



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

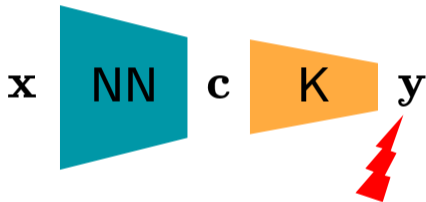


Figure: Neuro-Symbolic model: DeepProbLog (DPL) [1] & Logic Tensor Networks (LTN) [2]

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

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

Setting

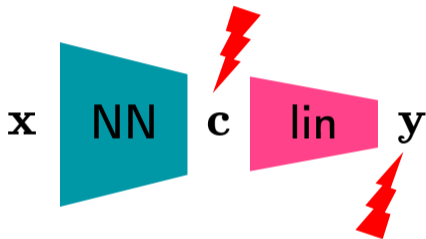


Figure: Concept bottleneck models (CBM) [3]

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

Setting

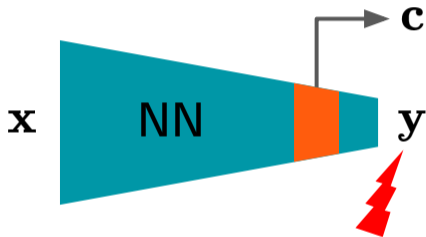
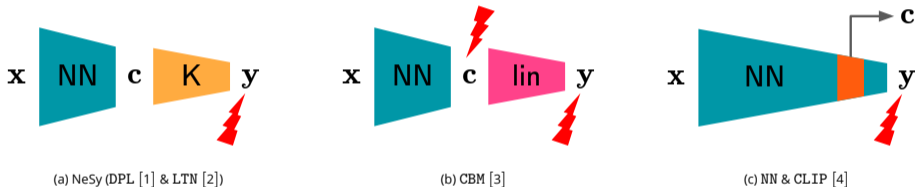


Figure: Neural Network (NN) & CLIP [4]

[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*.



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

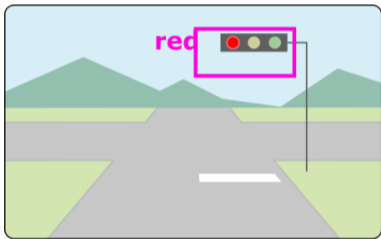
[2] Donadello et al., Logic Tensor Networks, IEEE (2018)

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

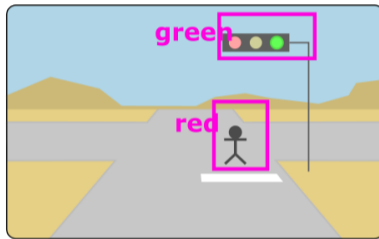
[4] Alec Radford et al., Learning Transferable Visual Models From Natural Language Supervision, ICML (2021)

Reasoning Shortcuts

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



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

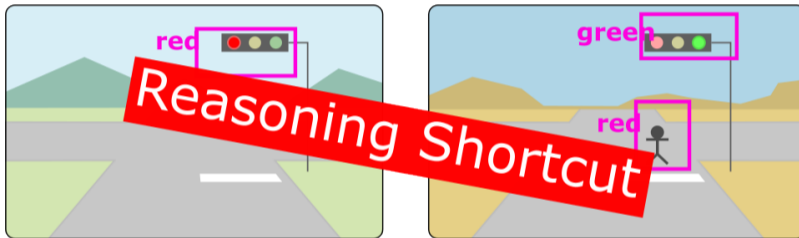


$y = \text{stop}$ $\hat{y} = \text{stop}$ ✓

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

Reasoning Shortcuts

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



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

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


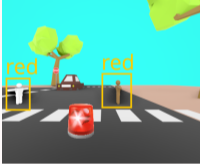
- Task: predict stop vs.go using concepts "pedestrian", "red", and "green".

Perfect accuracy by predicting pedestrians as red lights!

rsbench: *L&R tasks*

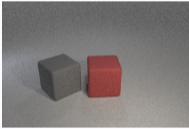
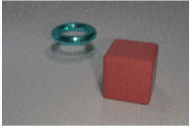
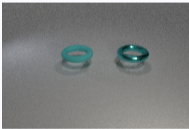
Task		Data			Properties		
		Gen	OOD	ConL	Cplx x	Cplx K	Amb K
Arithmetic	MNMath (<i>new</i>)	✓	✓	✓	✗	✓	✗
	MNAdd-Half	✗	✓✓	✗	✗	✗	-
	MNAdd-EvenOdd	✗	✓✓	✓✓	✗	✗	-
Logic	MNLogic (<i>new</i>)	✓	✓	✓	✗	✓	✗
	Kand-Logic	✓	✓	✓	✗	✓	✓
	CLE4EVR	✓	✓	✓✓	✓	✗	✓
High Stakes	BDD-OIA	✗	✗	✗	✓	✓	✓
	SDD-OIA (<i>new</i>)	✓	✓✓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

rsbench: *examples*

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

Knowledge \mathcal{K} = the traffic laws.

