

***Provably reliable***

***Deep learning***

***by design***

**antonio vergari** (he/him)

 @tetraduzione

joint work with Kareem Ahmed, Andreas Grivas, Lorenzo Loconte, Stefano Teso, Kai-Wei Chang,  
Guy Van den Broeck, Nicola di Mauro, Robert Peharz, Adam Lopez

27th Oct 2023 - **Workshop on Safe and Robust Machine Learning**

# *april*

*april is  
probably a  
recursive  
identifier of a  
lab*

# *april*

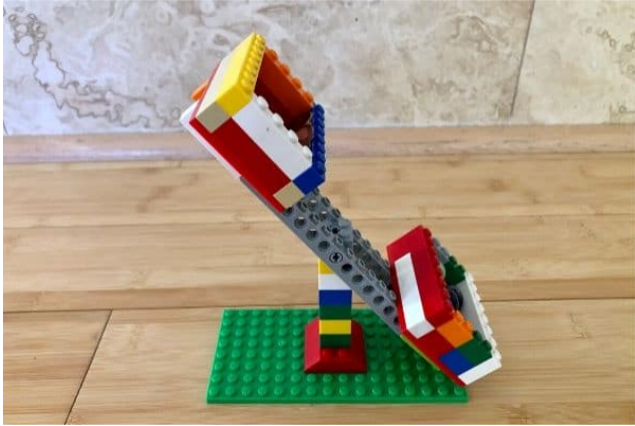
*autonomous &  
provably  
reliable  
intelligent  
learners*



***deep learning is differentiable lego***



***just stacking blocks can be unreliable though...***



*especially when dealing with **hard constraints***

**part I** how to satisfy wanted constraints  
in NNs by design

**part II** how the design of NNs implicitly poses  
unwanted constraints

**part I** how to satisfy wanted constraints  
in NNs by design

**part II** how the design of NNs implicitly poses  
unwanted constraints



**part I** how to satisfy wanted constraints  
in NNs by design

**Why?**

***“How can neural nets reason and learn with symbolic constraints reliably and efficiently?”***

## Why?

*“How can neural nets reason and learn with **symbolic** constraints reliably and efficiently?”*

integrate **hard (logical)** constraints

## Why?

**“How can neural nets  
reason and learn with  
symbolic constraints  
*reliably* and *efficiently*?”**

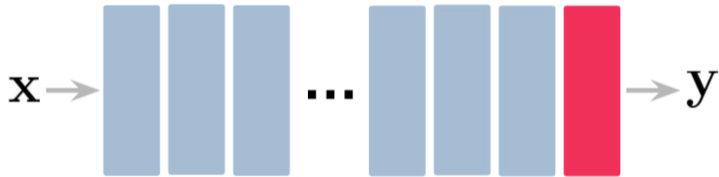
*guarantee* that predictions *always satisfy* constraints

## Why?

*“How can neural nets reason and learn with symbolic constraints reliably and **efficiently**?”*

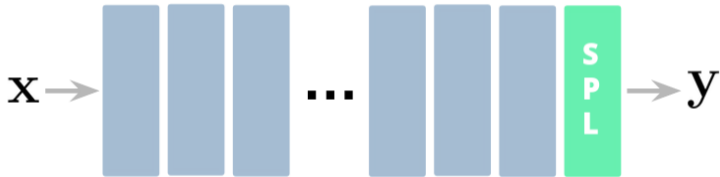
*fast* and *exact* gradients

**Why?**



***make any neural network architecture...***

# Why?



*...guarantee all predictions to conform to constraints*

# When?



Ground Truth

***e.g. predict shortest path in a map***



# When?



given  $x$  // e.g. a tile map

Ground Truth

***nesy structured output prediction (SOP) tasks***

# When?



Ground Truth

given  $\mathbf{x}$  // e.g. a tile map

find  $\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} p_{\theta}(\mathbf{y} \mid \mathbf{x})$  // e.g. a configurations of edges in a grid

## **nesy structured output prediction (SOP) tasks**

## When?



Ground Truth

given  $\mathbf{x}$  // e.g. a tile map

find  $\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} p_{\theta}(\mathbf{y} \mid \mathbf{x})$  // e.g. a configurations of edges in a grid  
s.t.  $\mathbf{y} \models \mathbf{K}$  // e.g., that form a valid path

## nesy structured output prediction (SOP) tasks

## When?



Ground Truth

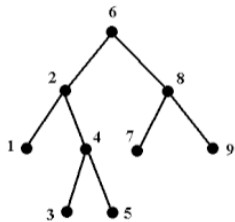
given  $\mathbf{x}$  // e.g. a tile map

find  $\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} p_{\theta}(\mathbf{y} \mid \mathbf{x})$  // e.g. a configurations of edges in a grid  
s.t.  $\mathbf{y} \models \mathbf{K}$  // e.g., that form a valid path

// for a  $12 \times 12$  grid,  $2^{144}$  states but only  $10^{10}$  valid ones!

## nesy structured output prediction (SOP) tasks

# When?



given  $\mathbf{x}$  // e.g. a feature map

find  $\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} p_{\theta}(\mathbf{y} \mid \mathbf{x})$  // e.g. labels of classes

s.t.  $\mathbf{y} \models \mathbf{K}$  // e.g., constraints over superclasses

$$\mathbf{K} : (Y_{\text{cat}} \implies Y_{\text{animal}}) \wedge (Y_{\text{dog}} \implies Y_{\text{animal}})$$

