

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

Guy Van den Broeck

**UCLA**

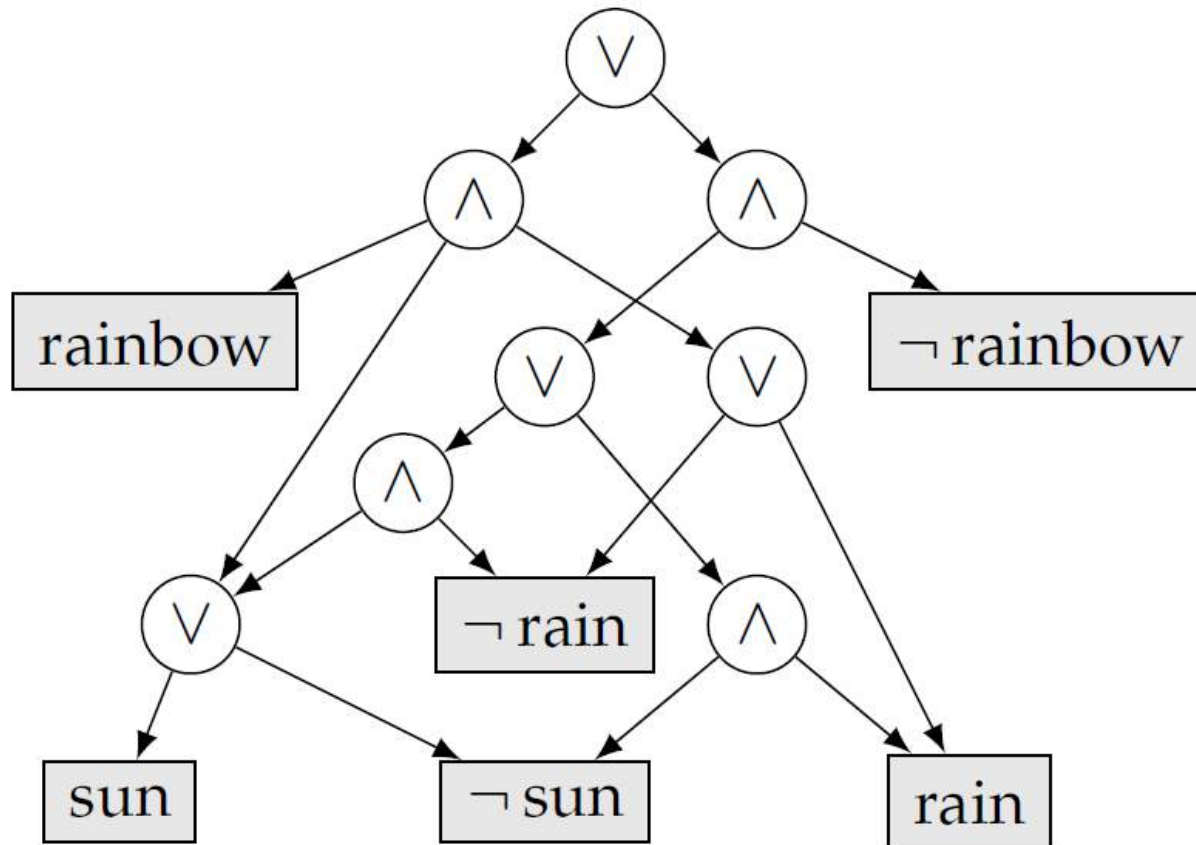
