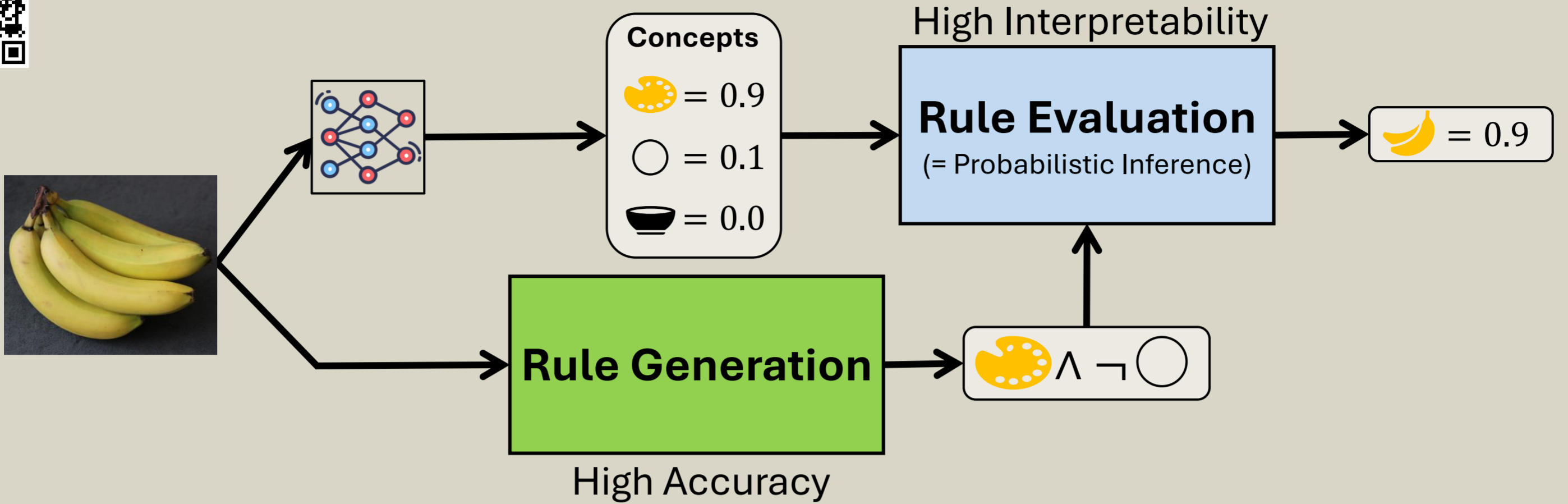


# Neurosymbolic Concept-Based Reasoners Go Beyond the Accuracy-Interpretability Trade-Off of Concept Bottleneck Models

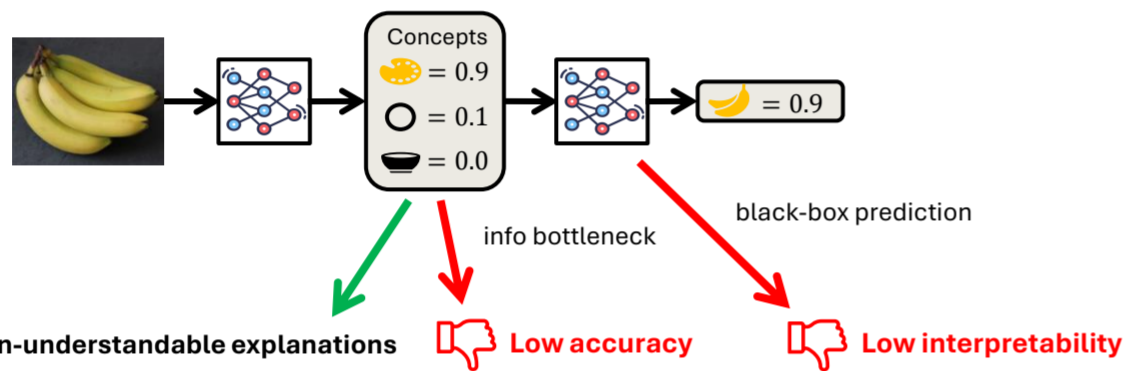
David Debot and Giuseppe Marra

Integrating Local and Global Interpretability for Deep Concept-Based Reasoning Models



