red ped

We propose **bears** to estimate concept uncertainty!

## OUR SOLUTION



$0.5$ (Entropy ✗ ✗ grn red ped) $+$ $0.5$ (Entropy ✗ ✗ grn red ped) $=$ (Entropy grn red ped)

**bears** combines **Deep Ensembles + diversification** ($\sim$ Bayesian NeSy) and provably optimizes for all desiderata:
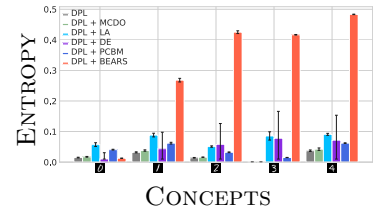
$$\mathcal{L}_{\text{bears}} = \mathcal{L}(\mathbf{x}, \mathbf{y}; \mathsf{K}, \theta_t)$$
$$+ \gamma_1 \cdot \text{KL}\left(p_{\theta_t}(\mathbf{C} \mid \mathbf{x}) \,\|\, \frac{1}{t}\sum_{j=1}^{t} p_{\theta_j}(\mathbf{C} \mid \mathbf{x})\right)$$
$$+ \gamma_2 \cdot H(p_{\theta_t}(\mathbf{C} \mid \mathbf{x}))$$
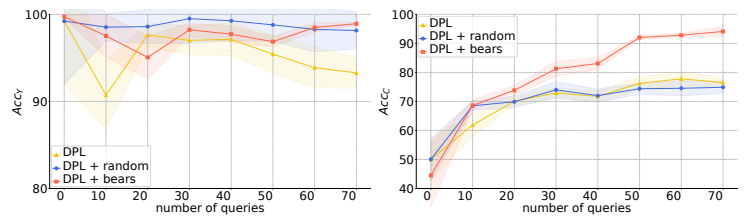
## EXPERIMENTS

① **An example from** `MNIST-Addition`

Solve the sum between two digits, *e.g.*, 2 + 3 = 5.

$\begin{cases} 0 + 0 = 0 \\ 0 + 1 = 1 \\ 2 + 3 = 5 \\ 2 + 4 = 6 \end{cases}$



② **Active learning with** `bears`



③ `bears` **in real-world:** BDD-OIA [4]

| | $\text{mECE}_C$ | $\text{ECE}_C(F,S)$ | $\text{ECE}_C(R)$ | $\text{ECE}_C(L)$ |
|---|---|---|---|---|
| DPL | $0.84 \pm 0.01$ | $0.75 \pm 0.17$ | $0.79 \pm 0.05$ | $0.59 \pm 0.32$ |
| + MCDO | $0.83 \pm 0.01$ | $0.72 \pm 0.19$ | $0.76 \pm 0.08$ | $0.55 \pm 0.33$ |
| + LA | $0.85 \pm 0.01$ | $0.84 \pm 0.10$ | $0.87 \pm 0.04$ | $0.67 \pm 0.19$ |
| + PCBM | $0.68 \pm 0.01$ | $0.26 \pm 0.01$ | $0.26 \pm 0.02$ | $0.11 \pm 0.02$ |
| + DE | $0.79 \pm 0.01$ | $0.62 \pm 0.03$ | $0.71 \pm 0.10$ | $0.37 \pm 0.12$ |
| + bears | $\mathbf{0.58 \pm 0.01}$ | $\mathbf{0.14 \pm 0.01}$ | $\mathbf{0.10 \pm 0.01}$ | $\mathbf{0.02 \pm 0.01}$ |

## REFERENCES

[1] Manhaeve *et al.*, DeepProbLog, NeurIPS (2018)
[2] Donadello *et al.*, Logic Tensor Networks, IEEE (2018)
[3] Marconato et al., Not All Neuro-Symbolic Concepts are Created Equal: Analysis and Mitigation of Reasoning Shortcuts, NeurIPS (2023)
[4] Xu *et al.*, BDD-OIA dataset, CVPR (2020).

PAPER   CODE