Stanford  
Nov 14, 2018



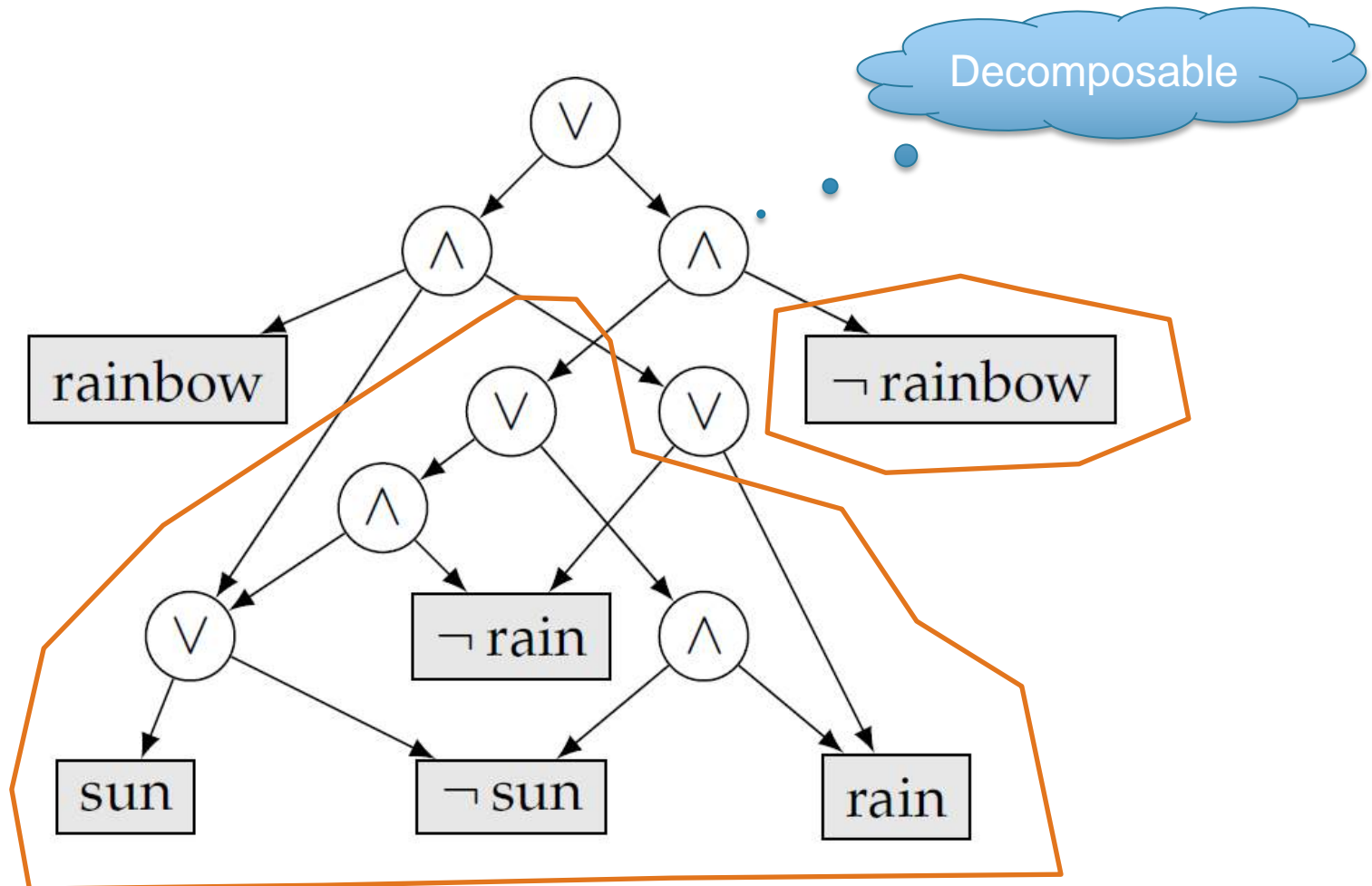
***Foundation:  
Logical Circuit Languages***