## hierarchical multi-label classification

# When?



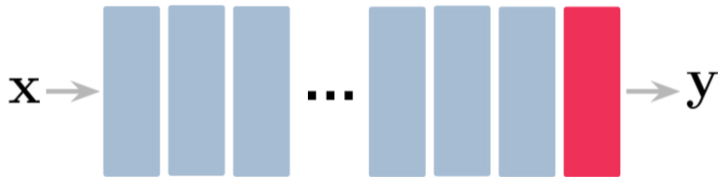
given  $\mathbf{x}$  // e.g. a user preference over  $K - N$  sushi types  
find  $\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} p_{\theta}(\mathbf{y} \mid \mathbf{x})$  // e.g. prefs over  $N$  more types  
s.t.  $\mathbf{y} \models \mathbf{K}$  // e.g., output valid rankings

## ***user preference learning***

---

Choi, Van den Broeck, and Darwiche, "Tractable learning for structured probability spaces: A case study in learning preference distributions", *IJCAI*, 2015

*e.g.,*



***sigmoid linear layers***

$$p(\mathbf{y} | \mathbf{x}) = \prod_{i=1}^N p(y_i | \mathbf{x})$$

# When?



Ground Truth



ResNet-18

***neural nets struggle to satisfy validity constraints!***



## ***the issues!***

- I) Logical constraints can be hard to represent in a unified way
  - ⇒ ***a single framework*** for matching, paths, hierarchies, ...
  
- II) How to integrate logic and probabilities in a single neural layer
  - ⇒ *combining soft and hard constraints*
  
- III) Logical constraints are piecewise constant functions!
  - ⇒ *differentiable almost everywhere but **gradient is zero!***

## ***the issues!***

- I) Logical constraints can be hard to represent in a unified way  
⇒ ***a single framework*** for matching, paths, hierarchies, ...
  
- II) How to integrate logic and probabilities in a single neural layer  
⇒ ***combining soft and hard constraints***
  
- III) Logical constraints are piecewise constant functions!  
⇒ *differentiable almost everywhere but **gradient is zero!***

## ***the issues!***

- I) Logical constraints can be hard to represent in a unified way  
⇒ ***a single framework*** for matching, paths, hierarchies, ...
  
- II) How to integrate logic and probabilities in a single neural layer  
⇒ ***combining soft and hard constraints***
  
- III) Logical constraints are piecewise constant functions!  
⇒ *differentiable almost everywhere but **gradient is zero!***

## ***Constraint losses***

$$\mathcal{L}(\theta; \mathbf{x}, \mathbf{y}) + \lambda \mathcal{L}_K(\mathbf{x}, \mathbf{y})$$

***losses improve consistency during training...***

## Constraint losses

$$\mathcal{L}(\theta; \mathbf{x}, \mathbf{y}) + \lambda \mathcal{L}_K(\mathbf{x}, \mathbf{y})$$

**losses improve consistency during training...**

e.g., the **semantic loss**:  $\mathcal{L}_{SL} := -\log \sum_{\mathbf{y} \models K} \prod_i p(Y_i | \mathbf{x})$

## ***Constraint losses***



Ground Truth



ResNet-18



Semantic Loss

***...but cannot guarantee consistency at test time!***

# what?

DESIDERATUM	LOSSES			LAYERS				
	DL2 [29]	SL [80]	NeSyENT [3]	FIL	EBM [43]	MULTIPLEXNET [38]	CCN [33]	SPL ( <i>ours</i> )
(D1) Probabilistic	✗	✓	✓	✓	✗	✓	✗	✓
(D2) Expressive	✗	✗	✗	✗	✓	✗	✗	✓
(D3) Consistent	✗	✗	✗	✗	✗	✓	✓	✓
(D4) General	✓	✓	✓	✗	✓	✓	✗	✓
(D5) Modular	✓	✓	✓	✓	✓	✓	✓	✓
(D6) Efficient	✓	✓	✓	✓	✗	✗	✓	✓

# SPL



Ground Truth



ResNet-18



Semantic Loss



SPL (ours)

***you can predict valid paths 100% of the time!***

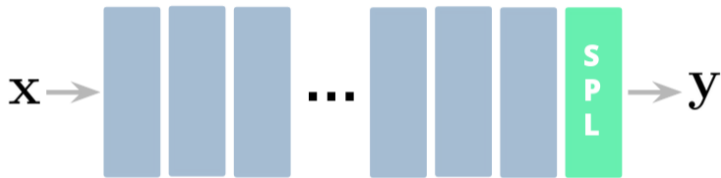


**How?**



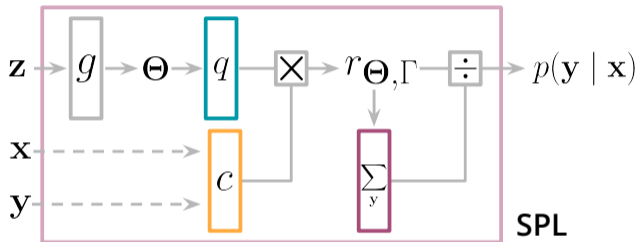
***take an unreliable neural network architecture...***

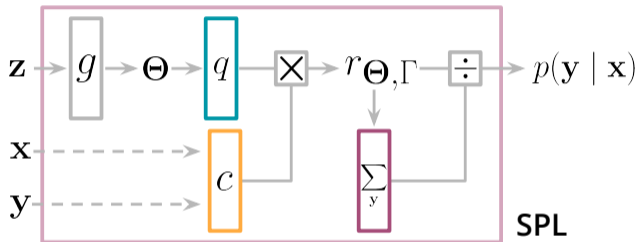
## How?



