

# Probabilistic and Logistic Circuits: A New Synthesis of Logic and Machine Learning

Guy Van den Broeck

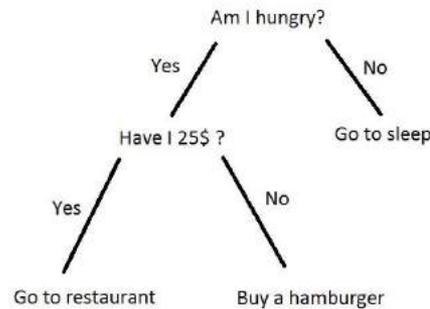
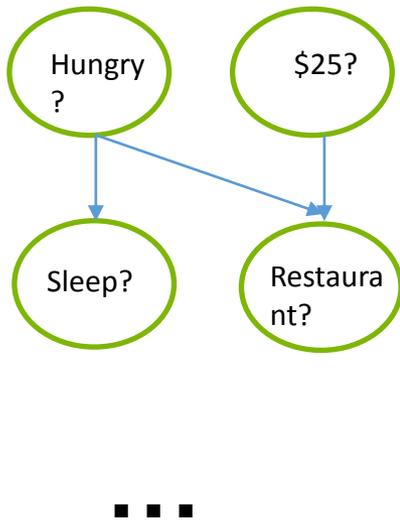
**UCLA**

RelationalAI ArrowCon  
Feb 5, 2019



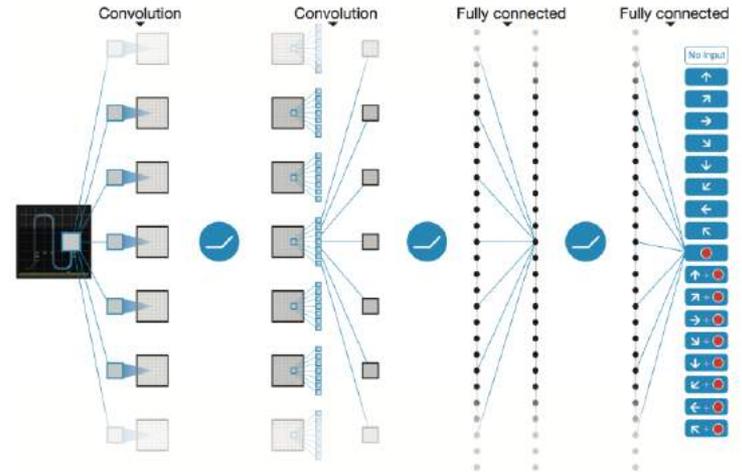
# Which method to choose?

## Classical AI Methods:



Clear Modeling Assumption  
Well-understood

## Neural Networks:



“Black Box”  
Good performance  
on Image Classification

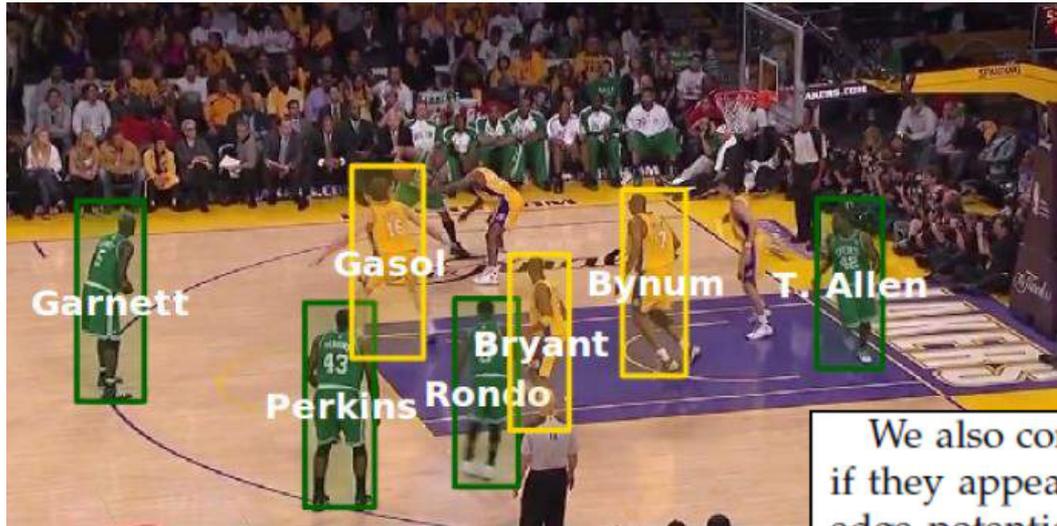
# Outline

- Adding knowledge to deep learning
- Probabilistic circuits
- Logistic circuits for image classification

# Outline

- **Adding knowledge to deep learning**
- Probabilistic circuits
- Logistic circuits for image classification

# Motivation: Video

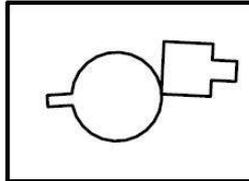
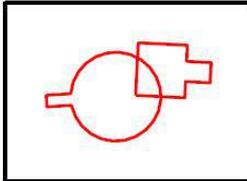
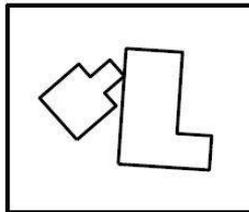
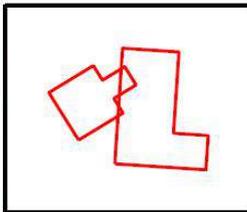
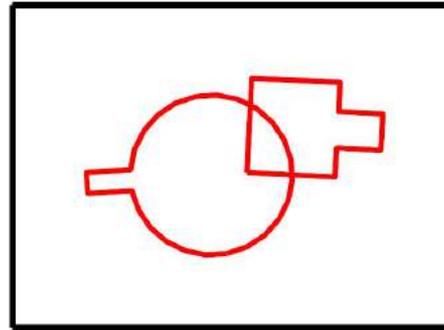


We also connect all pairs of identity nodes  $y_{t,i}$  and  $y_{t,j}$  if they appear in the same time  $t$ . We then introduce an edge potential that enforces mutual exclusion:

$$\psi_{\text{mutex}}(y_{t,i}, y_{t,j}) = \begin{cases} 1 & \text{if } y_{t,i} \neq y_{t,j} \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

This potential specifies the constraint that a player can be **appear only once in a frame**. For example, if the  $i$ -th detection  $y_{t,i}$  has been assign to Bryant,  $y_{t,j}$  cannot have the same identity because Bryant is impossible to appear twice in a frame.

# Motivation: Robotics



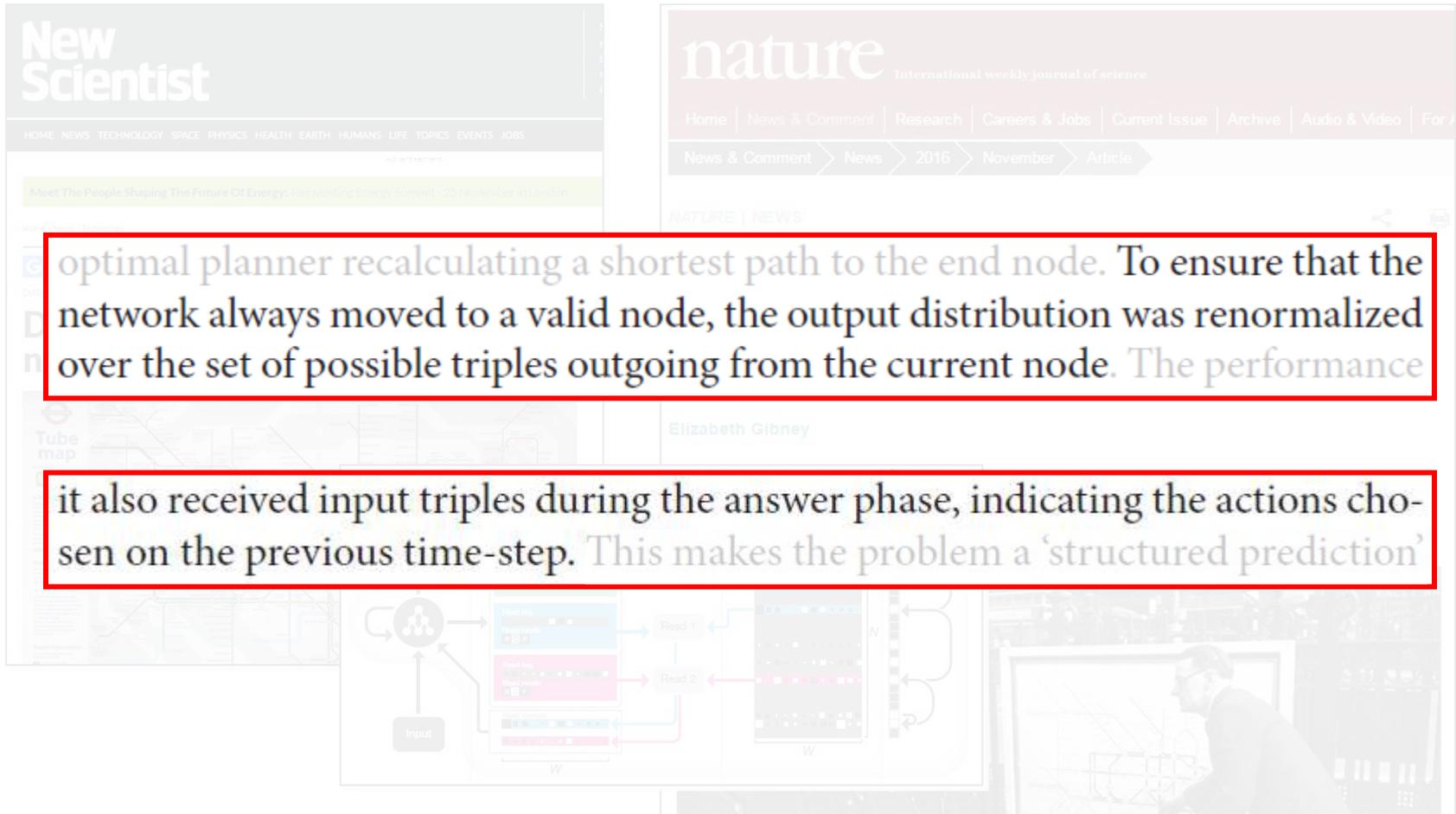
The method developed in this paper can be used in a broad variety of semantic mapping and object manipulation tasks, providing an efficient and effective way to incorporate collision constraints into a recursive state estimator, obtaining optimal or near-optimal solutions.

# Motivation: Language

- Non-local dependencies:  
*At least one verb in each sentence*
  - Sentence compression  
*If a modifier is kept, its subject is also kept*
  - Information extraction
  - Semantic role labeling
- ... and many more!

Citations	
Start	The citation must start with author or editor.
AppearsOnce	Each field must be a consecutive list of words, and can appear at most once in a citation.
Punctuation	State transitions must occur on punctuation marks.
BookJournal	The words <i>proc</i> , <i>journal</i> , <i>proceedings</i> , <i>ACM</i> are <i>JOURNAL</i> or <i>BOOKTITLE</i> .
...	...
TechReport	The words <i>tech</i> , <i>technical</i> are <i>TECH_REPORT</i> .
Title	Quotations can appear only in titles.
Location	The words <i>CA</i> , <i>Australia</i> , <i>NY</i> are <i>LOCATION</i> .

# Motivation: Deep Learning



optimal planner recalculating a shortest path to the end node. To ensure that the network always moved to a valid node, the output distribution was renormalized over the set of possible triples outgoing from the current node. The performance

it also received input triples during the answer phase, indicating the actions chosen on the previous time-step. This makes the problem a 'structured prediction'

The background features a collage of images: the New Scientist website on the left, the Nature website on the right, a Tube map in the bottom left, and a neural network diagram with an 'Input' node and 'Read 1', 'Read 2' nodes in the bottom center. A person is visible in the bottom right corner, looking at a screen.

[Graves, A., Wayne, G., Reynolds, M., Harley, T., Danihelka, I., Grabska-Barwińska, A., et al.. (2016). Hybrid computing using a neural network with dynamic external memory. *Nature*, 538(7626), 471-476.]

# Running Example

## Courses:

- Logic (L)
- Knowledge Representation (K)
- Probability (P)
- Artificial Intelligence (A)

## Constraints

- Must take at least one of Probability or Logic.
- Probability is a prerequisite for AI.
- The prerequisites for KR is either AI or Logic.

## Data

L	K	P	A	Students
0	0	1	0	6
0	0	1	1	54
0	1	1	1	10
1	0	0	0	5
1	0	1	0	1
1	0	1	1	0
1	1	0	0	17
1	1	1	0	4
1	1	1	1	3

# Structured Space

unstructured

L	K	P	A
0	0	0	0
0	0	0	1
0	0	1	0
0	0	1	1
0	1	0	0
0	1	0	1
0	1	1	0
0	1	1	1
1	0	0	0
1	0	0	1
1	0	1	0
1	0	1	1
1	1	0	0
1	1	0	1
1	1	1	0
1	1	1	1



structured

L	K	P	A
0	0	0	0
0	0	0	1
0	0	1	0
0	0	1	1
0	1	0	0
0	1	0	1
0	1	1	0
0	1	1	1
1	0	0	0
1	0	0	1
1	0	1	0
1	0	1	1
1	1	0	0
1	1	0	1
1	1	1	0
1	1	1	1

- Must take at least one of Probability (**P**) or Logic (**L**).
- Probability is a prerequisite for AI (**A**).
- The prerequisites for KR (**K**) is either AI or Logic.