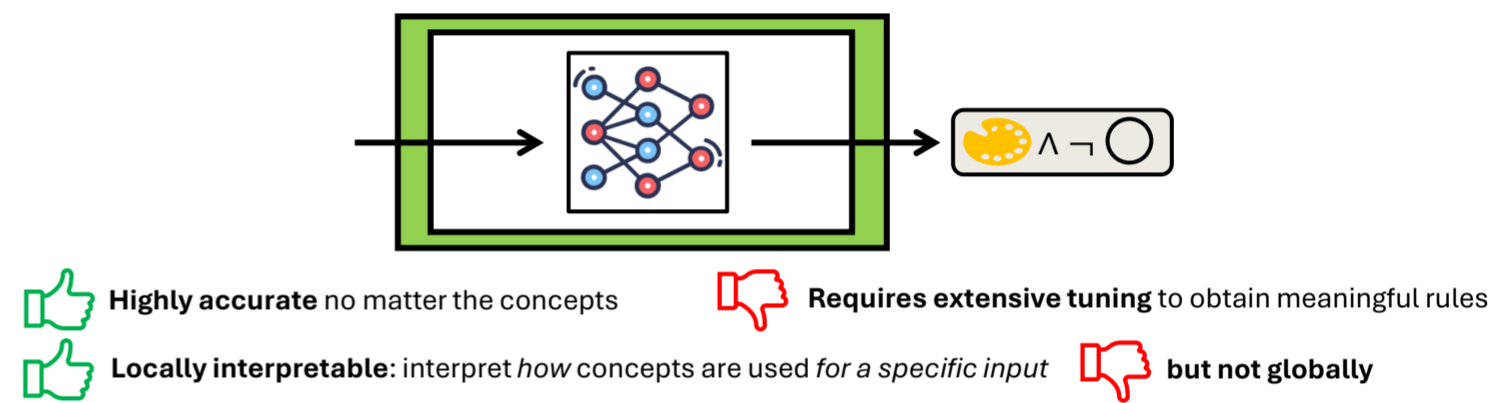
## Concept Bottleneck Models [1]

CBMs = intrinsically explainable models that first predict **concepts** and then predict a downstream task with them  
= high-level, human-understandable features related to the task



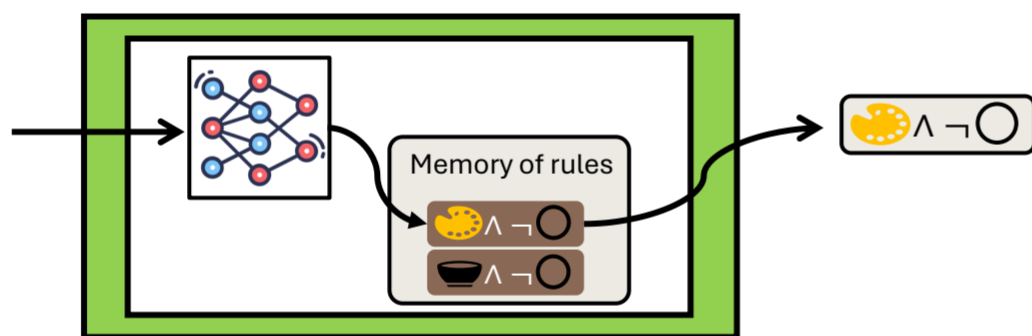
## Deep Concept Reasoner [2]

Rule generation = neural prediction of a rule



## Concept-based Memory Reasoner [3]

Rule generation = neural selection in a learned memory of rules



CMR is highly accurate no matter the employed concepts

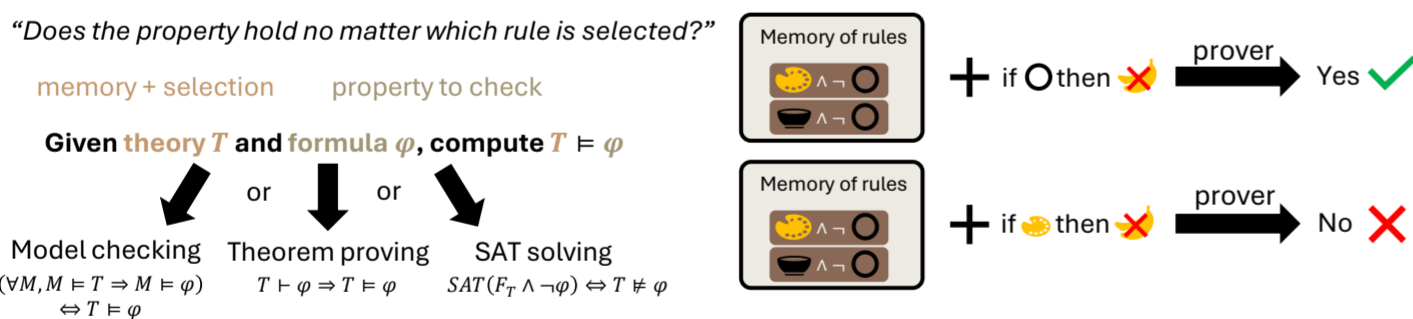
Theorem: CMR is a universal binary approximator if  $n_r \geq 3$