*.....and replace the last layer with  
a semantic probabilistic layer*

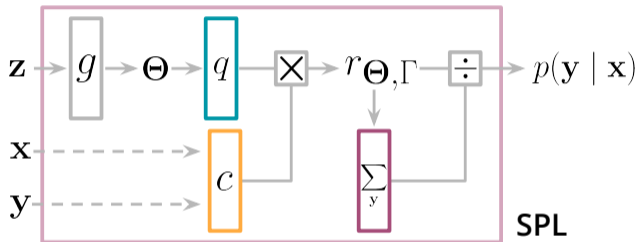
# SPL





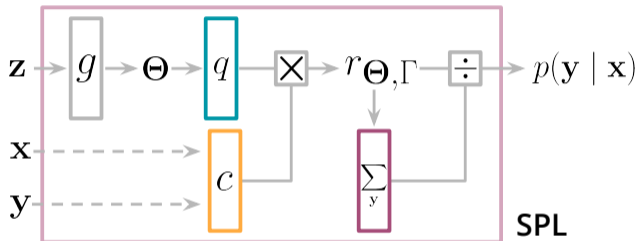
$$p(\mathbf{y} | \mathbf{x}) = \mathbf{q}_{\Theta}(\mathbf{y} | g(\mathbf{z}))$$

$\mathbf{q}_{\Theta}(\mathbf{y} | g(\mathbf{z}))$  is an expressive distribution over labels



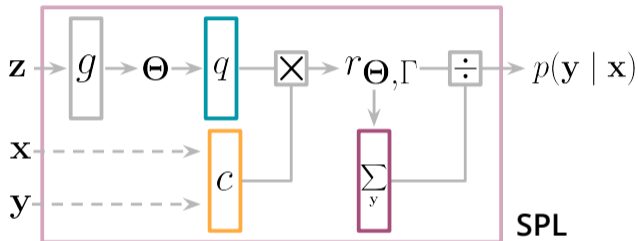
$$p(\mathbf{y} | \mathbf{x}) = q_{\Theta}(\mathbf{y} | g(\mathbf{z})) \cdot c_{\mathbf{K}}(\mathbf{x}, \mathbf{y})$$

$c_{\mathbf{K}}(\mathbf{x}, \mathbf{y})$  encodes the constraint  $\mathbb{1}\{\mathbf{x}, \mathbf{y} \models \mathbf{K}\}$



$$p(\mathbf{y} | \mathbf{x}) = q_{\Theta}(\mathbf{y} | g(\mathbf{z})) \cdot c_{\mathbf{K}}(\mathbf{x}, \mathbf{y})$$

***a product of experts : (***



$$p(\mathbf{y} | \mathbf{x}) = q_{\Theta}(\mathbf{y} | g(\mathbf{z})) \cdot c_{\mathbf{K}}(\mathbf{x}, \mathbf{y}) / \mathcal{Z}(\mathbf{x})$$

$$\mathcal{Z}(\mathbf{x}) = \sum_{\mathbf{y}} q_{\Theta}(\mathbf{y} | \mathbf{x}) \cdot c_{\mathbf{K}}(\mathbf{x}, \mathbf{y})$$

## Goal

*Can we design  $q$  and  $c$   
to be **expressive models**  
yet yielding a tractable product?*



# Probabilistic Circuits (PCs)

*A grammar for tractable computational graphs*

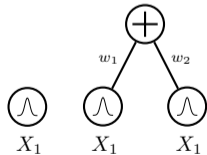
1. *A simple tractable function is a circuit*

  
 $X_1$

# Probabilistic Circuits (PCs)

*A grammar for tractable computational graphs*

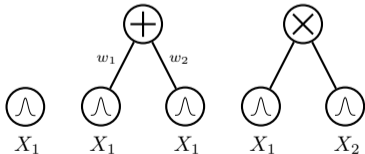
- I. *A simple tractable function is a circuit*
- II. *A weighted combination of circuits is a circuit*



# Probabilistic Circuits (PCs)

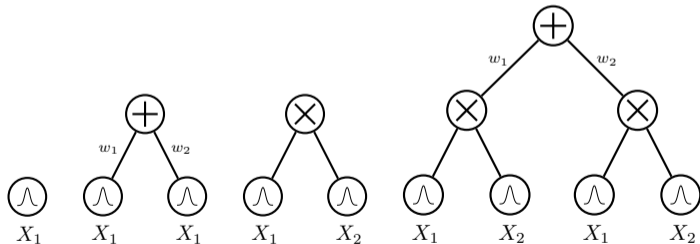
*A grammar for tractable computational graphs*

- I. *A simple tractable function is a circuit*
- II. *A weighted combination of circuits is a circuit*
- III. *A product of circuits is a circuit*



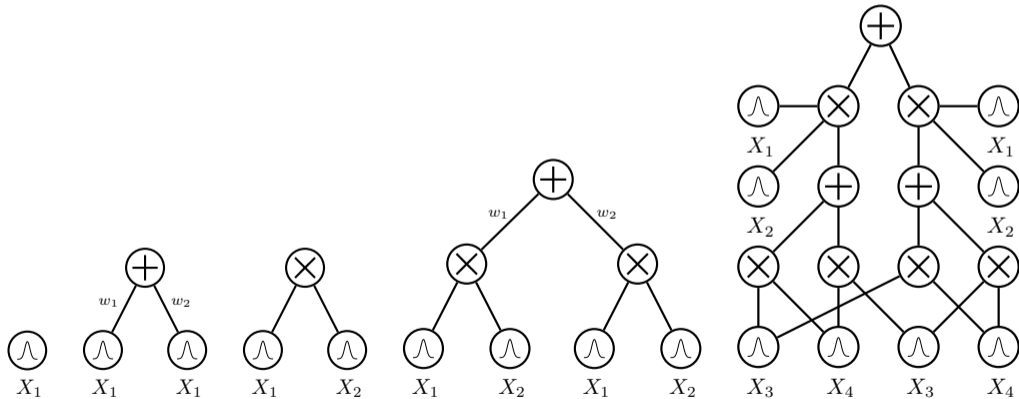
# Probabilistic Circuits (PCs)

*A grammar for tractable computational graphs*



# Probabilistic Circuits (PCs)

A grammar for tractable computational graphs



## **...why PCs?**

### **1. A grammar for tractable models**

One formalism to represent many models. *#HMMs #Trees #XGBoost, ...*

### **2. Expressiveness**

Competitive with intractable models, VAEs, Flow...*#hierachical #mixtures #polynomials*