**7 out of 16 instantiations  
are impossible**

# Boolean Constraints

unstructured

L	K	P	A
0	0	0	0
0	0	0	1
0	0	1	0
0	0	1	1
0	1	0	0
0	1	0	1
0	1	1	0
0	1	1	1
1	0	0	0
1	0	0	1
1	0	1	0
1	0	1	1
1	1	0	0
1	1	0	1
1	1	1	0
1	1	1	1



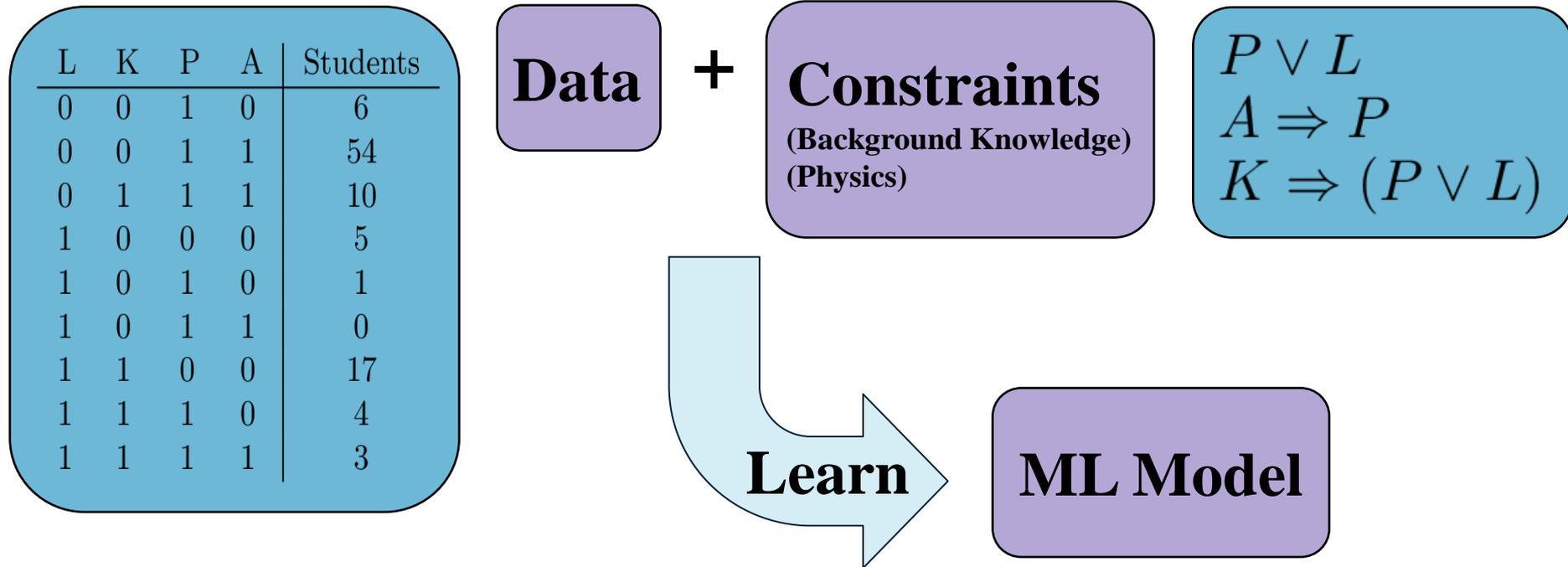
structured

L	K	P	A
0	0	0	0
0	0	0	1
0	0	1	0
0	0	1	1
0	1	0	0
0	1	0	1
0	1	1	0
0	1	1	1
1	0	0	0
1	0	0	1
1	0	1	0
1	0	1	1
1	1	0	0
1	1	0	1
1	1	1	0
1	1	1	1

$$\begin{aligned} P \vee L \\ A \Rightarrow P \\ K \Rightarrow (P \vee L) \end{aligned}$$

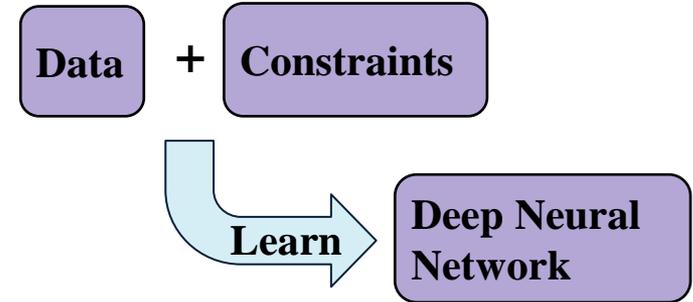
**7 out of 16 instantiations  
are impossible**

# Learning in Structured Spaces



Today's machine learning tools  
don't take knowledge as input! ☹️

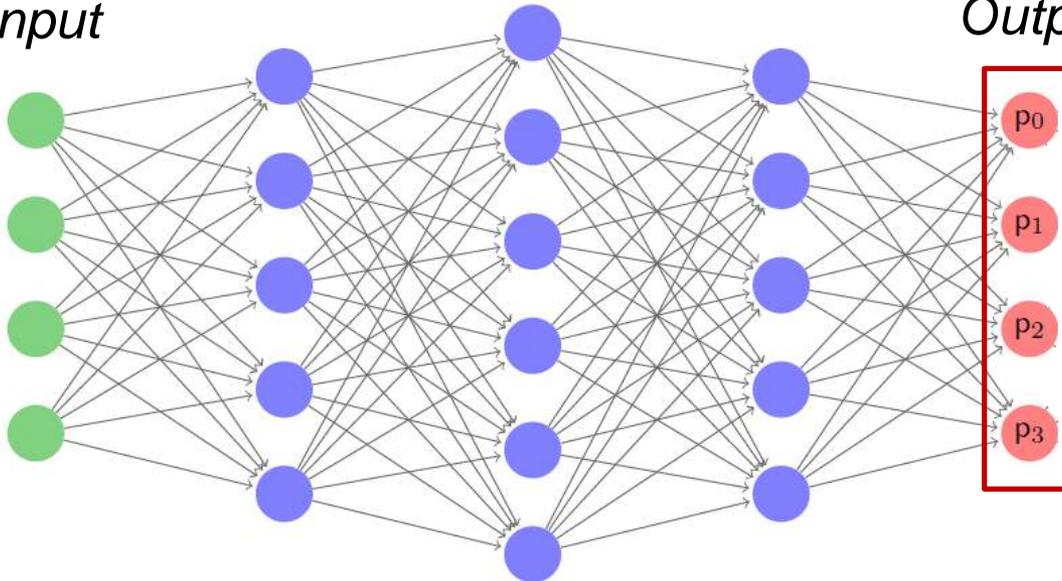
# Deep Learning with Logical Knowledge



*Neural Network*

*Input*

*Output*



Output is  
probability vector  $\mathbf{p}$ ,  
not Boolean logic!

# Semantic Loss

*Q: How close is output  $\mathbf{p}$  to satisfying constraint?*

Answer: Semantic loss function  $L(\alpha, \mathbf{p})$

- Axioms, for example:
  - If  $\mathbf{p}$  is Boolean then  $L(\mathbf{p}, \mathbf{p}) = 0$
  - If  $\alpha$  implies  $\beta$  then  $L(\alpha, \mathbf{p}) \geq L(\beta, \mathbf{p})$  ( *$\alpha$  more strict*)
- Properties:
  - If  $\alpha$  is equivalent to  $\beta$  then  $L(\alpha, \mathbf{p}) = L(\beta, \mathbf{p})$  SEMANTIC
  - If  $\mathbf{p}$  is Boolean and satisfies  $\alpha$  then  $L(\alpha, \mathbf{p}) = 0$  Loss!

# Semantic Loss: Definition

Theorem: Axioms imply unique semantic loss:

$$L^S(\alpha, \mathbf{p}) \propto -\log \sum_{\mathbf{x} \models \alpha} \prod_{i: \mathbf{x} \models X_i} p_i \prod_{i: \mathbf{x} \models \neg X_i} (1 - p_i)$$

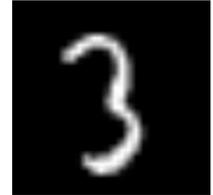
Probability of getting  $\mathbf{x}$  after  
flipping coins with prob.  $\mathbf{p}$

Probability of satisfying  $\alpha$  after  
flipping coins with prob.  $\mathbf{p}$

# Example: Exactly-One

- Data must have some label

*We agree this must be one of the 10 digits:*



- Exactly-one constraint  
→ For 3 classes: 
$$\begin{cases} x_1 \vee x_2 \vee x_3 \\ \neg x_1 \vee \neg x_2 \\ \neg x_2 \vee \neg x_3 \\ \neg x_1 \vee \neg x_3 \end{cases}$$

- Semantic loss:

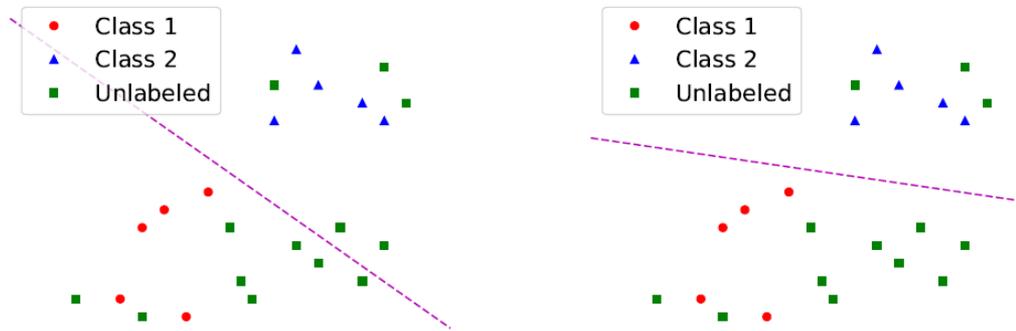
$$L^s(\text{exactly-one}, p) \propto -\log \sum_{i=1}^n p_i \prod_{j=1, j \neq i}^n (1 - p_j)$$

Only  $x_i = 1$  after flipping coins

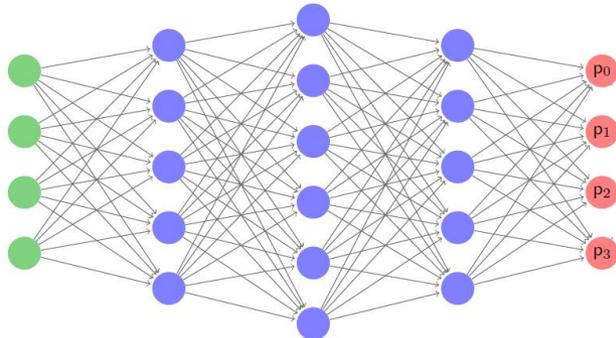
Exactly one true  $x$  after flipping coins

# Semi-Supervised Learning

- Intuition: Unlabeled data must have some label  
Cf. entropy constraints, manifold learning



- Minimize exactly-one semantic loss on unlabeled data



Train with  
*existing loss* +  $w \cdot$  *semantic loss*

# MNIST Experiment



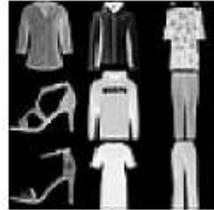
Accuracy % with # of used labels	100	1000	ALL
AtlasRBF (Pitelis et al., 2014)	91.9 ( $\pm 0.95$ )	96.32 ( $\pm 0.12$ )	98.69
Deep Generative (Kingma et al., 2014)	96.67( $\pm 0.14$ )	97.60( $\pm 0.02$ )	99.04
Virtual Adversarial (Miyato et al., 2016)	97.67	98.64	99.36
Ladder Net (Rasmus et al., 2015)	<b>98.94</b> ( $\pm 0.37$ )	<b>99.16</b> ( $\pm 0.08$ )	99.43 ( $\pm 0.02$ )
Baseline: MLP, Gaussian Noise	78.46 ( $\pm 1.94$ )	94.26 ( $\pm 0.31$ )	99.34 ( $\pm 0.08$ )
Baseline: Self-Training	72.55 ( $\pm 4.21$ )	87.43 ( $\pm 3.07$ )	
MLP with Semantic Loss	98.38 ( $\pm 0.51$ )	98.78 ( $\pm 0.17$ )	99.36 ( $\pm 0.02$ )

Competitive with state of the art  
in semi-supervised deep learning

# FASHION Experiment



(a) Confidently Correct



(b) Unconfidently Correct



(c) Unconfidently Incorrect



(d) Confidently Incorrect

Accuracy % with # of used labels	100	500	1000	ALL
Ladder Net (Rasmus et al., 2015)	81.46 ( $\pm 0.64$ )	85.18 ( $\pm 0.27$ )	86.48 ( $\pm 0.15$ )	90.46
Baseline: MLP, Gaussian Noise	69.45 ( $\pm 2.03$ )	78.12 ( $\pm 1.41$ )	80.94 ( $\pm 0.84$ )	89.87
MLP with Semantic Loss	<b>86.74</b> ( $\pm 0.71$ )	<b>89.49</b> ( $\pm 0.24$ )	89.67 ( $\pm 0.09$ )	89.81

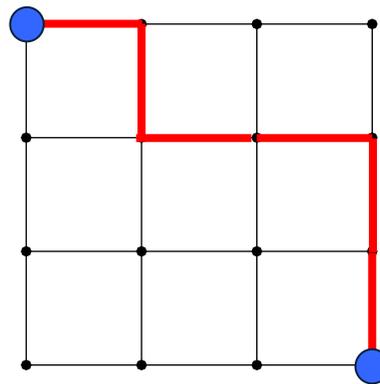
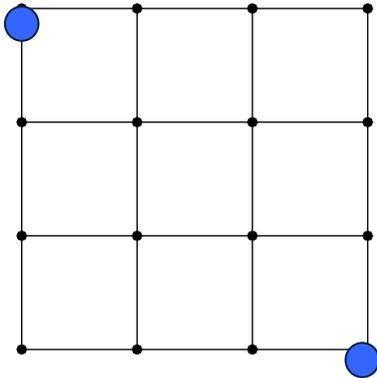
Outperforms Ladder Nets!

