



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

**Samuele Bortolotti**

*samuele.bortolotti@unitn.it*

**Emanuele Marconato**

*emanuele.marconato@unitn.it*

**Tommaso Carraro**

*tcarraro@fbk.eu*

**Paolo Morettin**

*paolo.morettin@unitn.it*

**Emile van Krieken**

*Emile.van.Krieken@ed.ac.uk*

**Antonio Vergari**

*avergari@ed.ac.uk*

**Stefano Teso**

*stefano.teso@unitn.it*

**Andrea Passerini**

*andrea.passerini@unitn.it*

## Setting

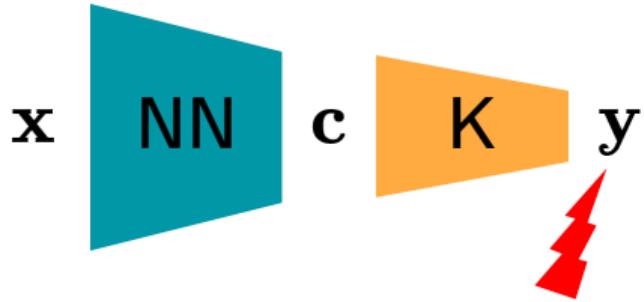


Figure: Neuro-Symbolic model: DeepProbLog (DPL) [1] & Logic Tensor Networks (LTN) [2]

[1] Manhaeve et al., DeepProbLog: Neural Probabilistic Logic Programming, NeurIPS (2018)

[2] Donadello et al., Logic Tensor Networks, IEEE (2018)

## *Setting*

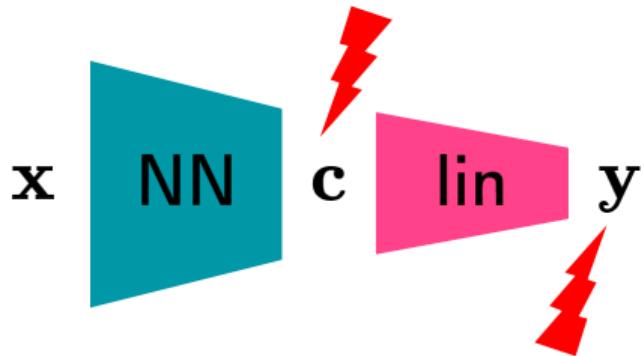


Figure: Concept bottleneck models (CBM) [3]

[3] Pang Wei Koh *et al.*, Concept bottleneck models, ICML (2020)

## *Setting*

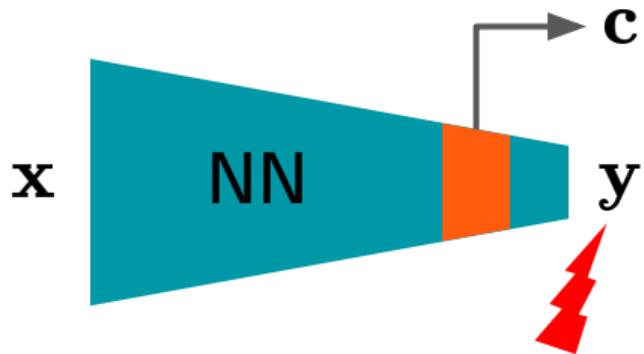
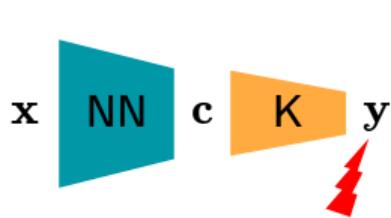


Figure: Neural Network (NN) & CLIP [4]

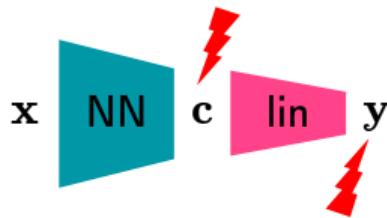
[4] Alec Radford *et al.*, Learning Transferable Visual Models From Natural Language Supervision, ICML (2021)

# Setting

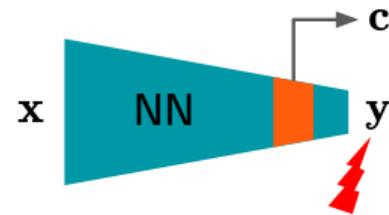
**Goal:** Study supervised models that classify samples *correctly* but for the *wrong concepts*.