## ...why PCs?

### 1. *A grammar for tractable models*

One formalism to represent many models. *#HMMs #Trees #XGBoost, ...*

### 2. *Expressiveness*

Competitive with intractable models, VAEs, Flow...*#hierachical #mixtures #polynomials*

### 3. *Tractability == Structural Properties!!!*

Exact computations of reasoning tasks are certified by guaranteeing certain structural properties. *#marginals #expectations #MAP, #product ...*

# ***Structural properties***

***smoothness***

***decomposability***

***determinism***

***compatibility***

---

*Vergari et al., "A Compositional Atlas of Tractable Circuit Operations: From Simple Transformations to Complex Information-Theoretic Queries", [NeurIPS](#), 2021*



# ***Structural properties***

***property A***

***property B***

***property C***

***property D***

---

*Vergari et al., "A Compositional Atlas of Tractable Circuit Operations: From Simple Transformations to Complex Information-Theoretic Queries", [NeurIPS](#), 2021*

# Structural properties

**property A**

**property B**

**property C**

**property D**

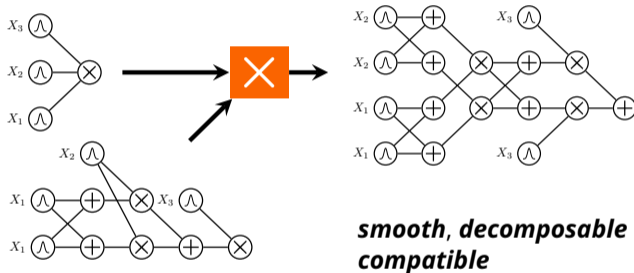
**tractable** computation of **arbitrary integrals**

$$p(\mathbf{y}) = \int_{\text{val}(\mathbf{Z})} p(\mathbf{z}, \mathbf{y}) d\mathbf{Z}, \quad \forall \mathbf{Y} \subseteq \mathbf{X}, \quad \mathbf{Z} = \mathbf{X} \setminus \mathbf{Y}$$

$\Rightarrow$  **sufficient** and **necessary** conditions  
for a single feedforward evaluation

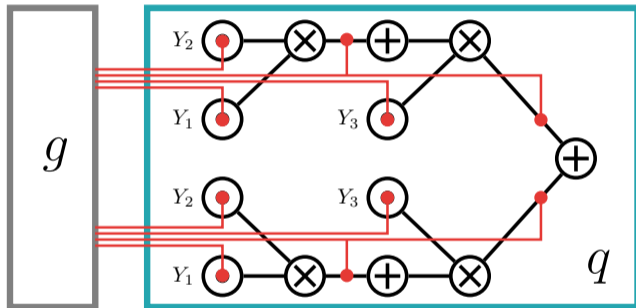
$\Rightarrow$  tractable partition function

# Tractable products

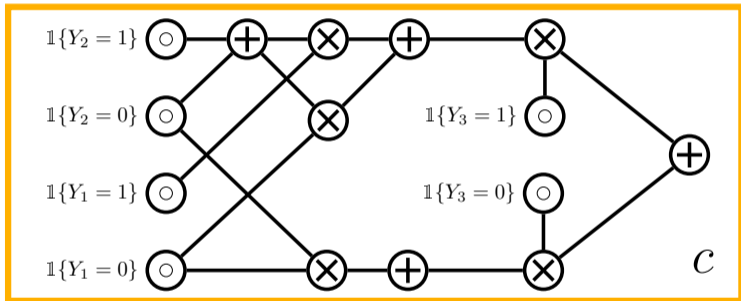


exactly compute  $\mathcal{Z}$  in time  $O(|q||c|)$

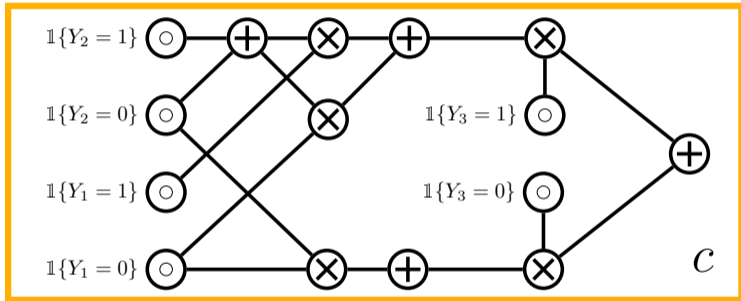
Vergari et al., "A Compositional Atlas of Tractable Circuit Operations: From Simple Transformations to Complex Information-Theoretic Queries", NeurIPS, 2021