Same conclusion on CIFAR10

Accuracy % with # of used labels	4000	ALL
CNN Baseline in Ladder Net	76.67 ( $\pm 0.61$ )	90.73
Ladder Net (Rasmus et al., 2015)	79.60 ( $\pm 0.47$ )	
Baseline: CNN, Whitening, Cropping	77.13	90.96
CNN with Semantic Loss	<b>81.79</b>	90.92

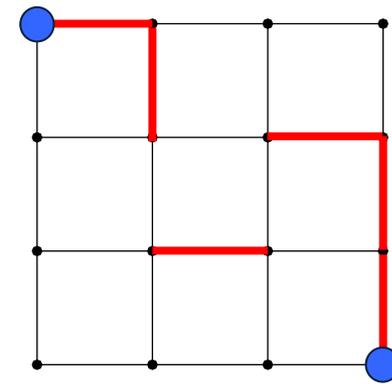
# What about real constraints? Paths

cf. Nature paper



Good variable assignment  
(represents route)

184



Bad variable assignment  
(does not represent route)

16,777,032

Unstructured probability space:  $184 + 16,777,032 = 2^{24}$

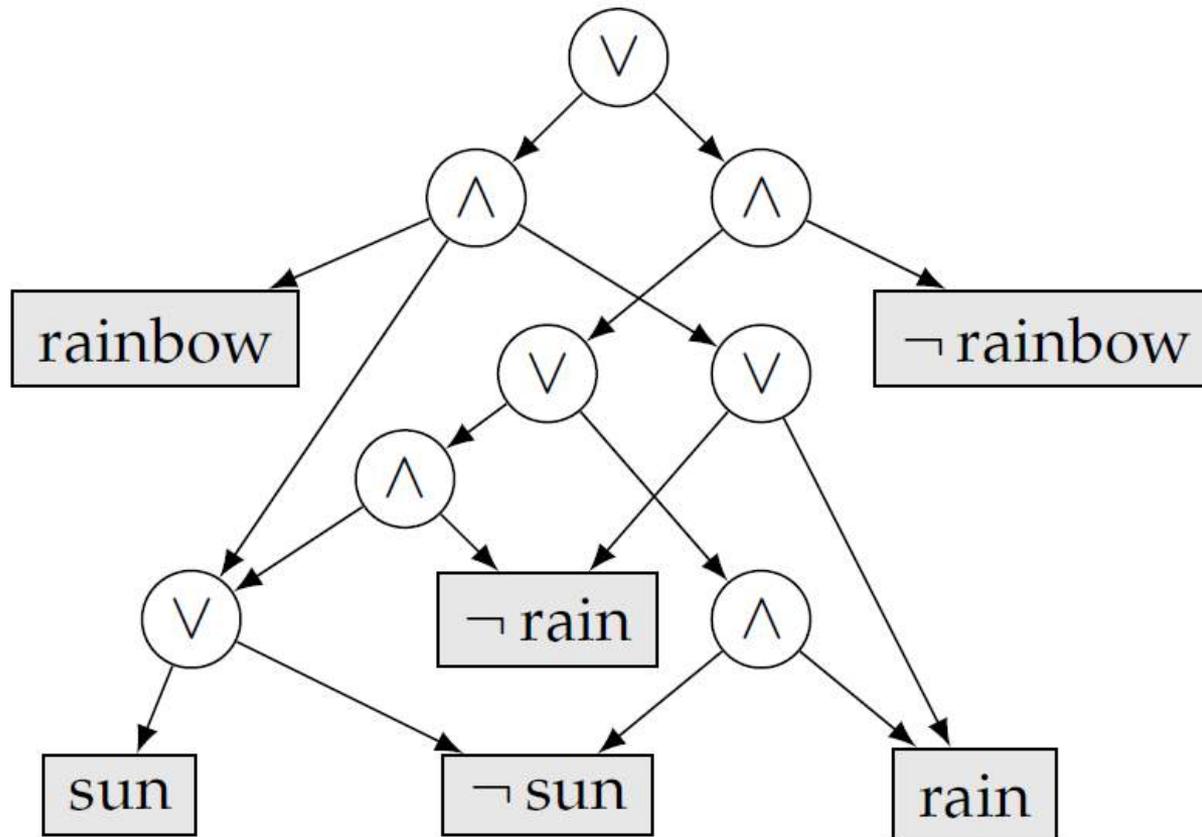
Space easily encoded in logical constraints 😊 [Nishino et al.]

# How to Compute Semantic Loss?

- In general: #P-hard ☹️

# Negation Normal Form Circuits

$$\Delta = (\text{sun} \wedge \text{rain} \Rightarrow \text{rainbow})$$



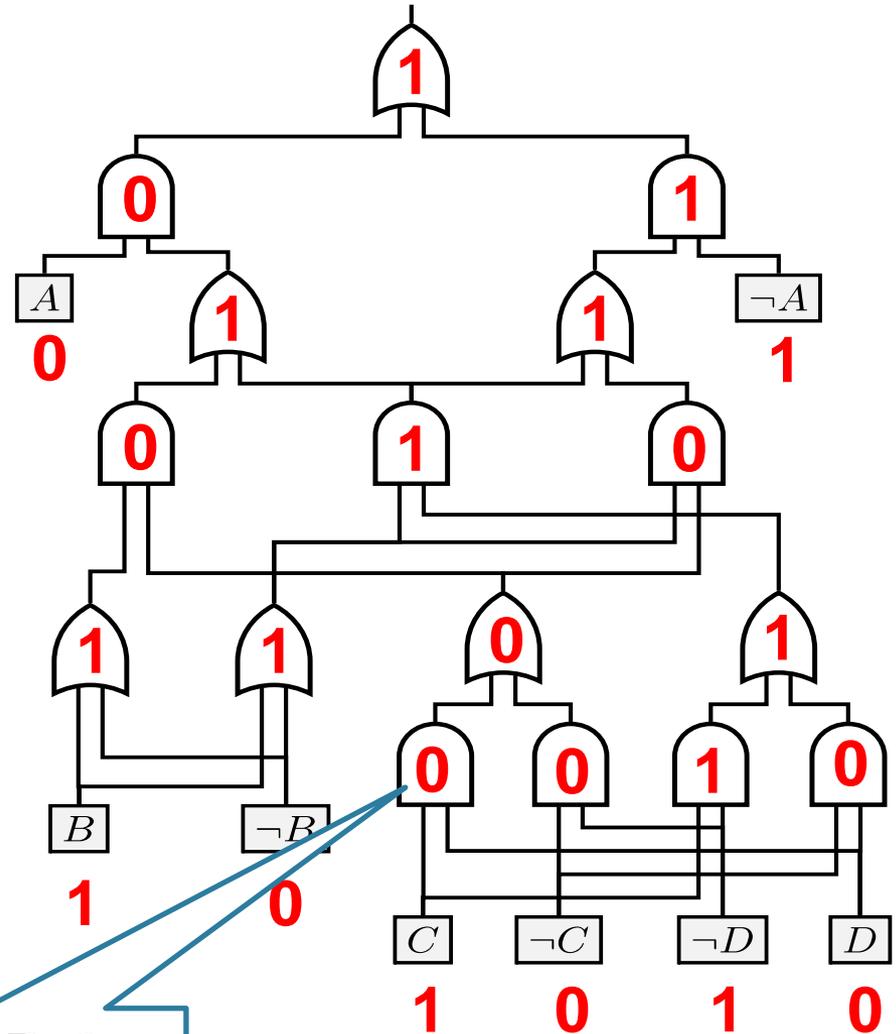
# Logical Circuits

Input:

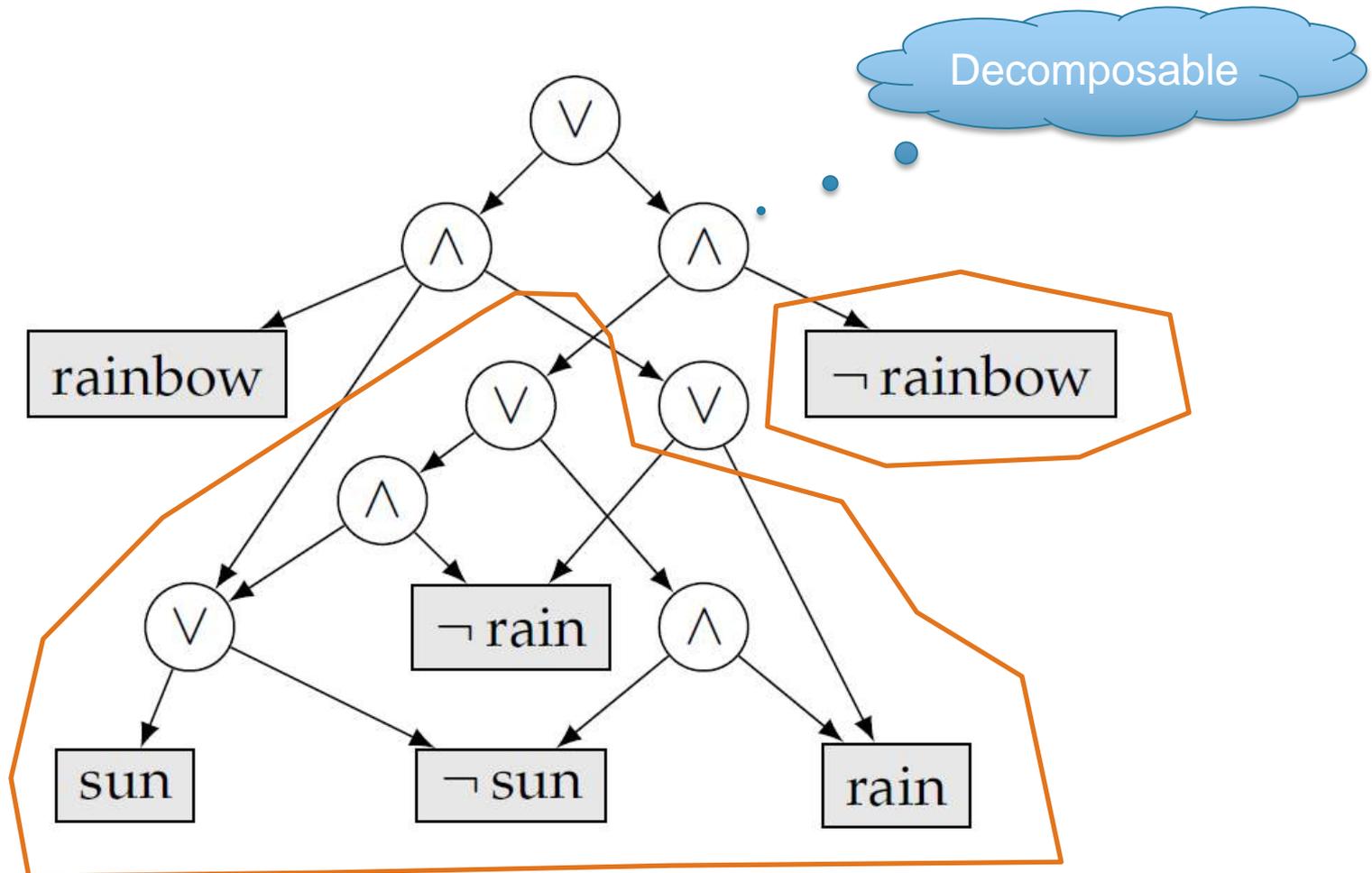
<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>
0	1	1	0