(a) NeSy (DPL [1] & LTN [2])



(b) CBM [3]



(c) NN & CLIP [4]

[1] Manhaeve et al., DeepProbLog: Neural Probabilistic Logic Programming, NeurIPS (2018)

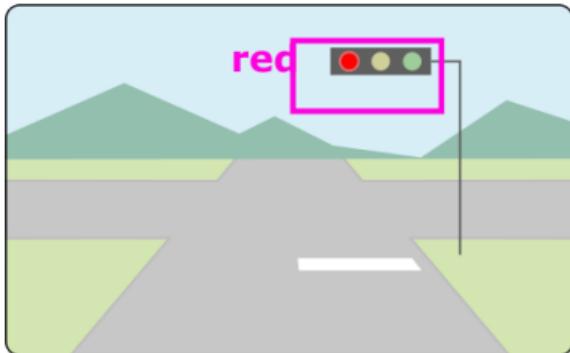
[2] Donadello et al., Logic Tensor Networks, IEEE (2018)

[3] Pang Wei Koh et al., Concept bottleneck models, ICML (2020)

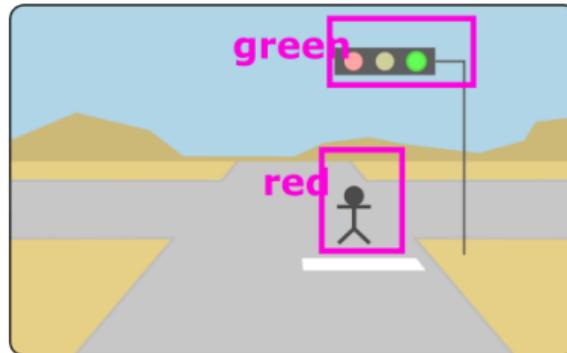
[4] Alec Radford et al., Learning Transferable Visual Models From Natural Language Supervision, ICML (2021)

## Reasoning Shortcuts

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



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

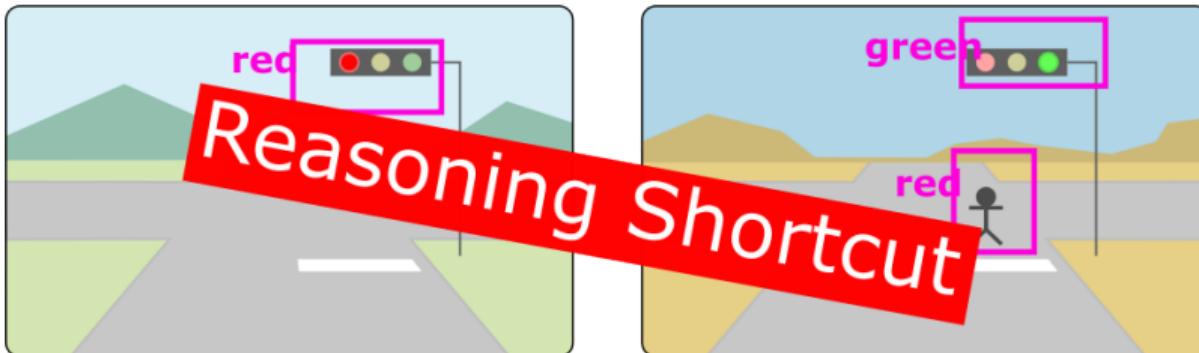


$y = \text{stop}$     $\hat{y} = \text{stop}$    ✓

- Task: predict stop vs.go using concepts “pedestrian”, “red”, and “green”.

## Reasoning Shortcuts

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



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

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

- Task: predict stop vs.go using concepts "pedestrian", "red", and "green".

Perfect accuracy by predicting pedestrians as red lights!

## rsbench: **L&R tasks**

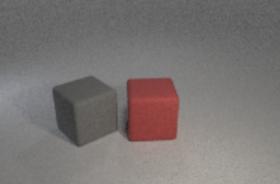
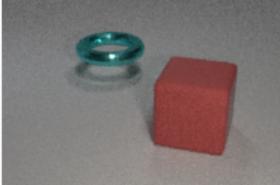
	Task	Data			Properties		
		Gen	OOD	ConL	Cplx x	Cplx K	Amb K
Arithmetic	MNMath ( <u>new</u> )	✓	✓	✓	✗	✓	✗
	MNAdd-Half	✗	✓	✗	✗	✗	-
	MNAdd-EvenOdd	✗	✓	✓	✗	✗	-
Logic	MNLogic ( <u>new</u> )	✓	✓	✓	✗	✓	✗
	Kand-Logic	✓	✓	✓	✗	✓	✓
	CLE4EVR	✓	✓	✓	✓	✗	✓
High Stakes	BDD-OIA	✗	✗	✗	✓	✓	✓
	SDD-OIA ( <u>new</u> )	✓	✓	✓	✓	✓	✓