# Negation Normal Form Circuits

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



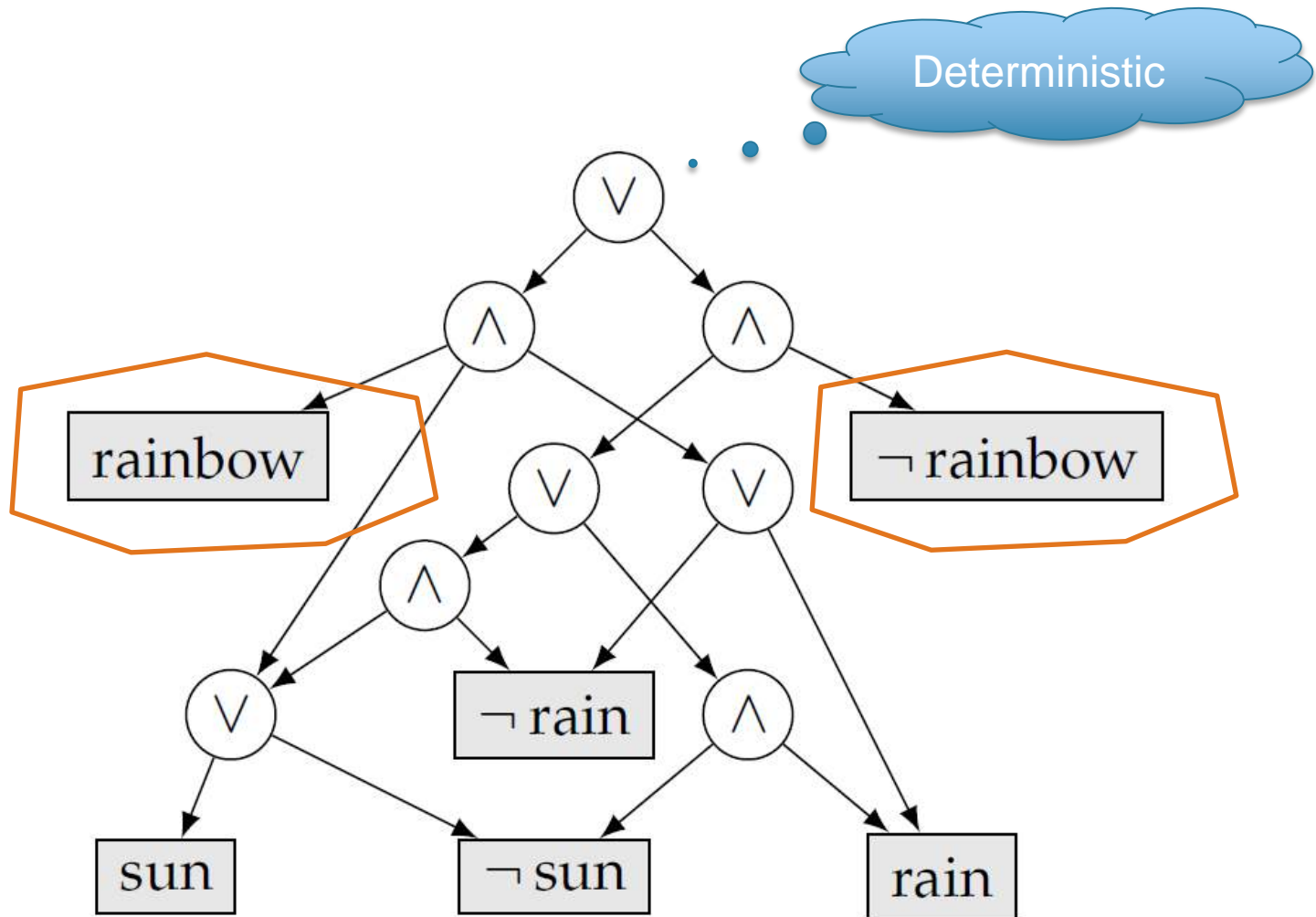
# Decomposable Circuits



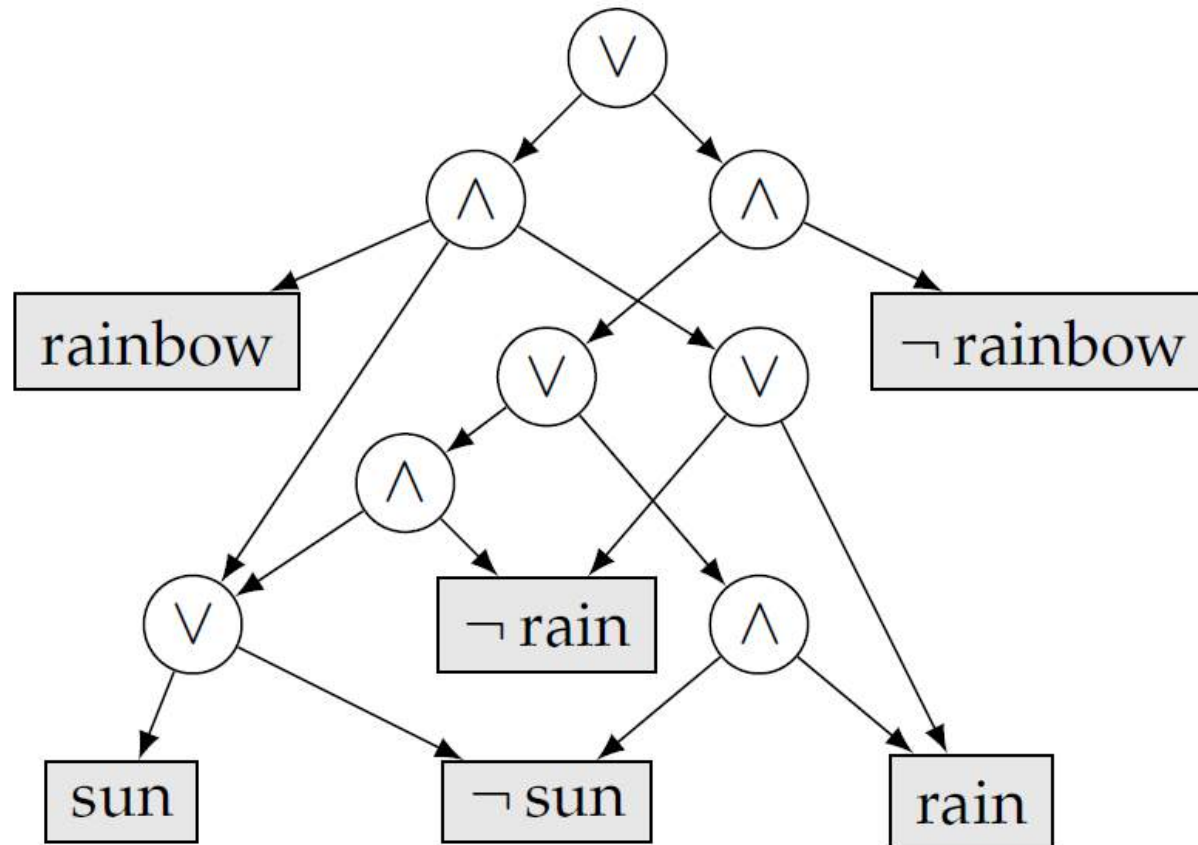
# Tractable for Logical Inference

- 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

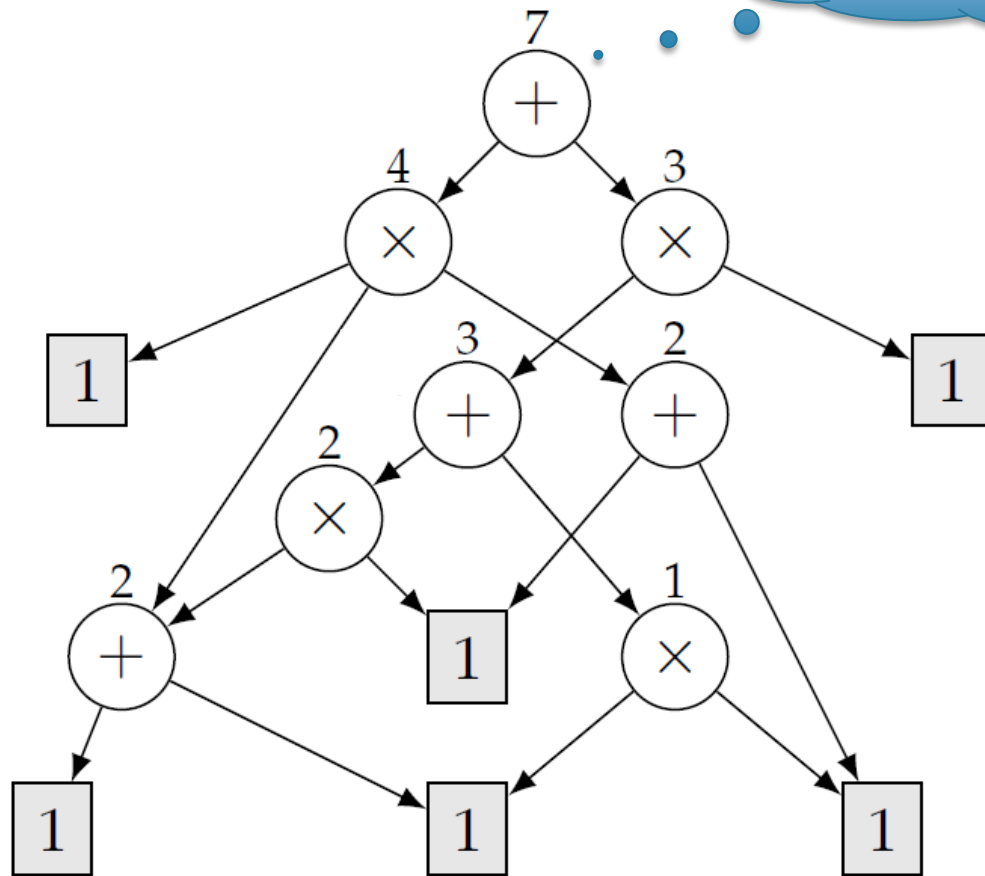