Bottom-up Evaluation

**0 = 1 AND 0**



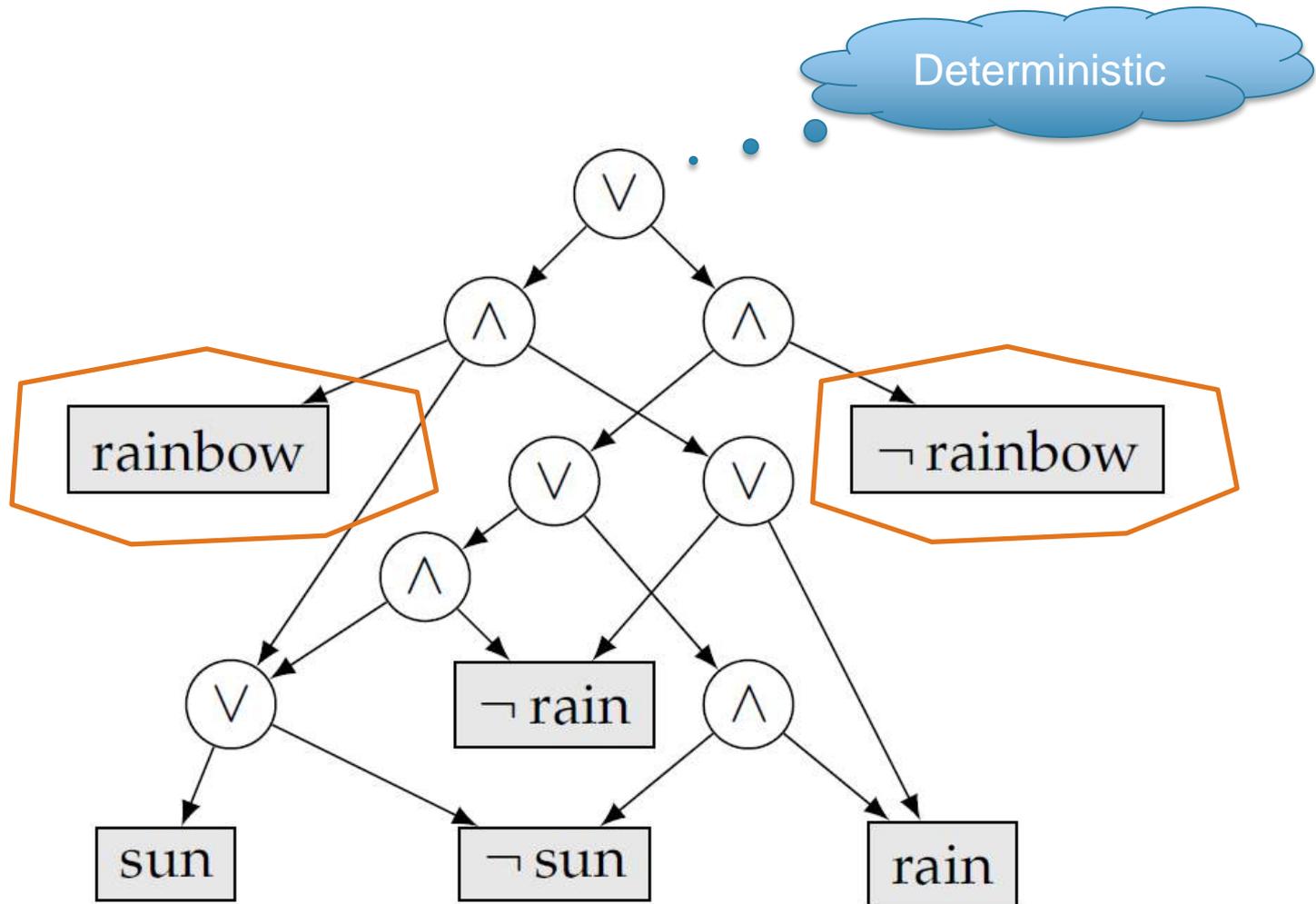
# Decomposable Circuits



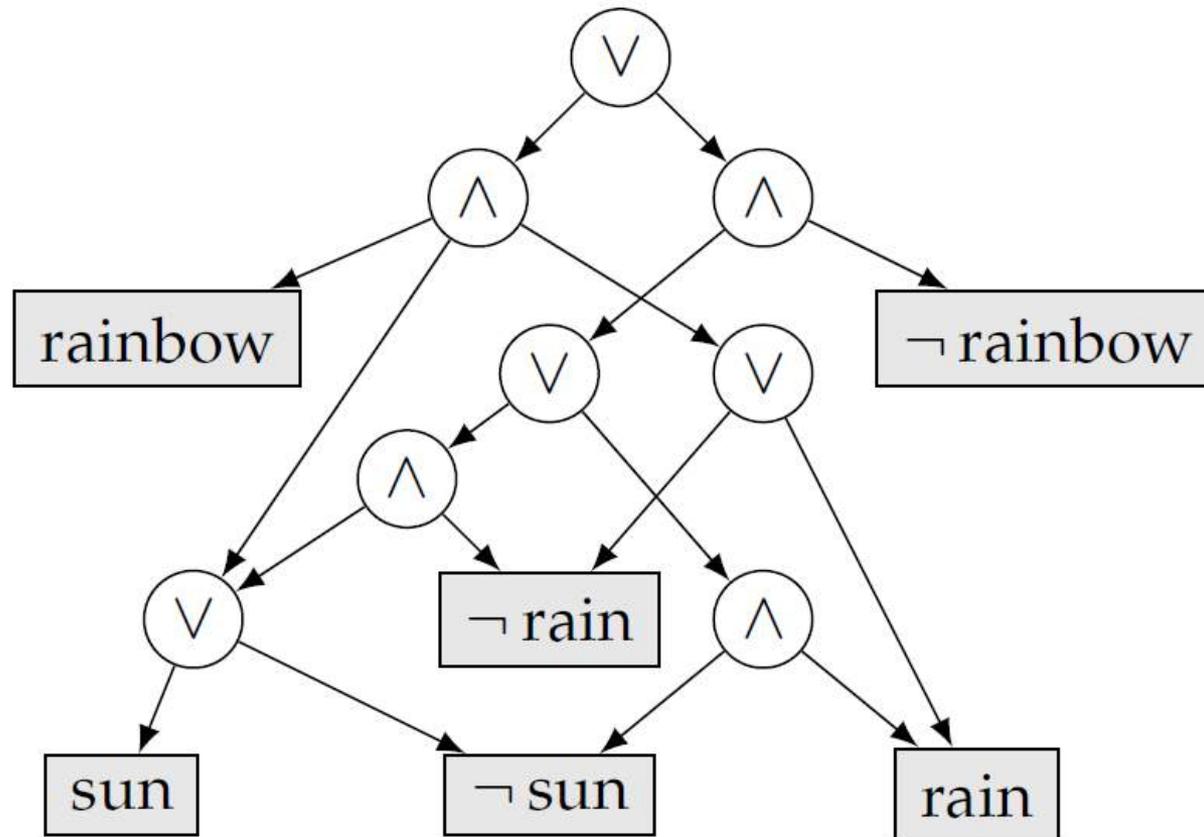
# Tractable for Logical Inference

- Is there a solution? (SAT) ✓
  - $\text{SAT}(\alpha \vee \beta)$  iff  $\text{SAT}(\alpha)$  or  $\text{SAT}(\beta)$  (*always*)
  - $\text{SAT}(\alpha \wedge \beta)$  iff  $\text{SAT}(\alpha)$  and  $\text{SAT}(\beta)$  (*decomposable*)
- How many solutions are there? (#SAT)
- Complexity linear in circuit size 😊

# Deterministic Circuits

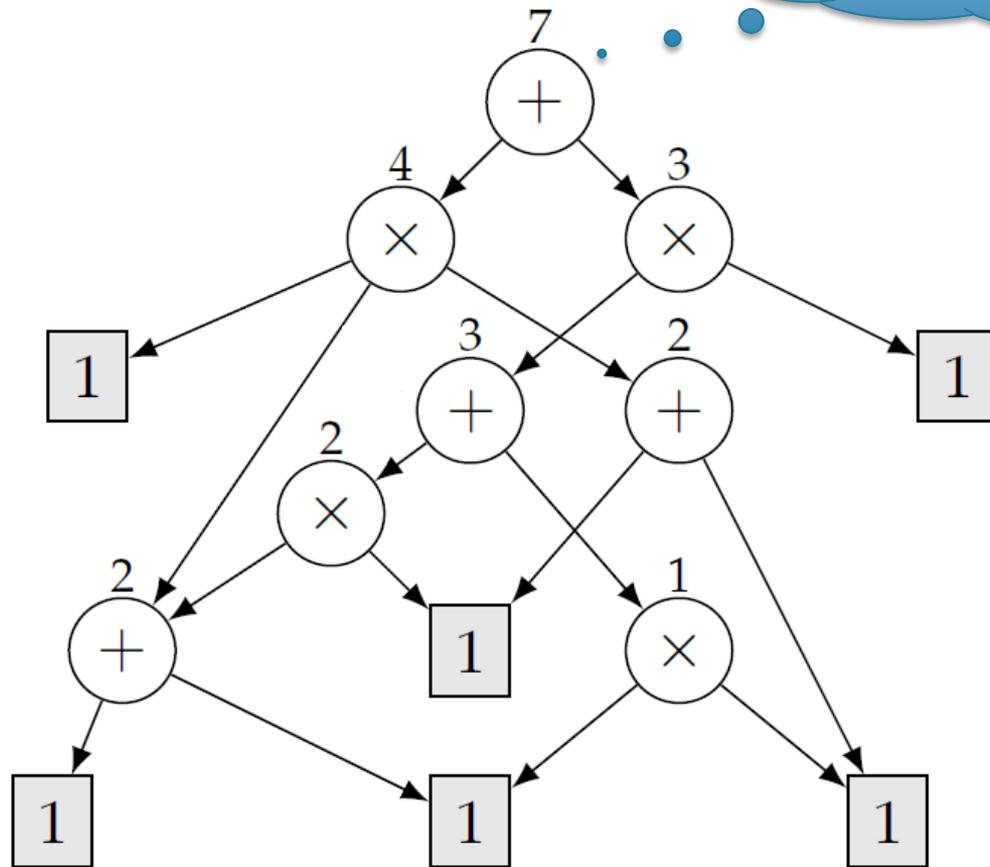


# How many solutions are there? (#SAT)



# How many solutions are there? (#SAT)

Arithmetic Circuit



# Tractable for Logical Inference

- Is there a solution? (SAT) ✓
- How many solutions are there? (#SAT) ✓
- Stricter languages (e.g., BDD, SDD):
  - Equivalence checking ✓
  - Conjoin/disjoint/negate circuits ✓
- Complexity linear in circuit size 😊
- Compilation into circuit language by either
  - ↓ exhaustive SAT solver
  - ↑ conjoin/disjoin/negate

# How to Compute Semantic Loss?

- In general: #P-hard ☹️
- With a logical circuit for  $\alpha$ : Linear!
- Example: exactly-one constraint:

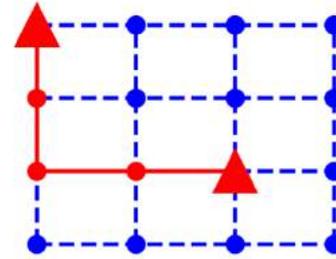
$$L(\alpha, \mathbf{p}) = L(\text{Circuit}, \mathbf{p}) = -\log(\text{Sum of Products})$$

The diagram illustrates the decomposition of the semantic loss for an exactly-one constraint. On the left, a logic circuit is shown with three AND gates and one OR gate. The inputs are  $x_1$ ,  $\neg x_2$ ,  $\neg x_3$ ,  $\neg x_1$ ,  $x_2$ , and  $x_3$ . The gates are:  $G_1$  (inputs  $x_1, \neg x_2$ ),  $G_2$  (inputs  $\neg x_3, \neg x_1$ ),  $G_3$  (inputs  $x_2, x_3$ ), and  $G_4$  (inputs  $G_1, G_2, G_3$ ). On the right, a sum-of-products tree is shown where the root is a plus sign (+) and the leaves are the probabilities  $\Pr(x_1)$ ,  $\Pr(\neg x_2)$ ,  $\Pr(\neg x_3)$ ,  $\Pr(\neg x_1)$ ,  $\Pr(x_2)$ , and  $\Pr(x_3)$ .

- *Why?* Decomposability and determinism!

# Predict Shortest Paths

Add semantic loss  
for path constraint



Test accuracy %	Coherent	Incoherent	Constraint
5-layer MLP	5.62	<b>85.91</b>	6.99
Semantic loss	<b>28.51</b>	83.14	<b>69.89</b>

*Is prediction  
the shortest path?*  
**This is the real task!**

*Are individual  
edge predictions  
correct?*

*Is output  
a path?*

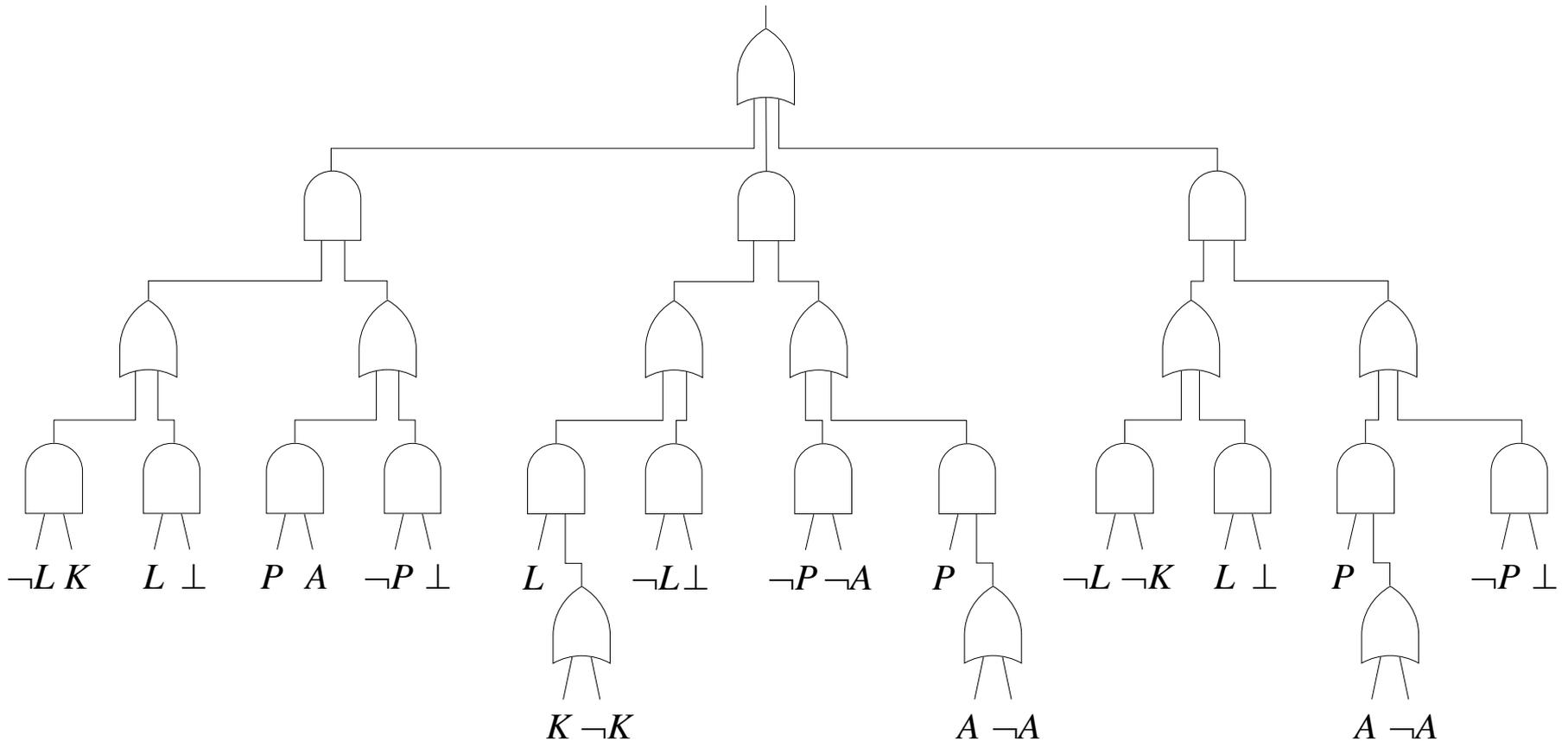
(same conclusion for predicting sushi preferences, see paper)