rsbench: *examples*

Task	Example	Shortcut	OOD Pred.
CLE4EVR	 = 0	$\left\{ \begin{array}{l} \color{red}\blacksquare \rightarrow \color{red}\blacksquare \\ \color{gray}\blacksquare \rightarrow \color{gray}\blacksquare \\ \color{blue}\circ \rightarrow \color{red}\blacksquare \end{array} \right.$	 = 1
	 = 1		

Knowledge \mathcal{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. Bortolotti¹ E. Marconato^{1,2} T. Carraro^{4,5} P. Morettin¹ E. v. Krieken³ A. Vergari³ S. Teso¹ A. Passerini¹

¹University of Trento ²University of Pisa ³University of Edinburgh ⁴Fondazione Bruno Kessler ⁵University of Padova



PAPER



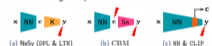
CODE



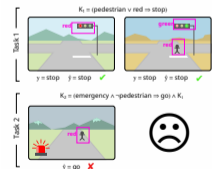
WEBSITE

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	CONV	CPLEX	CPLEX	K	ASS	K
MMath (small)	✓	✓	✓	✓	✓	✓	✓	✓
+> MMAdd-EvenOdd	✓	✓	✓	✓	✓	✓	✓	✓
MMLogic (small)	✓	✓	✓	✓	✓	✓	✓	✓
+> Kand-Logic	✓	✓	✓	✓	✓	✓	✓	✓
CLEAVER	✓	✓	✓	✓	✓	✓	✓	✓
SDD-DIA	✓	✓	✓	✓	✓	✓	✓	✓
SDD-DIA (small)	✓	✓	✓	✓	✓	✓	✓	✓

FEATURES

① **Challenging**: the # of RSs can be chosen a priori and counted using counters, 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 `MMLOGIC`

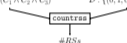
```
dataset = MMLOGIC(args)
train(model, dataset)
test(model, dataset)
```



ASSESSING RS

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

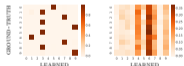
$$K: Y \leftrightarrow (C_1 \wedge C_2 \wedge C_3) \quad \mathcal{D}: \{(0,1,0), (1,0,0)\}$$



Example: with 3 concepts and an exhaustive training set, `MMLogic` 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 `MMAdd-EvenOdd`



Quantitatively: Concept F1, accuracy and collapse

REFERENCES

- Marconato et al., Analysis and Mitigation of RSCs, *NeurIPS* (2023)
- Marconato et al., DeepPruning, *NeurIPS* (2018)
- Fang Wei Koh et al., Concept bottleneck models, *ICML* (2020)
- Alec Radford et al., CLIP, *ICML* (2021)

EXAMPLES

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

Knowledge K = the traffic laws.

TASK	EXAMPLE	SHORTCUT	OOD PRED.
MMMath	$\begin{cases} 2 + 2 = 6 \\ 3 + 3 = 7 \end{cases}$	$\begin{cases} 2 \rightarrow 2 \\ 3 \rightarrow 3 \end{cases}$	$\begin{cases} 2 \rightarrow 1 \\ 3 \rightarrow 5 \end{cases}$

Knowledge K = equations must hold.

TASK	EXAMPLE	SHORTCUT	OOD PRED.
MMLogic ¹	$\begin{cases} \text{A} \rightarrow 1 \\ \text{B} \rightarrow 0 \end{cases}$	$\begin{cases} \text{A} \rightarrow 1 \\ \text{B} \rightarrow 0 \end{cases}$	$\begin{cases} \text{A} \rightarrow 1 \\ \text{B} \rightarrow 1 \end{cases}$
Kand-Logic ²	$\begin{cases} \text{A} \rightarrow 1 \\ \text{B} \rightarrow 1 \end{cases}$	$\begin{cases} \square \rightarrow \text{red} \\ \triangle \rightarrow \text{yel} \\ \circ \rightarrow \text{blu} \end{cases}$	$\begin{cases} \square \rightarrow 1 \\ \triangle \rightarrow 0 \end{cases}$
CLEAVER ³	$\begin{cases} \text{A} \rightarrow 0 \\ \text{B} \rightarrow 1 \end{cases}$	$\begin{cases} \square \rightarrow \text{red} \\ \triangle \rightarrow \text{blu} \end{cases}$	$\begin{cases} \square \rightarrow 1 \\ \triangle \rightarrow 1 \end{cases}$

- Knowledge K = formula must hold
- Knowledge K = pattern must hold
- Knowledge K = same color and shape?