# A Neuro-Symbolic Benchmark Suite for Concept Quality and Reasoning Shortcuts

**S. Bortolotti**[1]  **E. Marconato**[1,2]  T. Carraro[4,5]  P. Morettin[1]  E. v. Krieken[3]  A. Vergari[3]  S. Teso[1]  A. Passerini[1]

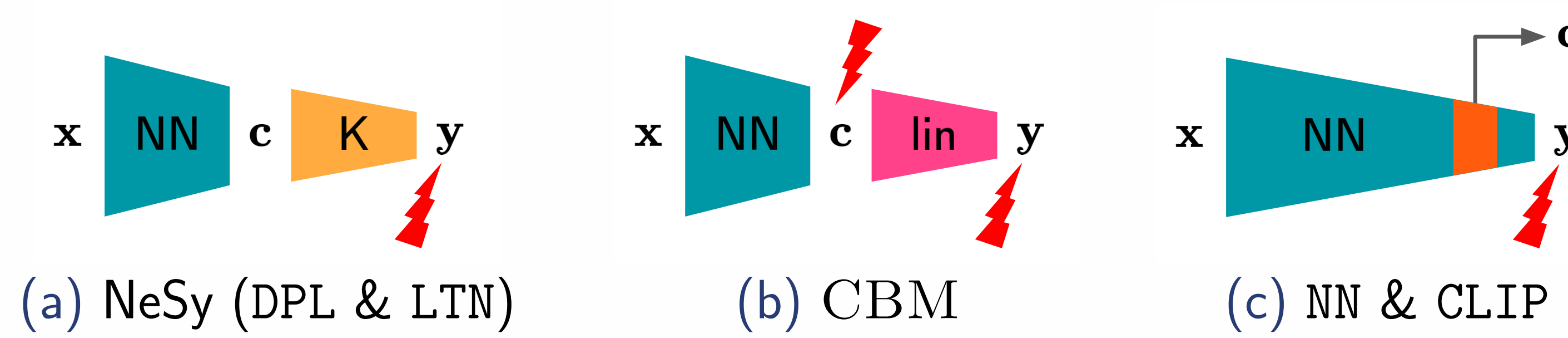[1]University of Trento   [2]University of Pisa   [3]University of Edinburgh   [4]Fondazione Bruno Kessler   [5]University of Padova
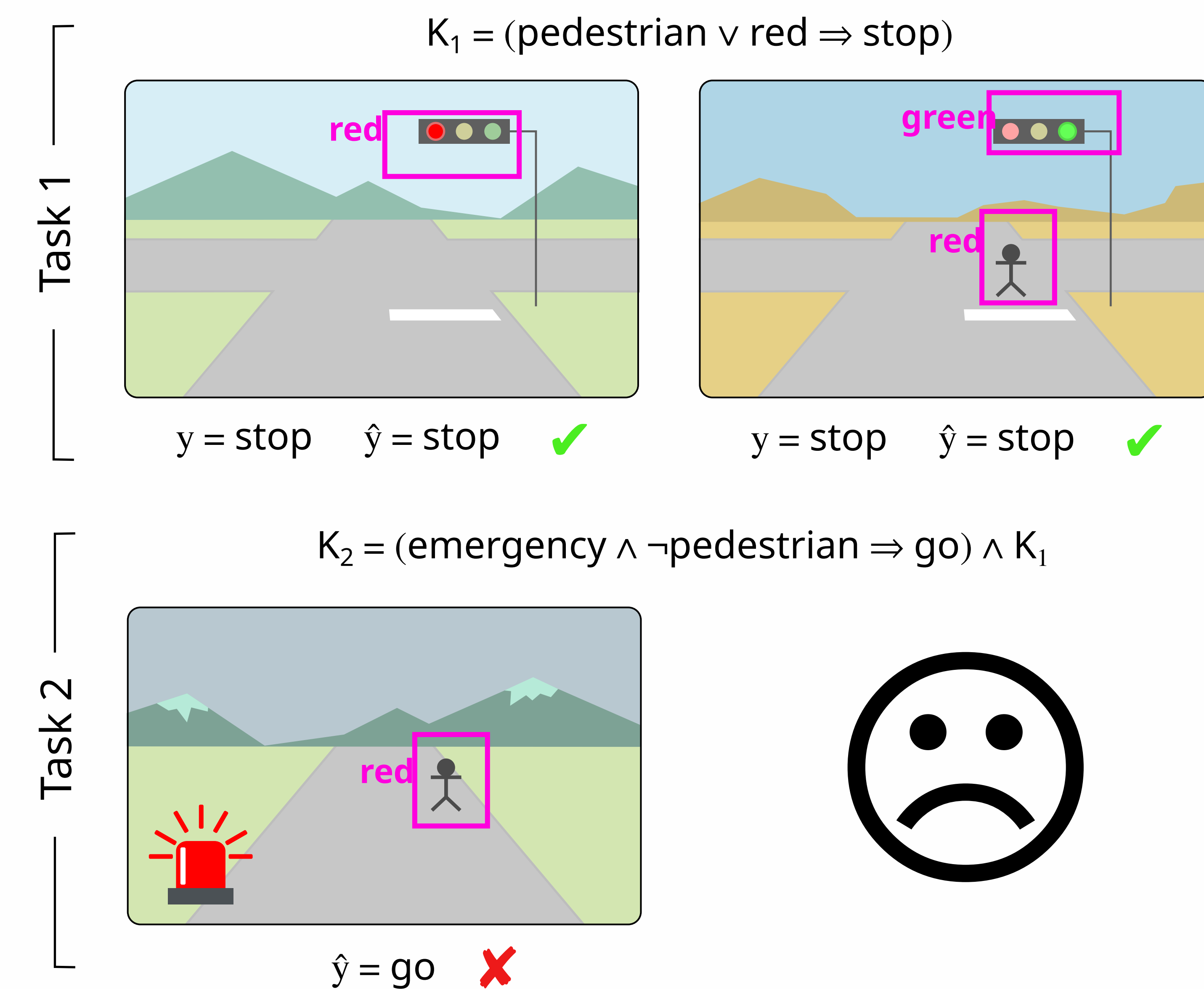
PAPER   CODE