# Outline

- Adding knowledge to deep learning
- **Probabilistic circuits**
- Logistic circuits for image classification

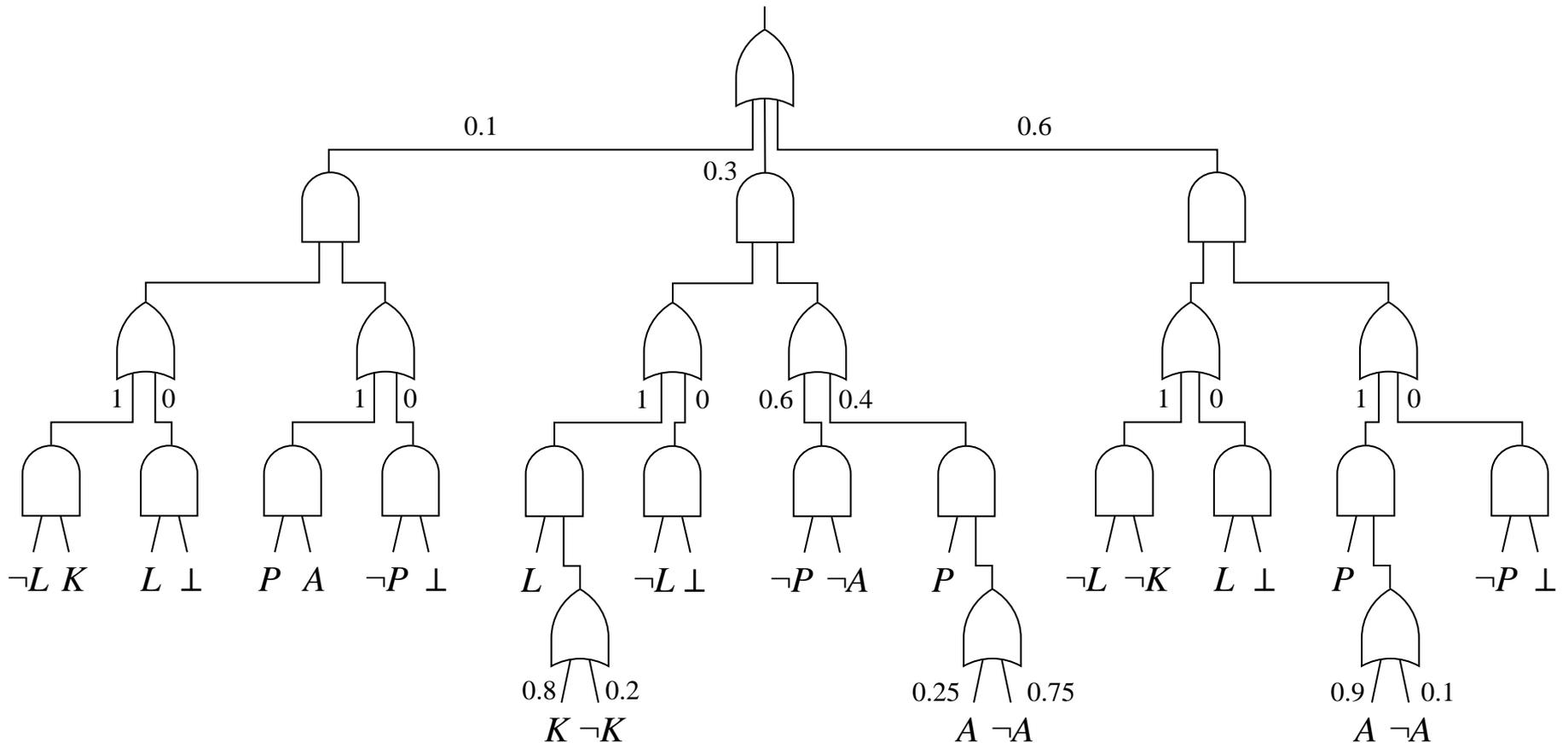
# Logical Circuits

$$P \vee L$$
$$A \Rightarrow P$$
$$K \Rightarrow (P \vee L)$$



Can we represent a **distribution** over the solutions to the constraint?

# Probabilistic Circuits



Syntax: assign a normalized probability to each OR gate input

# Bottom-Up Evaluation of PSDDs

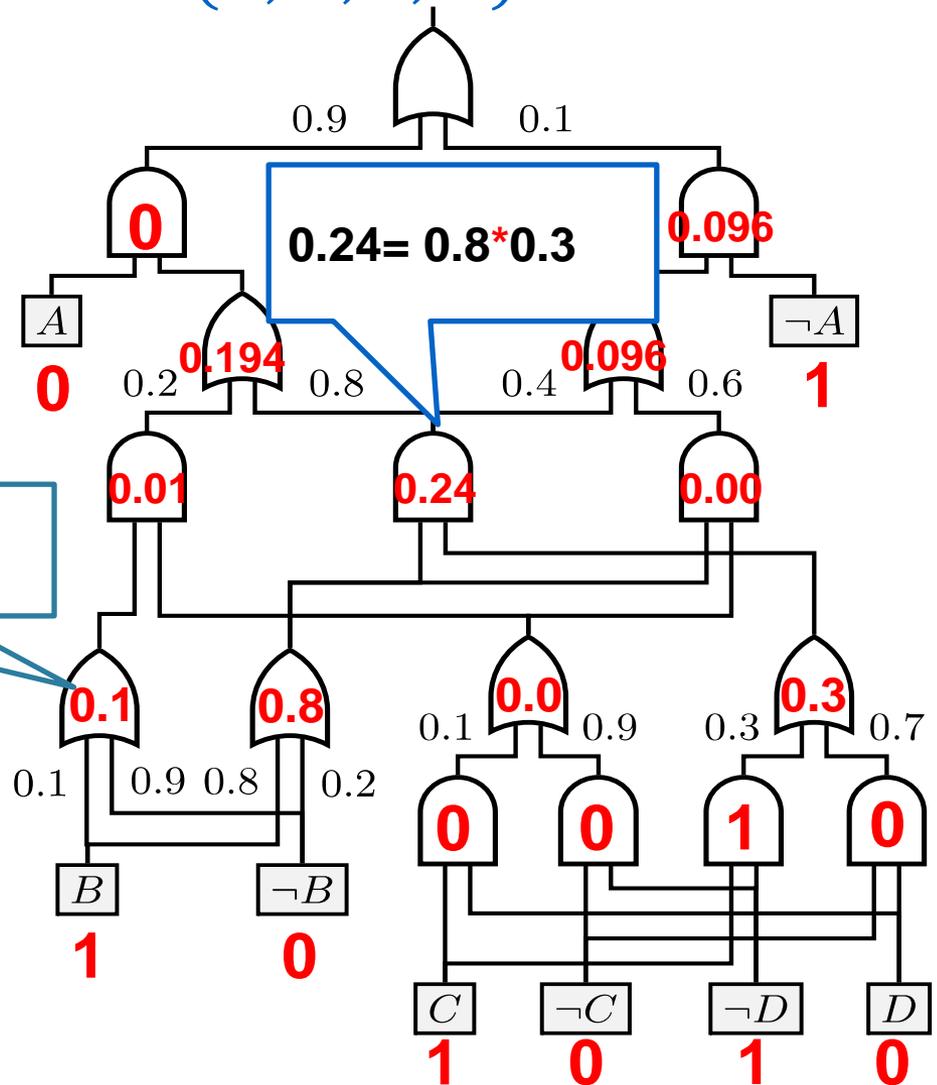
Input:

$A$	$B$	$C$	$D$	$\Pr(A, B, C, D)$
0	1	1	0	?