# How many solutions are there? (#SAT)



# How many solutions are there? (#SAT)

Arithmetic Circuit



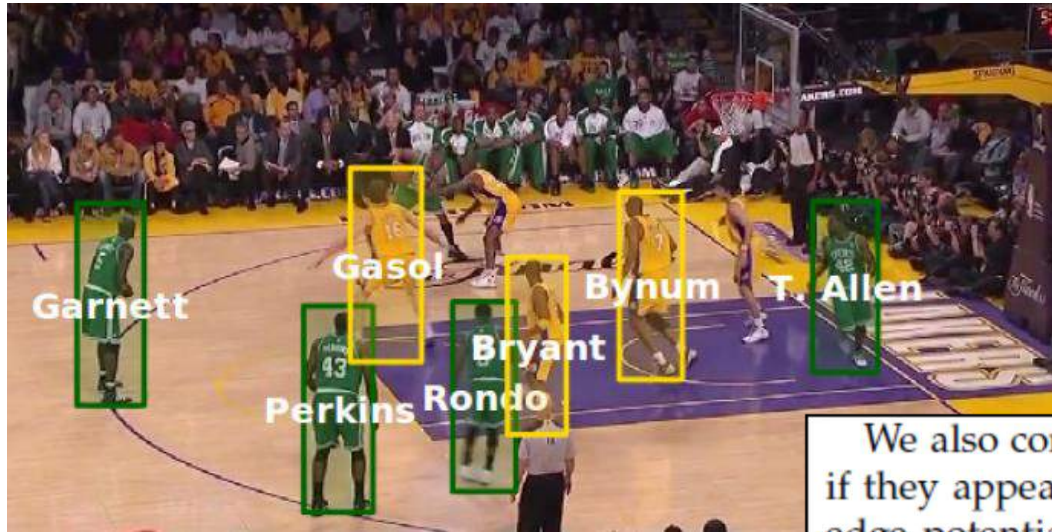


# Tractable for Logical Inference

- 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

***Learning with  
Logical Constraints***

# Motivation: Video

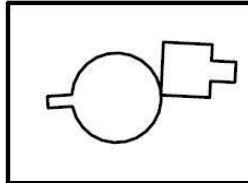
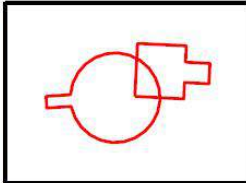
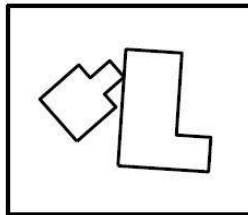