**a conditional circuit  $q(y; \Theta = g(z))$**



*and a logical circuit  $c(y, x)$  encoding  $K$*



***compiling logical formulas into circuits***

# Knowledge compilation

$$K : (Y_1 = 1 \implies Y_3 = 1)$$

$$\wedge (Y_2 = 1 \implies Y_3 = 1)$$

$$\mathbb{1}\{Y_1 = 0\} \odot$$

$$\mathbb{1}\{Y_1 = 1\} \odot$$

$$\mathbb{1}\{Y_2 = 0\} \odot$$

$$\mathbb{1}\{Y_2 = 1\} \odot$$

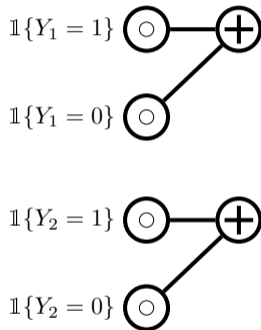
$$\mathbb{1}\{Y_3 = 0\} \odot$$

$$\mathbb{1}\{Y_3 = 1\} \odot$$

# Knowledge compilation

$$K : (Y_1 = 1 \implies Y_3 = 1)$$

$$\wedge (Y_2 = 1 \implies Y_3 = 1)$$

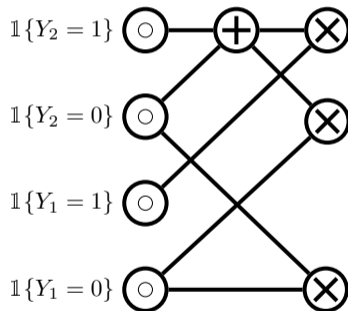




# Knowledge compilation

$K : (Y_1 = 1 \implies Y_3 = 1)$

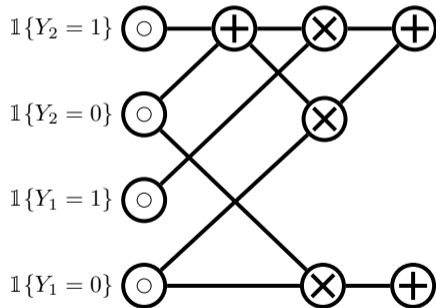
$\wedge (Y_2 = 1 \implies Y_3 = 1)$



# Knowledge compilation

$K : (Y_1 = 1 \implies Y_3 = 1)$

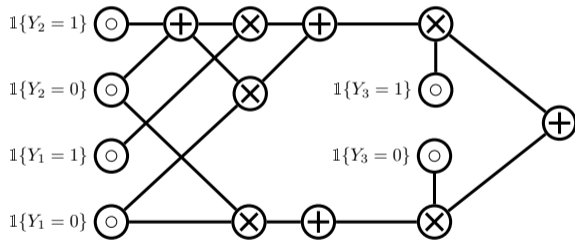
$\wedge (Y_2 = 1 \implies Y_3 = 1)$



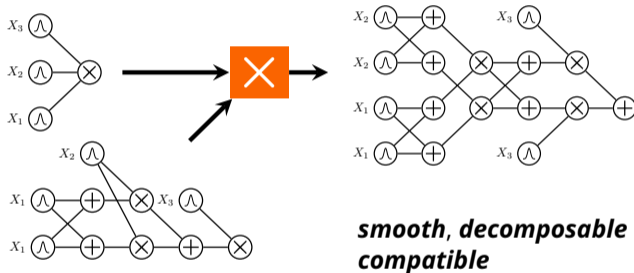
# Knowledge compilation

$$K : (Y_1 = 1 \implies Y_3 = 1)$$

$$\wedge (Y_2 = 1 \implies Y_3 = 1)$$



# Tractable products



exactly compute  $\mathcal{Z}$  in time  $O(|q||c|)$

Vergari et al., "A Compositional Atlas of Tractable Circuit Operations: From Simple Transformations to Complex Information-Theoretic Queries", NeurIPS, 2021

## ***SPL recipe***

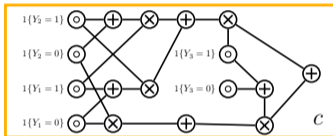
$$K : (Y_1 = 1 \implies Y_3 = 1)$$

$$\wedge (Y_2 = 1 \implies Y_3 = 1)$$

**1)** Take a  
logical constraint

## SPL recipe

$$K : (Y_1 = 1 \implies Y_3 = 1) \\ \wedge (Y_2 = 1 \implies Y_3 = 1)$$

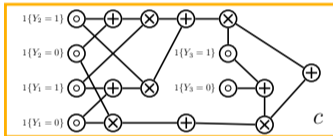


**1)** Take a  
logical constraint

**2)** Compile it into  
a constraint circuit

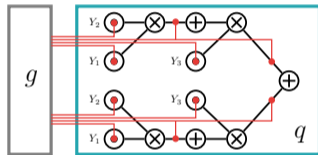
# SPL recipe

$$K : (Y_1 = 1 \implies Y_3 = 1) \\ \wedge (Y_2 = 1 \implies Y_3 = 1)$$



**1)** Take a logical constraint

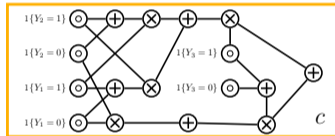
**2)** Compile it into a constraint circuit