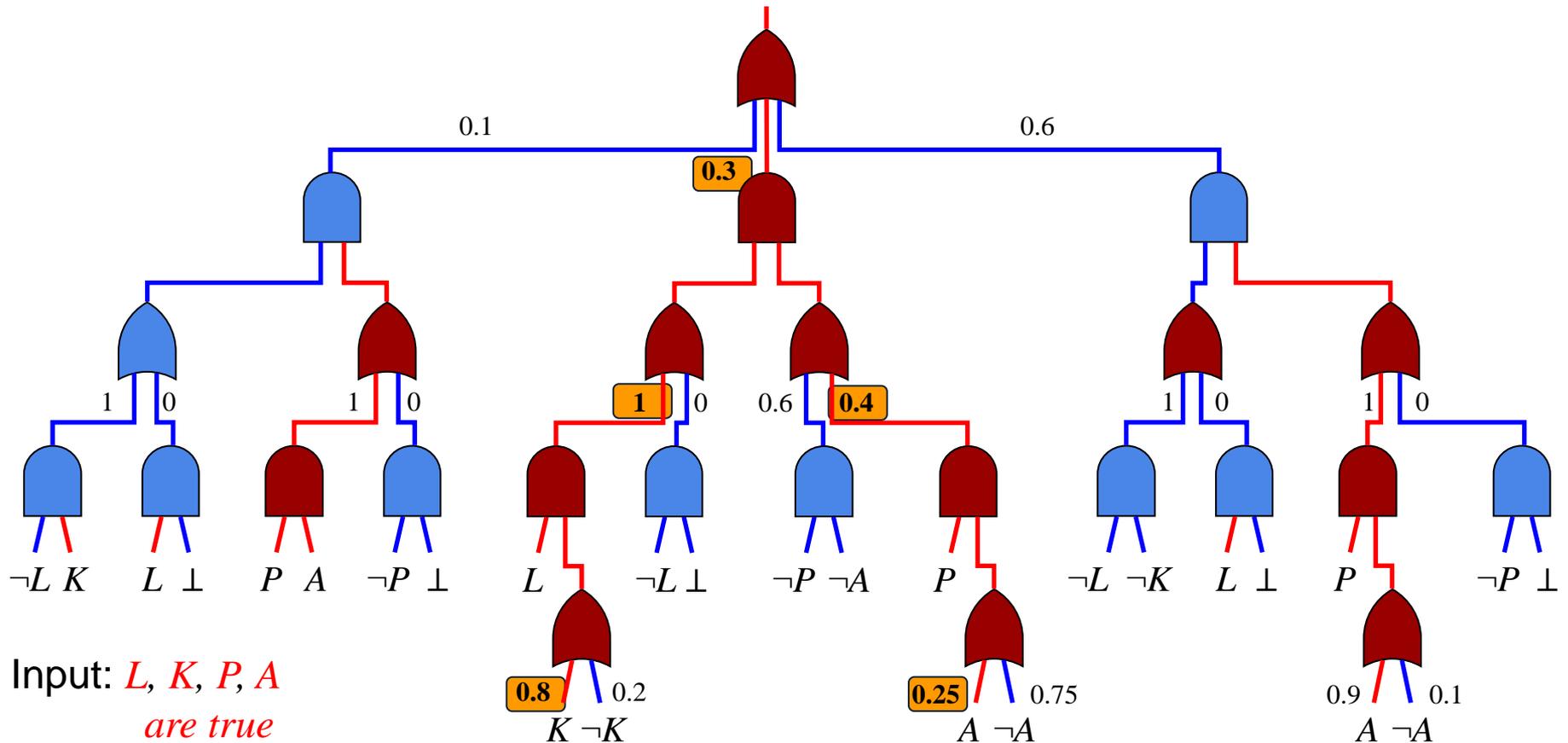
$0.1 = 0.1 * 1 + 0.9 * 0$

Multiply the parameters  
bottom-up

$\Pr(A, B, C, D) = 0.096$

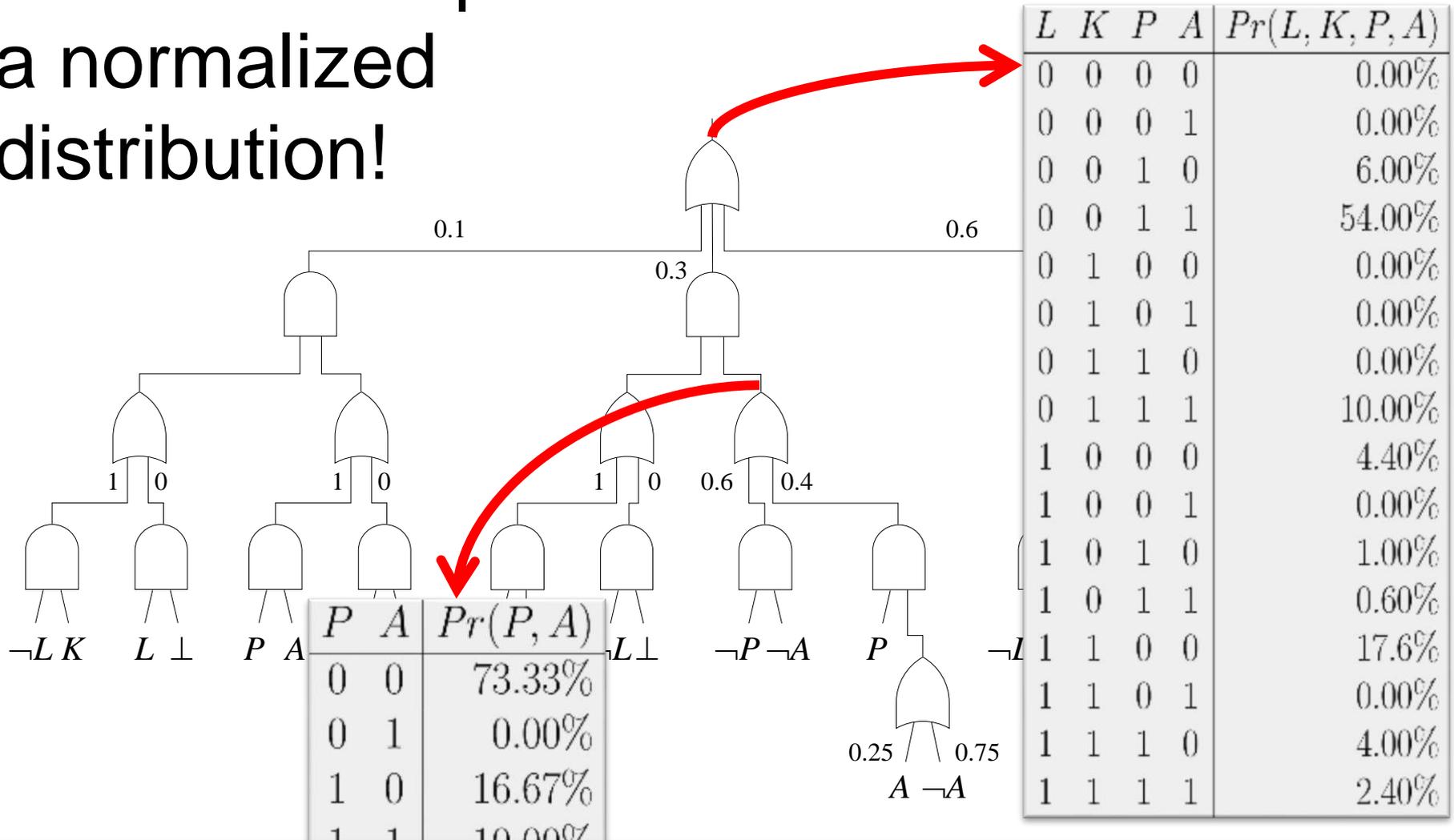


# Alternative View of PSDDs



$$\Pr(L, K, P, A) = 0.3 \times 1 \times 0.8 \times 0.4 \times 0.25 = \mathbf{0.024}$$

# Each node represents a normalized distribution!



Can read probabilistic independences off the circuit structure!

Can interpret every parameter as a conditional probability! (XAI)

# Tractable for Probabilistic Inference

- **MAP inference:**  
Find most-likely assignment to  $x$  given  $y$   
(otherwise NP-hard)
- Computing **conditional probabilities**  $\Pr(x|y)$   
(otherwise #P-hard)
- **Sample** from  $\Pr(x|y)$
- Algorithms linear in circuit size 😊  
(pass up, pass down, similar to backprop)

# Parameter Learning Algorithms

- Closed form  
max likelihood  
from complete data

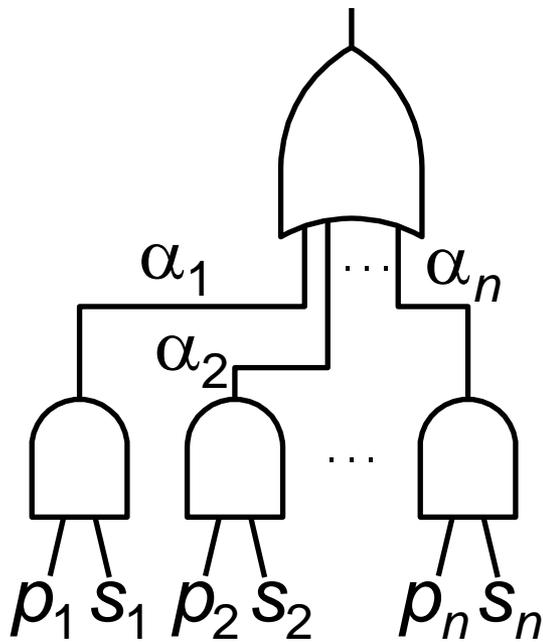
L	K	P	A	Students
0	0	1	0	6
0	0	1	1	54
0	1	1	1	10
1	0	0	0	5
1	0	1	0	1
1	0	1	1	0
1	1	0	0	17
1	1	1	0	4
1	1	1	1	3

- One pass over data to estimate  $\Pr(x|y)$

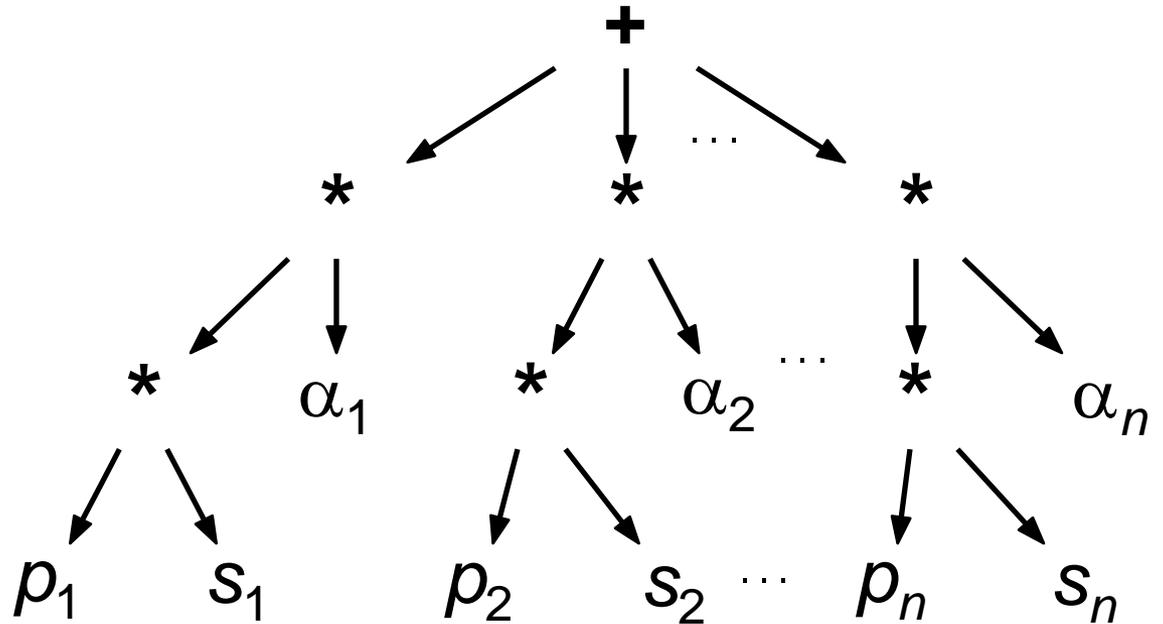
Not a lot to say: very easy! 😊

# PSDDs

...are Sum-Product Networks  
...are Arithmetic Circuits



**PSDD**



**AC**

# Learn Mixtures of PSDD Structures

Datasets	Var	LearnPSDD Ensemble	Best-to-Date
NLTCs	16	-5.99 <sup>†</sup>	-6.00
MSNBC	17	-6.04 <sup>†</sup>	-6.04 <sup>†</sup>
KDD	64	-2.11 <sup>†</sup>	-2.12
Plants	69	-13.02	-11.99 <sup>†</sup>
Audio	100	-39.94	-39.49 <sup>†</sup>
Jester	100	-51.29	-41.11 <sup>†</sup>
Netflix	100	-55.71 <sup>†</sup>	-55.84
Accidents	111	-30.16	-24.87 <sup>†</sup>
Retail	135	-10.72 <sup>†</sup>	-10.78
Pumsb-Star	163	-26.12	-22.40 <sup>†</sup>
DNA	180	-88.01	-80.03 <sup>†</sup>
Kosarek	190	-10.52 <sup>†</sup>	-10.54
MSWeb	294	-9.89	-9.22 <sup>†</sup>
Book	500	-34.97	-30.18 <sup>†</sup>
EachMovie	500	-58.01	-51.14 <sup>†</sup>
WebKB	839	-161.09	-150.10 <sup>†</sup>
Reuters-52	889	-89.61	-80.66 <sup>†</sup>
20NewsGrp.	910	-155.97	-150.88 <sup>†</sup>
BBC	1058	-253.19	-233.26 <sup>†</sup>
AD	1556	-31.78	-14.36 <sup>†</sup>

State of the art  
on 6 datasets!

Q: “Help! I need to learn a discrete probability distribution...”

A: Learn mixture of PSDDs!

Strongly outperforms

- Bayesian network learners
- Markov network learners

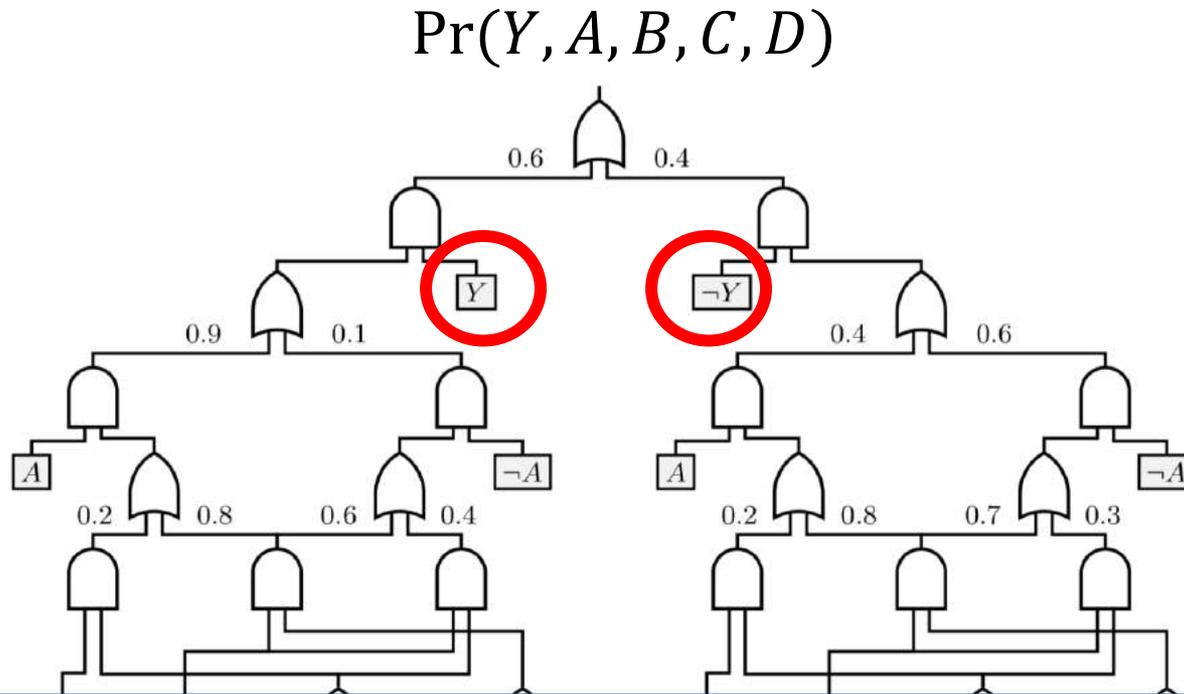
Competitive with

- SPN learners
- Cutset network learners

# Outline

- Adding knowledge to deep learning
- Probabilistic circuits
- **Logistic circuits for image classification**