## rsbench: *examples*

Task	Example	Shortcut	OOD Pred.
SDD-OIA	 = STOP	$\left\{ \begin{array}{l} \text{person} \rightarrow \text{red} \\ \text{green} \rightarrow \text{green} \end{array} \right.$	 = GO

*Knowledge K = the traffic laws.*

## rsbench: *examples*

Task	Example	Shortcut	OOD Pred.
CLE4EVR		$= 0$ $\left\{ \begin{array}{l} \text{■} \rightarrow \text{■} \\ \text{■} \rightarrow \text{■} \\ \text{○} \rightarrow \text{■} \end{array} \right.$	 = 1
		$= 1$	

*Knowledge K = same color and shape?*

## rsbench: *features*

■ **Challenging:** require complex perception and/or reasoning.

■ **Versatile:** supports NeSy models, CBMs, post-hoc explainers.

■ **Configurable:** can be easily configured with YAML/JSON files.

■ **Intuitive:** straightforward to use:

```
from rsbench import MNLOGIC
```

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

■ **Model-level metrics:**

- ▶ F1 and Accuracy
- ▶ Concept level confusion matrix
- ▶ Concept Collapse

■ **Task-level metrics:**

- ▶ Formally counts the # of potential RSs in any L&R task!

# If you care about your concepts, come to our poster!



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

S. Bortolotti<sup>1</sup> E. Marconato<sup>1,2</sup> T. Carraro<sup>4,5</sup> P. Morettin<sup>1</sup> E. V. Krieken<sup>3</sup> A. Vergari<sup>3</sup> S. Teso<sup>1</sup> A. Passerini<sup>1</sup>

<sup>1</sup>University of Trento

<sup>2</sup>University of Pisa

<sup>3</sup>University of Edinburgh

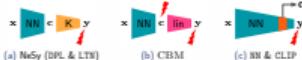
<sup>4</sup>Fondazione Bruno Kessler

<sup>5</sup>University of Padova

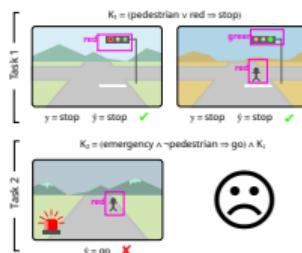


### REASONING SHORTCUTS

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



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



### L&R TASKS

TASK	DATA		PROPERTIES					
	GEN	OOD	CONST	CPLX	K	AMB	K	RS
MNMath	✓	✓	✓	✓	✓	✓	✓	✗
+/ $\times$ MNAdd-Half	✗	✗	✗	✗	✗	✗	✗	✗
MNAdd-EvenDid	✗	✗	✗	✗	✗	✗	✗	✗
MNLogic [2]	✓	✓	✓	✓	✓	✓	✓	✗
A/V Kand-Logic	✓	✓	✓	✓	✓	✓	✓	✓
CLE4EVN	✓	✓	✓	✓	✗	✓	✓	✓
SDG-DIA	✗	✗	✗	✓	✓	✓	✓	✓
SDG-DIA [2]	✗	✗	✗	✓	✓	✓	✓	✓

### FEATURES

- ① **Challenging:** the # of RSs can be chosen *a priori* and counted using `countrs`, 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)
```



### EXAMPLES

TASK	EXAMPLE	SHORTCUT	OOD PRED.
SDG-DIA	STOP	STOP	GO
SDG-DIA	STOP	STOP	GO

Knowledge K – the traffic laws.

TASK	EXAMPLE	SHORTCUT	OOD PRED.
MNMath	$\frac{1}{2}, \frac{1}{2} = 6$	$\frac{1}{2} \rightarrow 2$	✓
MNMath	$\frac{1}{2}, \frac{1}{2} = 7$	$\frac{1}{2} \rightarrow 4$	✓

Knowledge K – equations must hold.

### ASSESSING RS

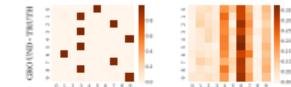
- ④ **Task-level:** `countrs` counts the # of potential RSs in any L&R task!



**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!

- ⑤ **Model-level:** rsbench tasks induce RSs in all models!

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



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

### REFERENCES

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



PAPER



CODE



WEBSITE

<https://unitn-sml.github.io/rsbench/>



samuele.bortolotti samuele.bortolotti@unitn.it