**3)** Multiply it by a circuit distribution

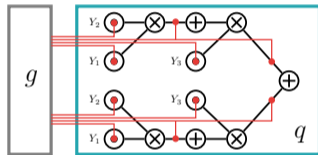
# SPL recipe

$$K : (Y_1 = 1 \implies Y_3 = 1) \\ \wedge (Y_2 = 1 \implies Y_3 = 1)$$



1) Take a logical constraint

2) Compile it into a constraint circuit

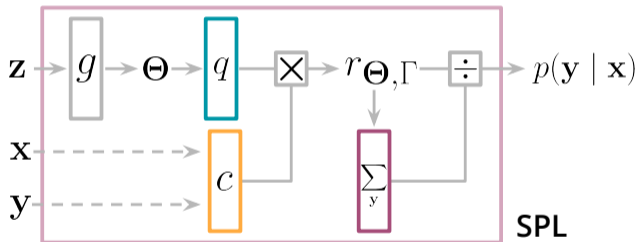


3) Multiply it by a circuit distribution

4) *train end-to-end by sgd!*

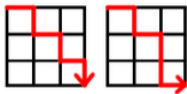


# Experiments



*how good are SPLs?*

# Experiments



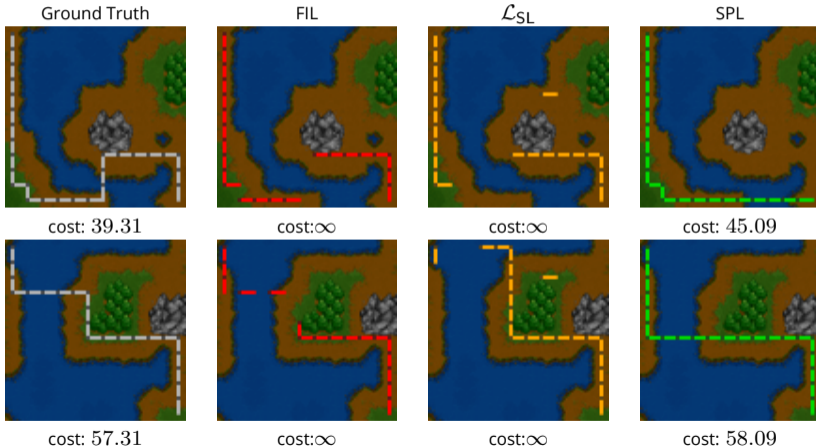
Architecture	Simple Path			Preference Learning		
	Exact	Hamming	Consistent	Exact	Hamming	Consistent
MLP+FIL	5.6	85.9	7.0	1.0	<b>75.8</b>	2.7
MLP+ $\mathcal{L}_{SL}$	28.5	83.1	75.2	15.0	72.4	69.8
MLP+NeSyEnt	30.1	83.0	91.6	18.2	71.5	96.0
MLP+SPL	<b>37.6</b>	<b>88.5</b>	<b>100.0</b>	<b>20.8</b>	72.4	<b>100.0</b>

# Experiments



Architecture	Exact	Hamming	Consistent
ResNet-18+FIL	55.0	<b>97.7</b>	56.9
ResNet-18+ $\mathcal{L}_{SL}$	59.4	<b>97.7</b>	61.2
ResNet-18+SPL	<b>78.2</b>	96.3	<b>100.0</b>

# Experiments



# Experiments

DATASET	EXACT MATCH	
	HMCNN	MLP+SPL
CELLCYCLE	3.05 ± 0.11	<b>3.79 ± 0.18</b>
DERISI	1.39 ± 0.47	<b>2.28 ± 0.23</b>
EISEN	5.40 ± 0.15	<b>6.18 ± 0.33</b>
EXPR	4.20 ± 0.21	<b>5.54 ± 0.36</b>
GASCH1	3.48 ± 0.96	<b>4.65 ± 0.30</b>
GASCH2	3.11 ± 0.08	<b>3.95 ± 0.28</b>
SEQ	5.24 ± 0.27	<b>7.98 ± 0.28</b>
SPO	<b>1.97 ± 0.06</b>	<b>1.92 ± 0.11</b>
DIATOMS	48.21 ± 0.57	<b>58.71 ± 0.68</b>
ENRON	5.97 ± 0.56	<b>8.18 ± 0.68</b>
IMCLEF07A	79.75 ± 0.38	<b>86.08 ± 0.45</b>
IMCLEF07D	76.47 ± 0.35	<b>81.06 ± 0.68</b>

---

## How to Turn Your Knowledge Graph Embeddings into Generative Models via Probabilistic Circuits

---

**Lorenzo Loconte**

University of Edinburgh, UK  
l.loconte@sms.ed.ac.uk

**Nicola Di Mauro**

University of Bari, Italy  
nicola.dimauro@uniba.it

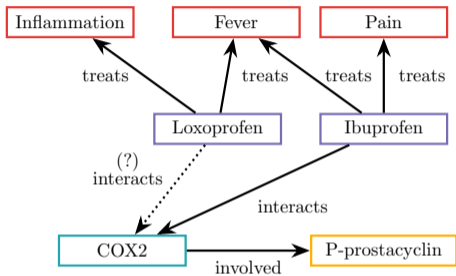
**Robert Peharz**

TU Graz, Austria  
robert.pehartz@tugraz.at

**Antonio Vergari**

University of Edinburgh, UK  
avergari@ed.ac.uk

***SPL meets knowledge graph embedding models***



- Drugs
- Symptoms
- Proteins
- Functions

⟨Loxoprofen, treats, Inflammation⟩

⟨Ibuprofen, interacts, COX2⟩

⟨COX2, involved, P-prostacyclin⟩

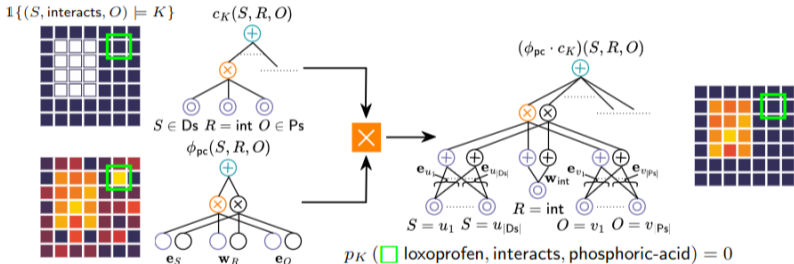
...

Q: ⟨Loxoprofen, interacts, ?⟩

A: ⟨Loxoprofen, interacts, **phosphoric-acid**⟩ !!!

***neural link predictors can violate domain constraints***

# SPL+KGE







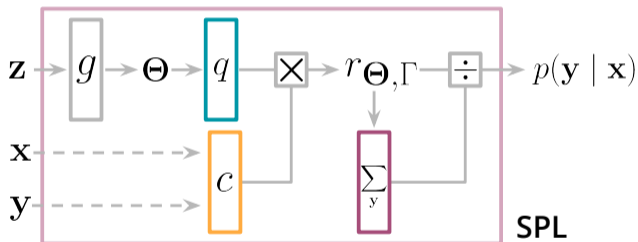
**MACHINE  
LEARNERS**

**COMPLEX  
REASONING**

A close-up photograph of a person's hand holding a black eraser against a chalkboard. The word "APPROXIMATIONS" is written in large, bold, white capital letters across the center of the image. The chalkboard has some faint, illegible markings and a yellowish smudge. In the background, a sign with the letters "s la" is partially visible.