# What if I only want to classify Y?



**What if we only want to learn a classifier  $\Pr(Y|X)$**



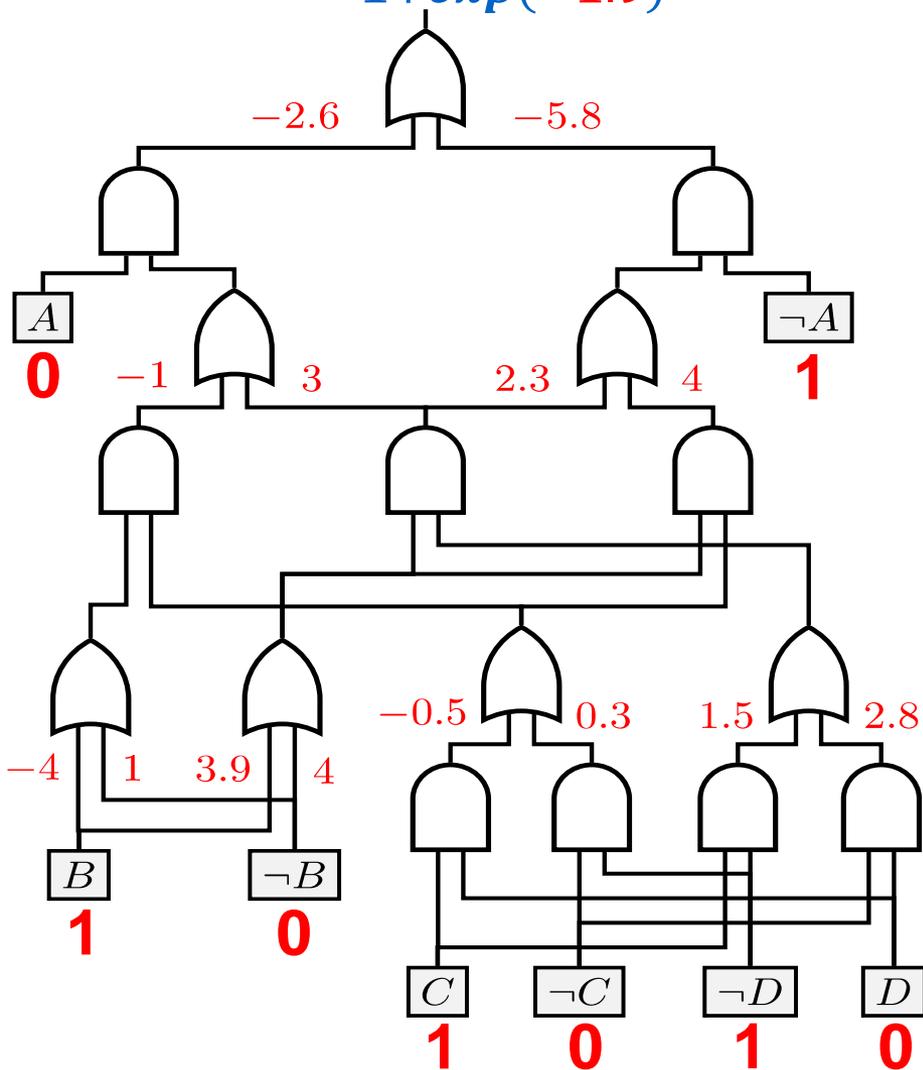
# Logistic Circuits: Evaluation

$$\Pr(Y = 1 \mid A, B, C, D) = \frac{1}{1 + \exp(-1.9)} = 0.869$$

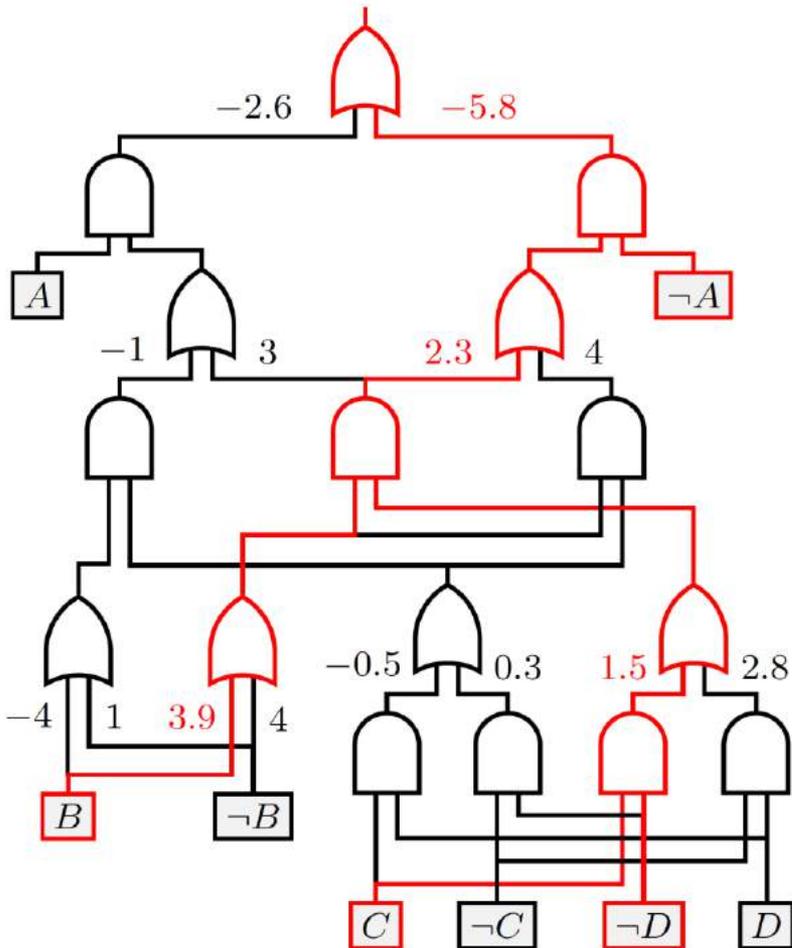
Input:

$A$	$B$	$C$	$D$	$\Pr(Y \mid A, B, C, D)$
0	1	1	0	?

Aggregate the parameters  
bottom-up  
Logistic function on final  
output



# Alternative View on Logistic Circuits

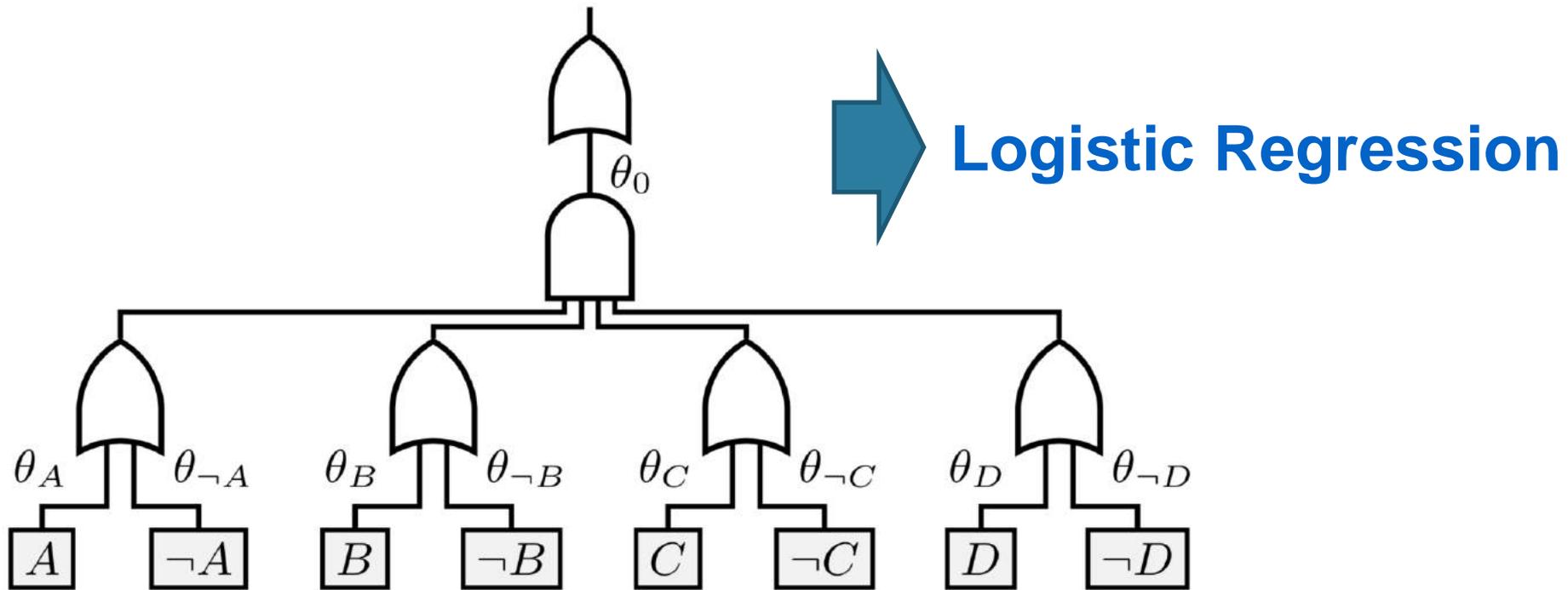


Represents  $\Pr(Y | A, B, C, D)$

- Take all 'hot' wires
- Sum their weights
- Push through logistic function

$A$	$B$	$C$	$D$	$g_r(ABCD)$	$\Pr(Y = 1   ABCD)$
1	0	1	1	-3.1	4.31%
0	1	1	0	1.9	86.99%
1	1	1	0	5.8	99.70%

# Special Case: Logistic Regression



$$\Pr(Y = 1|A, B, C, D) = \frac{1}{1 + \exp(-A * \theta_A - \neg A * \theta_{\neg A} - B * \theta_B - \dots)}$$

**What about other logistic circuits in more general forms?**

# Parameter Learning

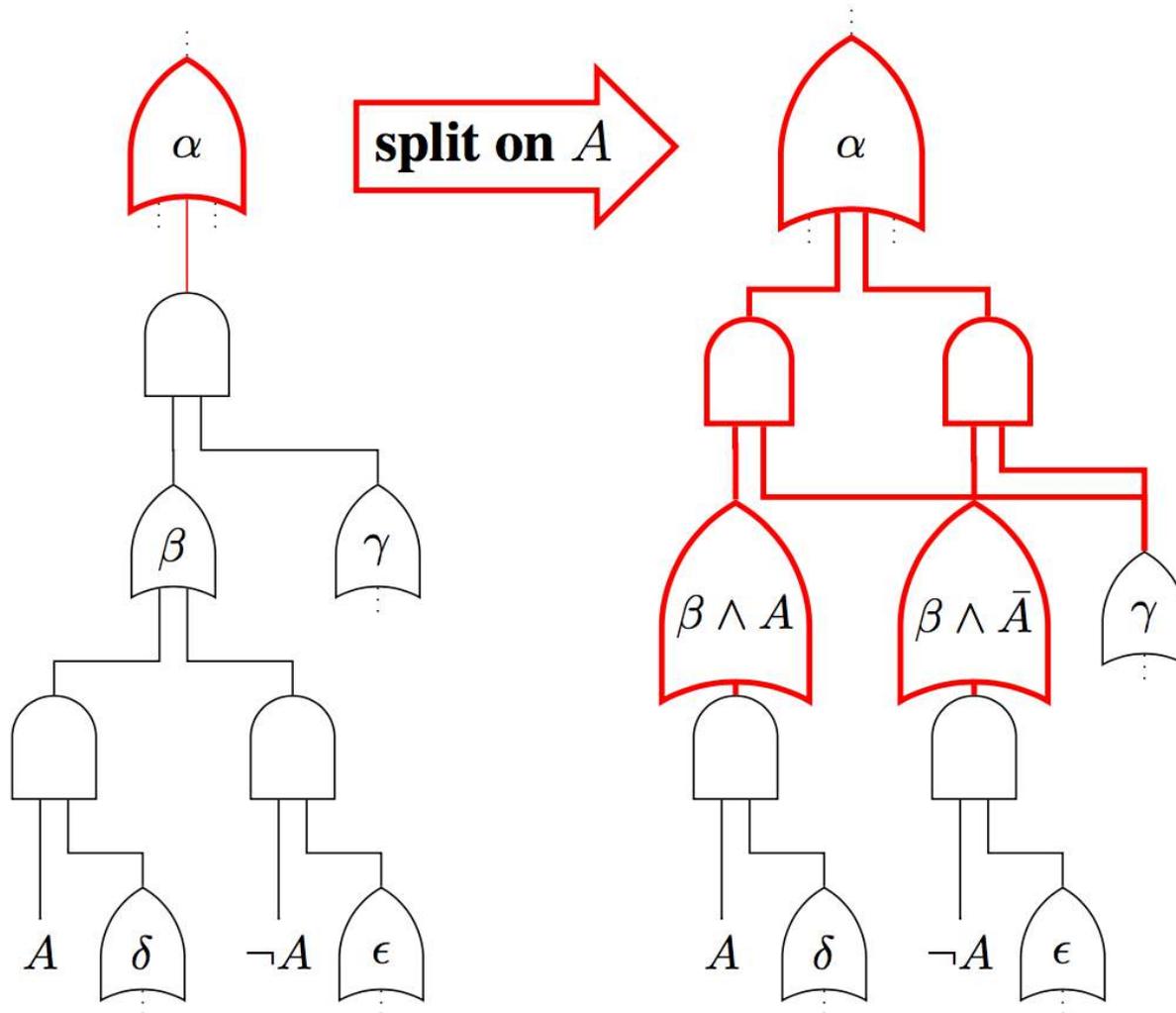
Reduce to logistic regression:

$$\Pr(Y = 1 \mid \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{x} \cdot \boldsymbol{\theta})}$$