## REASONING SHORTCUTS

**Goal:** Study supervised models that classify samples *correctly* but for the *wrong concepts*.



(a) NeSy (DPL & LTN)   (b) CBM   (c) NN & CLIP

**Reasoning Shorcuts** [1]: **NeSy predictors** [2], **Concept-based Models** [3] and **VLMs** like CLIP [4] solve Learning & Reasoning tasks by exploiting semantically misleading concepts.

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

Task 1



y = stop   ŷ = stop ✓     y = stop   ŷ = stop ✓

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

Task 2



ŷ = go ✗

## L&R TASKS

| Task | Data | | | Properties | | |
|------|------|------|------|------|------|------|
| | Gen | OOD | Cont | Cplx x | Cplx K | Amb K |
| MNMath (new) | ✓ | ✓ | ✓ | ✗ | ✓ | ✗ |
| +/× MNAdd-Half | ✗ | ✓✓ | ✗ | ✗ | ✗ | – |
| MNAdd-EvenOdd | ✗ | ✓✓ | ✓✓ | ✗ | ✗ | – |
| MNLogic (new) | ✓ | ✓ | ✓ | ✗ | ✓ | ✗ |
| ∧/∨ Kand-Logic | ✓ | ✓ | ✓ | ✗ | ✗ | ✓ |
| CLE4EVR | ✓ | ✓ | ✓✓ | ✓ | ✗ | ✓ |
| ⚠ BDD-OIA | ✗ | ✗ | ✗ | ✓ | ✓ | ✓ |
| SDD-OIA (new) | ✓ | ✓✓ | ✓ | ✓ | ✓ | ✓ |

## EXAMPLES

| Task | Example | Shortcut | OOD Pred. |
|------|---------|----------|-----------|
| SDD-OIA |  | STOP { 🚶→🔴 , 🟢→🔴 } |  GO |
| BDD-OIA |  | STOP { 🚶→🔴 , 🟢→🔴 } |  GO |

*Knowledge K = the traffic laws.*

| Task | Example | Shortcut | OOD Pred. |
|------|---------|----------|-----------|
| MNMath | $2 \cdot 2 + 2 = 6$ $3 + 4 = 7$ | { 2→2 , 3→4 , 4→3 } | 2 + 4 = 5 |

*Knowledge K = equations must hold.*

| Task | Example | Shortcut | OOD Pred. |
|------|---------|----------|-----------|
| MNLogic [1] | $0 \oplus 1 = 1$ | { 0→1 , 1→0 } | 0 ∧ 0 = 1 |
| Kand-Logic [2] |  = 1 | { □→red , △→yel , ○→blu } |  = 0 |
| CLE4EVR [3] |  = 0  = 1 | { 🟥→🟥 , ⬜→⬜ , ⭕→🟥 } |  = 1 |

1: *Knowledge K = formula must hold.*
2: *Knowledge K = pattern must hold.*
3: *Knowledge K = same color and shape?*

## FEATURES

① **Challenging**: the # of RSs can be chosen **a priori** and counted using `countrss`, allows to control task difficulty.