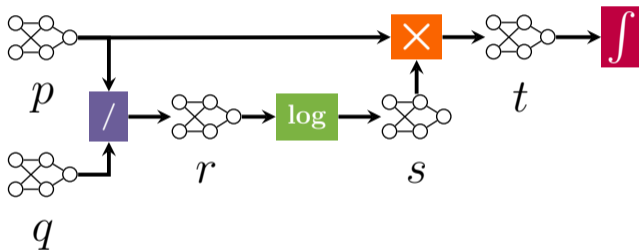
**APPROXIMATIONS**

$$p(\mathbf{y} \mid \mathbf{x}) = \mathbf{q}_{\Theta}(\mathbf{y} \mid g(\mathbf{z})) \cdot \mathbf{c}_{\mathbf{K}}(\mathbf{x}, \mathbf{y}) / \mathbf{Z}(\mathbf{x})$$



***efficient and reliable reasoning over constraints***

$$\int p(\mathbf{x}) \times \log \left( p(\mathbf{x}) / q(\mathbf{x}) \right) d\mathbf{X}$$



**build a LEGO-like query calculus!**

*Vergari et al., "A Compositional Atlas of Tractable Circuit Operations: From Simple Transformations to Complex Information-Theoretic Queries", NeurIPS, 2021*

**part I** how to satisfy wanted constraints  
in NNs by design

**part II** how the design of NNs implicitly poses  
unwanted constraints

# **Taming the Sigmoid Bottleneck: Provably Argmaxable Sparse Multi-Label Classification**

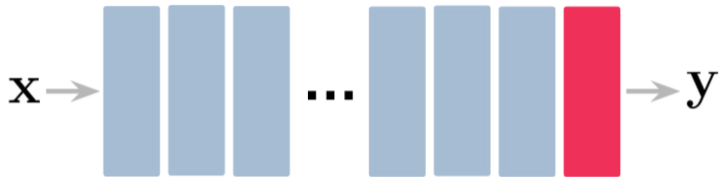
**Andreas Grivas<sup>1</sup>, Antonio Vergari<sup>†2</sup>, Adam Lopez<sup>†1</sup>**

<sup>1</sup> Institute for Language, Cognition, and Computation

<sup>2</sup> Institute for Adaptive and Neural Computation

School of Informatics, University of Edinburgh

{agrivas, avergari, alopez}@ed.ac.uk



***sigmoid linear layers***

$$p(\mathbf{y} \mid \mathbf{x}) = \prod_{i=1}^n p(y_i \mid \mathbf{x})$$



***low-rank classifiers***

# labels ( $n$ )  $\gg$  embedding size ( $d$ )



*e.g.,*



*clinically annotated text*

$n \approx 9000$

$d \approx 200 - 500$



*large-scale biomedical  
question answering*

$n \approx 20000$

$d \approx 200 - 500$



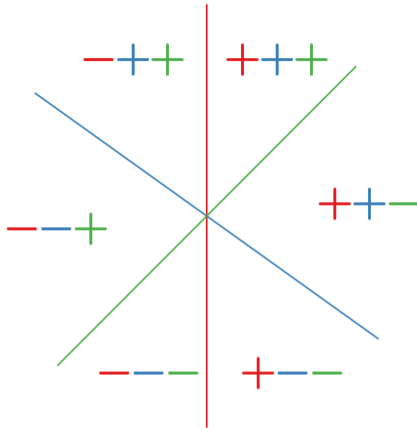
OpenImages Dataset  
*object recognition*

$n \approx 9000$

$d \approx 200 - 500$

$$n = 3$$

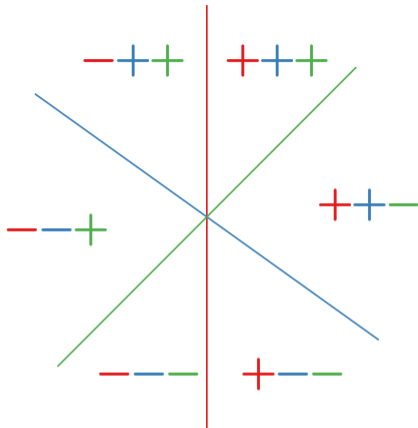
$$d = 2$$



some label configurations are *unargmaxable!*

$$n = 3$$

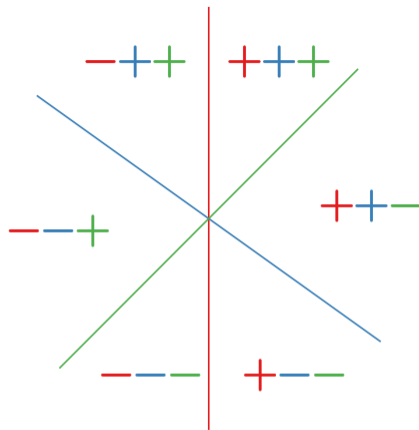
$$d = 2$$



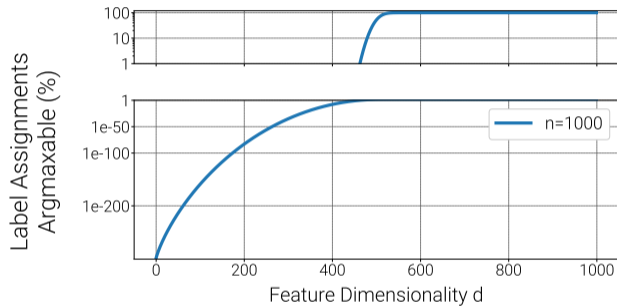
$$\nexists \mathbf{x} : \operatorname{argmax}_{\mathbf{y}'} p(\mathbf{y}' | \mathbf{x}; \mathbf{W}) = \mathbf{y}$$