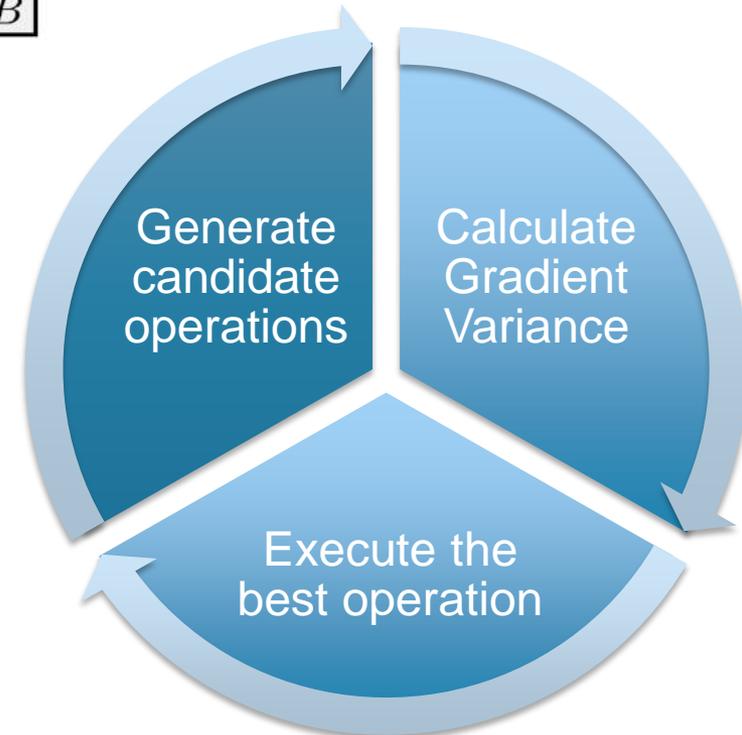
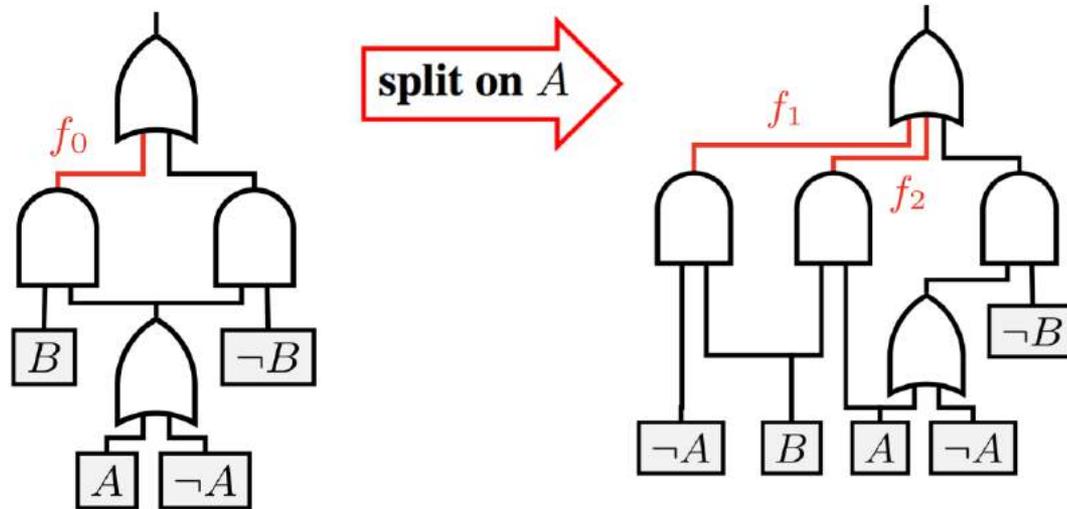
Features associated with each wire  
“Global Circuit Flow” features

Learning parameters  $\theta$  is convex optimization!

# Structure Learning Primitive



# Logistic Circuit Structure Learning



# Comparable Accuracy with Neural Nets

ACCURACY % ON DATASET	MNIST	FASHION
BASELINE: LOGISTIC REGRESSION	85.3	79.3
BASELINE: KERNEL LOGISTIC REGRESSION	97.7	88.3
RANDOM FOREST	97.3	81.6
3-LAYER MLP	97.5	84.8
RAT-SPN (PEHARZ ET AL. 2018)	98.1	89.5
SVM WITH RBF KERNEL	98.5	87.8
5-LAYER MLP	99.3	89.8
LOGISTIC CIRCUIT (BINARY)	97.4	87.6
LOGISTIC CIRCUIT (REAL-VALUED)	99.4	91.3
CNN WITH 3 CONV LAYERS	99.1	90.7
RESNET (HE ET AL. 2016)	99.5	93.6

# Significantly Smaller in Size

NUMBER OF PARAMETERS	MNIST	FASHION
BASELINE: LOGISTIC REGRESSION	<1K	<1K
BASELINE: KERNEL LOGISTIC REGRESSION	1,521 K	3,930K
LOGISTIC CIRCUIT (REAL-VALUED)	182K	467K
LOGISTIC CIRCUIT (BINARY)	268K	614K
3-LAYER MLP	1,411K	1,411K
RAT-SPN (PEHARZ ET AL. 2018)	8,500K	650K
CNN WITH 3 CONV LAYERS	2,196K	2,196K
5-LAYER MLP	2,411K	2,411K
RESNET (HE ET AL. 2016)	4,838K	4,838K

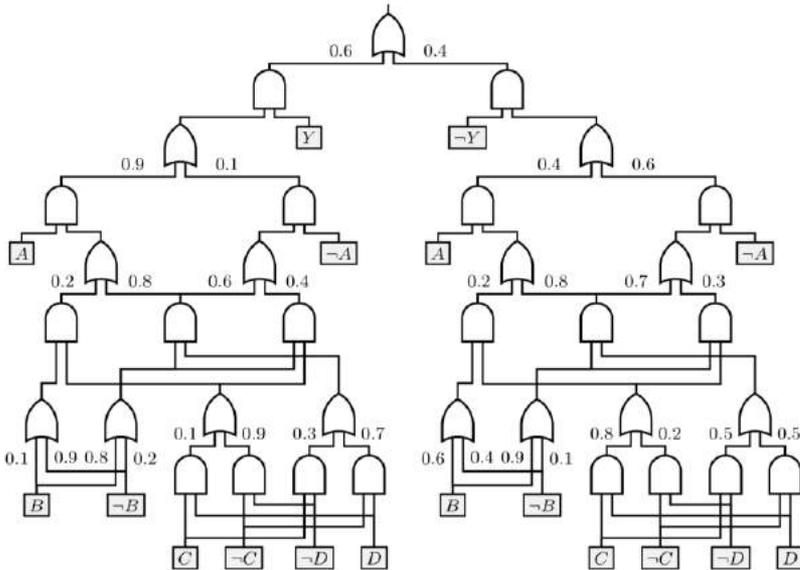
# Better Data Efficiency

ACCURACY % WITH % OF TRAINING DATA	MNIST			FASHION		
	100%	10%	2%	100%	10%	2%
5-LAYER MLP	99.3	<b>98.2</b>	94.3	89.8	86.5	80.9
CNN WITH 3 CONV LAYERS	99.1	98.1	95.3	90.7	87.6	83.8
LOGISTIC CIRCUIT (BINARY)	97.4	96.9	94.1	87.6	86.7	83.2
LOGISTIC CIRCUIT (REAL-VALUED)	<b>99.4</b>	97.6	<b>96.1</b>	<b>91.3</b>	<b>87.8</b>	<b>86.0</b>

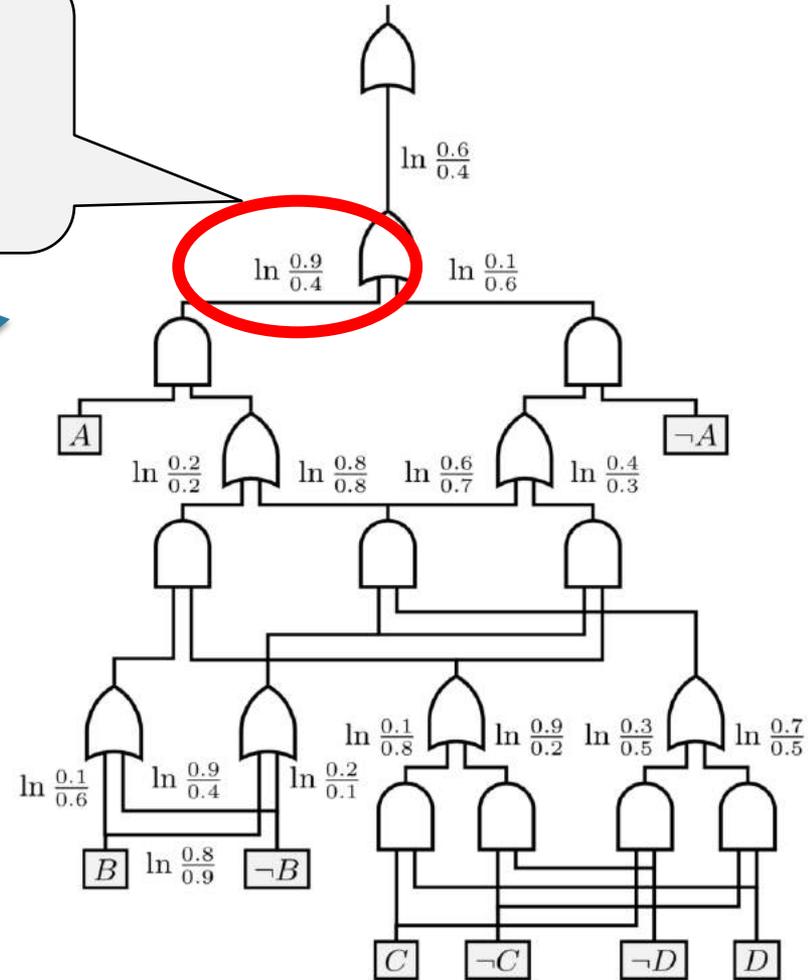
# Logistic vs. Probabilistic Circuits

Probabilities become log-odds

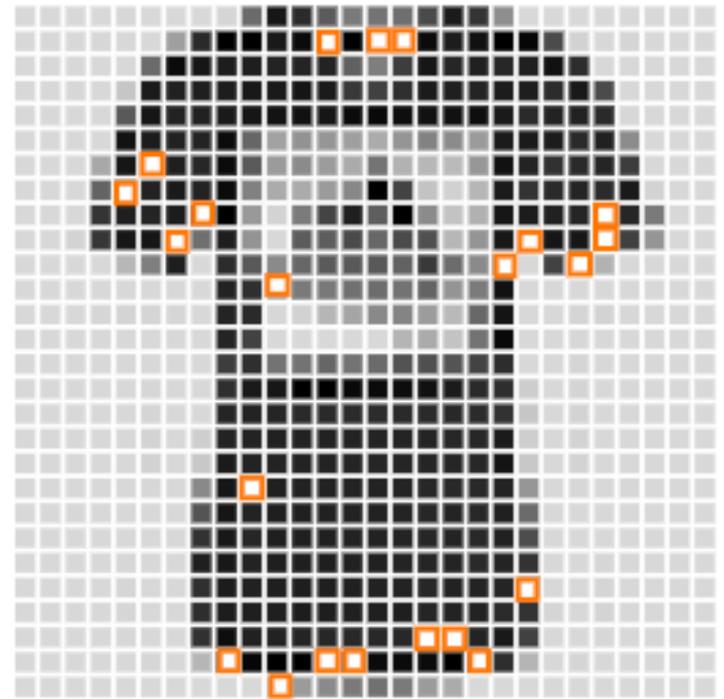
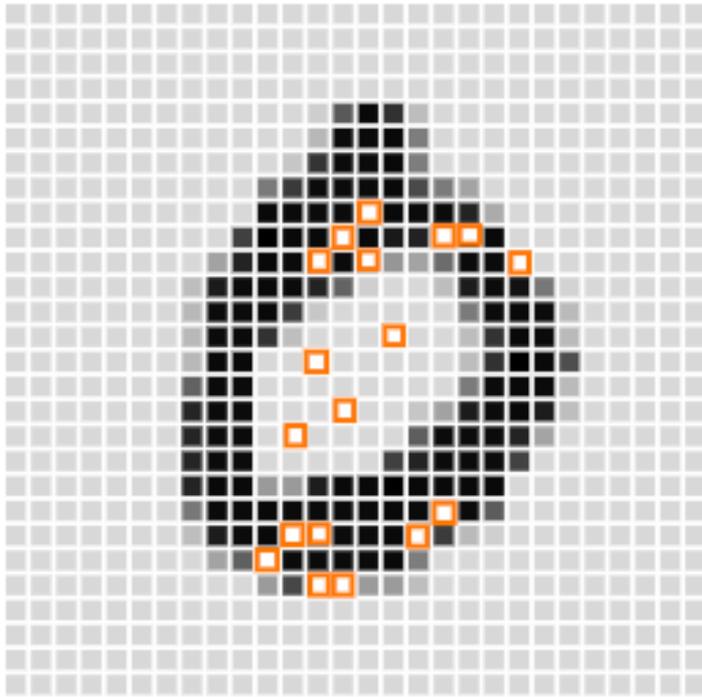
$\Pr(Y, A, B, C, D)$



$\Pr(Y | A, B, C, D)$



# Interpretable?



# Logistic Circuits: Conclusions

- Synthesis of symbolic AI and statistical learning
- Discriminative counterparts of probabilistic circuits
- Convex parameter learning
- Simple heuristic for structure learning
- Good performance
- Easy to interpret

# Conclusions