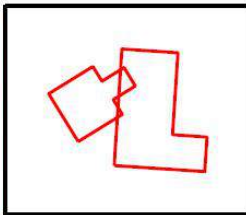
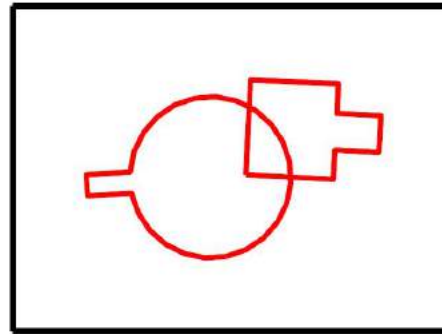
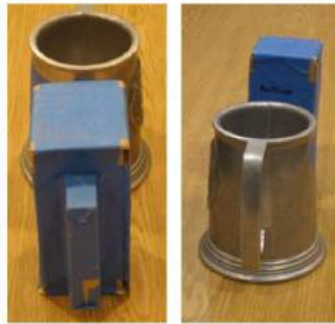


We also connect all pairs of identity nodes  $y_{t,i}$  and  $y_{t,j}$  if they appear in the same time  $t$ . We then introduce an edge potential that enforces mutual exclusion:

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

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



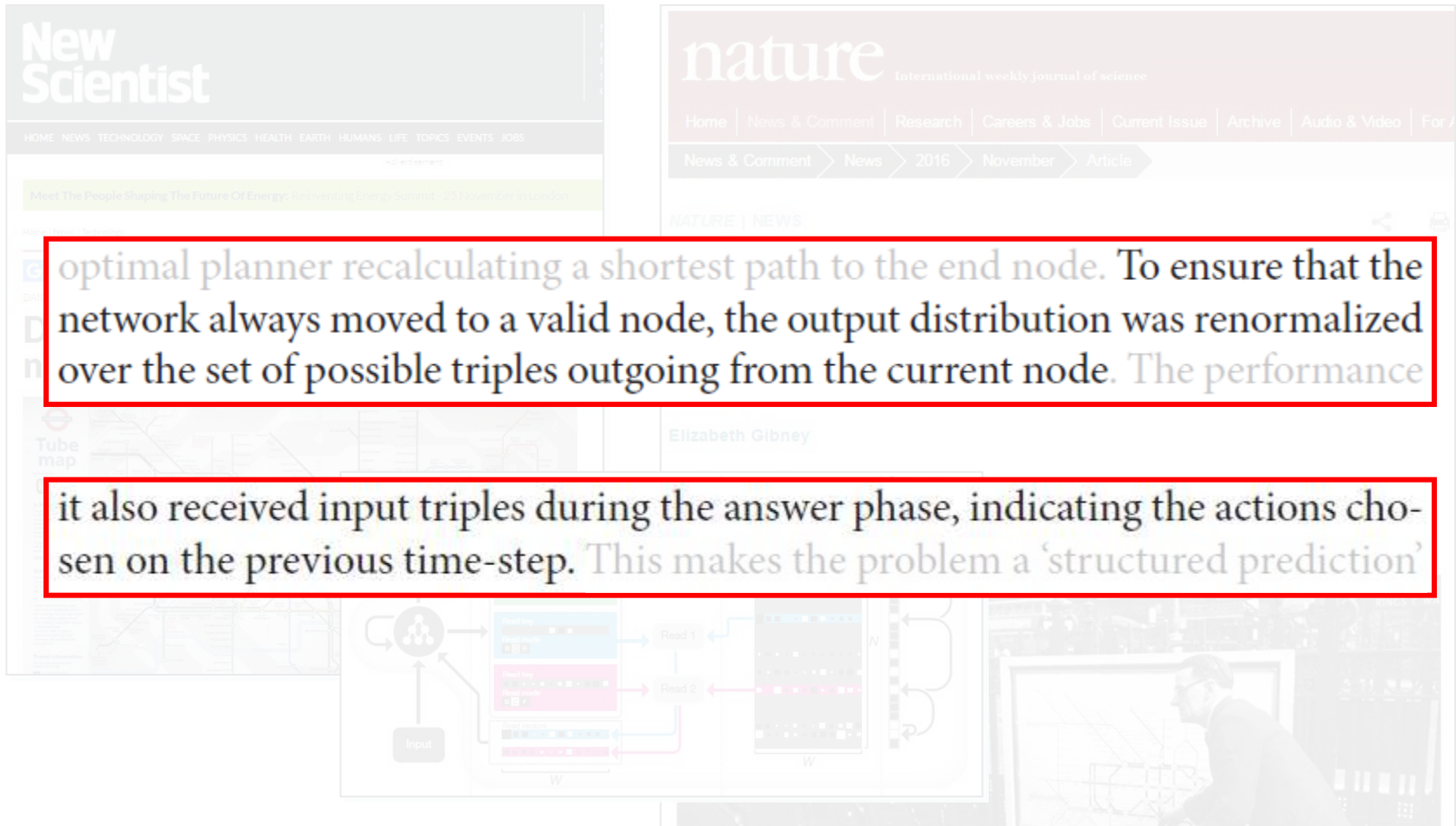
The method developed in this paper can be used in a broad variety of semantic mapping and object manipulation tasks, providing an efficient and effective way to incorporate collision constraints into a recursive state estimator, obtaining optimal or near-optimal solutions.