$$n = 3$$

$$d = 2$$



$-+-$  and  $+-+$  are *unargmaxable!*



***exponentially*** many configurations are unargmaxable

***but real data is sparse...***

*K-active labels*



*clinically annotated text*

$n \approx 9000$

$K = 80$



*large-scale biomedical  
question answering*

$n \approx 20000$

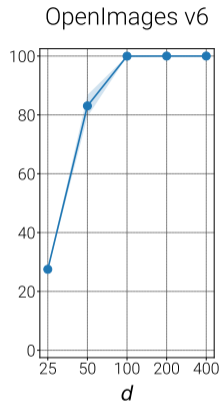
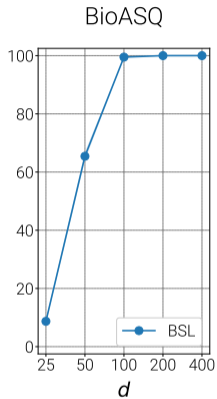
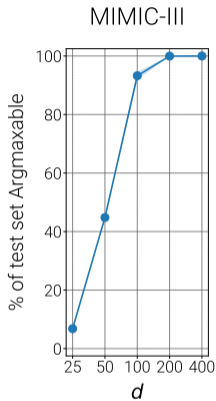
$K = 50$



OpenImages Dataset  
*object recognition*

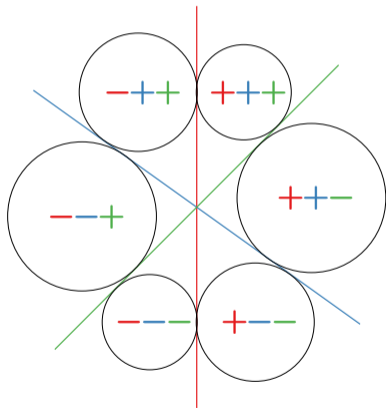
$n \approx 9000$

$K = 50$



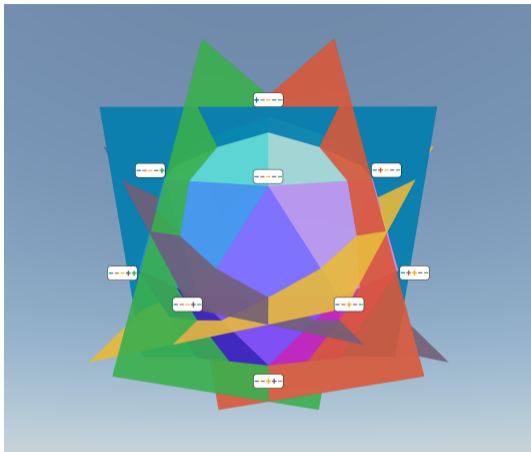
**even sparse label configurations are unargmaxable**

***how do we know?***

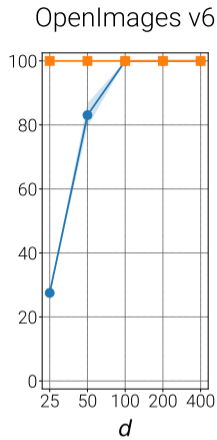
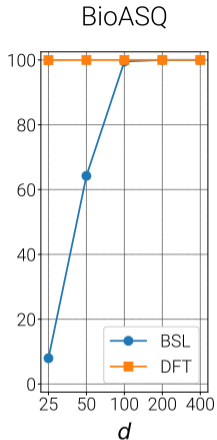
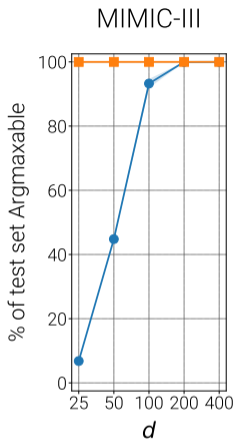


**we provide a Chebyshev LP to *verify argmaxability***



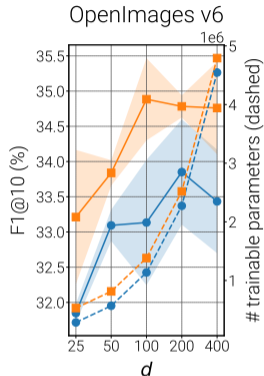
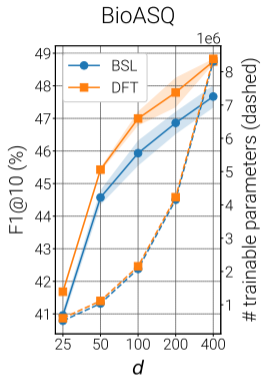
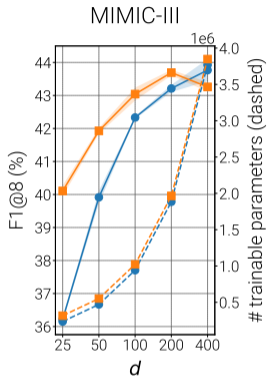


and a new differentiable layer to ***guarantee argmaxability***  
for  $K$ -active label configurations



**and a new differentiable layer to guarantee argmaxability**

***based on the DFT***



**with comparable or better performance**

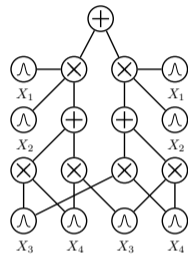
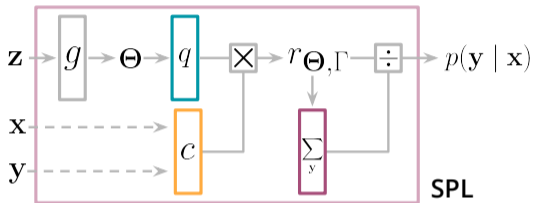


**part I** how to satisfy wanted constraints  
in NNs by design



**part II**

**how the design of NNs implicitly poses unwanted constraints**



**Ask me anything!**