	MNIST+	MNIST+*	CELEBA	CEBAB
Best CBM	97.41 ± 0.55	77.63 ± 0.44	50.24 ± 0.34	83.80 ± 0.01
Black box	83.26 ± 8.71	83.26 ± 8.71	65.72 ± 0.70	88.67 ± 0.19
CMR	97.52 ± 0.30	95.47 ± 0.47	63.56 ± 0.48	85.14 ± 0.43

ACCURACY vs NUMBER OF CONCEPTS (1, 12, 24, 37). Legend: Best CBM (orange), CMR (red), Black box (grey).

CMR's global interpretability allows verification of properties

All decision rules in memory are transparent  $\Rightarrow$  model properties can be verified before deployment

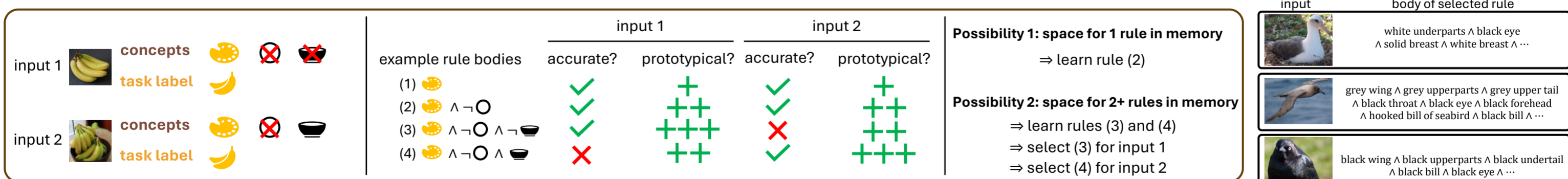


CMR's rule learning allows human interaction during training

The way CMR learns rules allows for human interaction in multiple ways = "rule interventions"

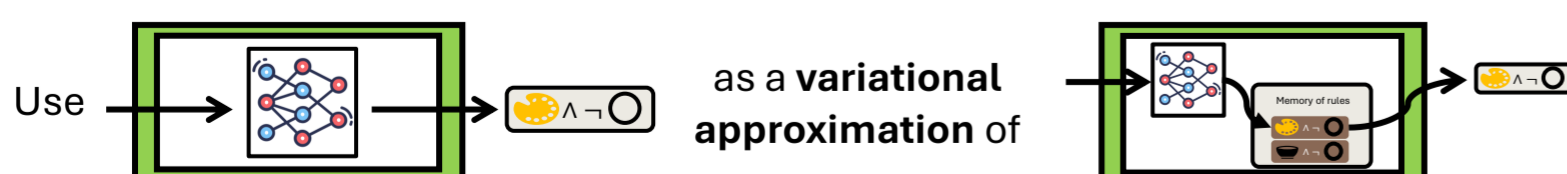
Rule intervention	Can be used for	Example
Add rules manually to the memory	Incorporating expert knowledge	Add $\text{Banana} \leftarrow \text{Banana} \wedge \text{Bowl}$
Forbid a concept from being in a rule	Debiasing	Avoid using $\text{Bowl}$ in rule 1
Force a concept to be in a rule	Enforcing safety	Use $\text{Banana}$ in all rules

CMR learns meaningful rules that are both accurate and prototypical of concept activations



Highly accurate no matter the concepts Globally interpretable: all decision rules transparent Prototypes as meaningful rules Difficult optimization problem

## Unified Concept Reasoner [4]



Training objective: 2 components

- Apply CMR's objective on DCR's rule generation (i.e. accurate + prototypical)
- Maximize rule correspondence between DCR's and CMR's rule generation (KL divergence)

	Rule correspondence
Positive class	98.9 ± 0.3
Negative class	74.2 ± 2.0
	MNIST+ accuracy
UCR (DCR head)	97.8 ± 0.2
UCR (CMR head)	95.7 ± 0.8

Highly accurate no matter the concepts Globally interpretable Prototypes as meaningful rules Easier optimization

## References

- Concept Bottleneck Models Koh et al.
- Interpretable Neural-Symbolic Concept Reasoning Barbiero et al.
- Interpretable Concept-Based Memory Reasoning Debot et al.
- Integrating Local and Global Interpretability for Deep Concept-Based Reasoning Models Debot et al.