② **Configurable**: data sets & generators can be easily configured with YAML/JSON files.

③ **Intuitive**: straightforward to use:

```
from rsbench import MNLOGIC

dataset = MNLOGIC(args)
train(model, dataset)
test(model, dataset)
```

NOTEBOOK

## ASSESSING RS

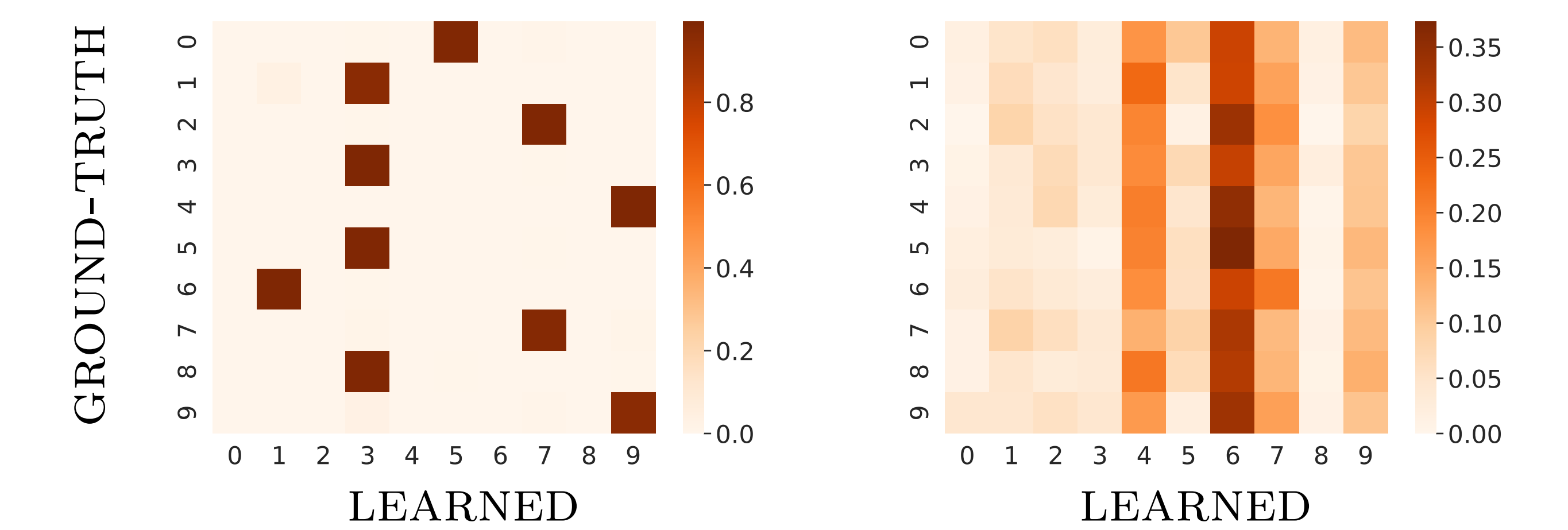① **Task-level**: `countrss` counts the # of potential RSs in any L&R task!

$$K : Y \leftrightarrow (C_1 \land C_2 \land C_3) \qquad \mathcal{D} : \{(0,1,0), (1,0,0)\}$$

countrss

$\#RSs$

**Example**: *with 3 concepts and an exhaustive training set, MNLogic has 6 RSs if K is a conjunction and 24 if K is a XOR. This grows **exponentially** with the # of concepts!*

② **Model-level**: `rsbench` tasks induce RSs in all models!

Table 1. (L) DPL and (R) NN concept confusion matrix on MNAdd-EvenOdd



**Quantitatively**: Concept F1, accuracy and *collapse*

## REFERENCES

[1] Marconato *et al.*, Analysis and Mitigation of RSs, NeurIPS (2023)
[2] Manhaeve *et al.*, DeepProblog, NeurIPS (2018)
[3] Pang Wei Koh *et al.*, Concept bottleneck models, ICML (2020)
[4] Alec Radford *et al.*, CLIP, ICML (2021)