# 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)

## 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

## 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.

# 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



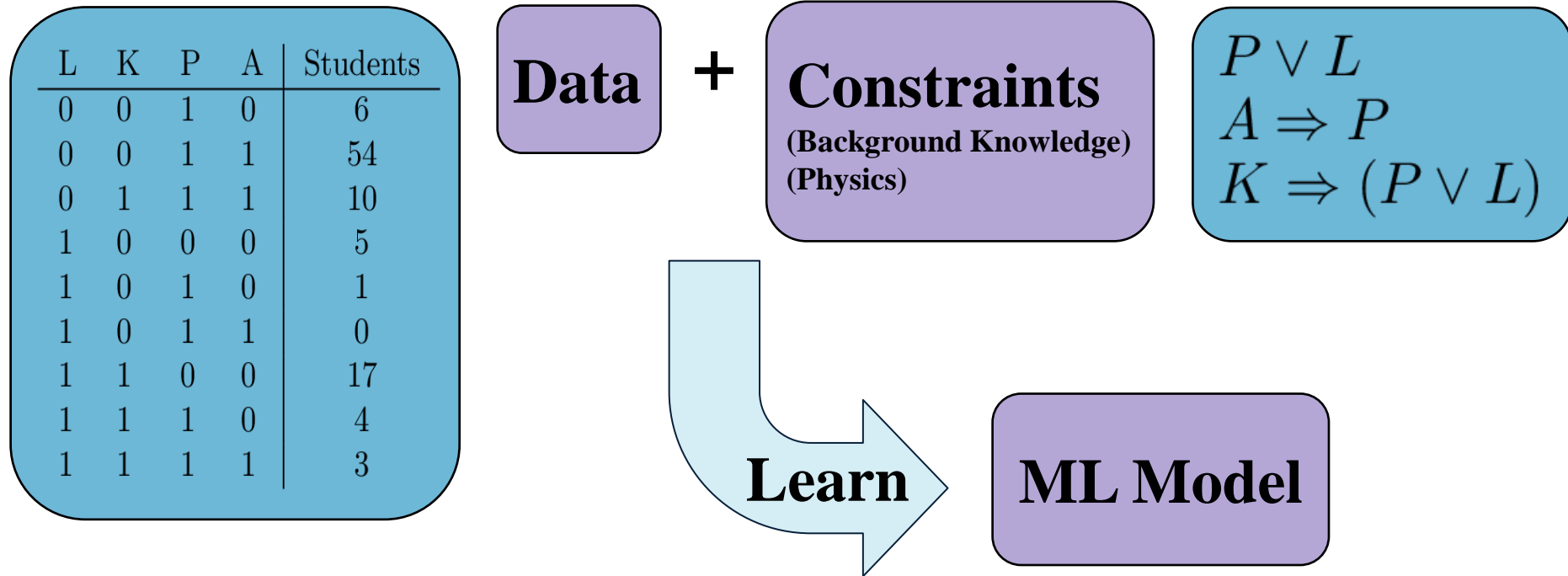
structured

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**

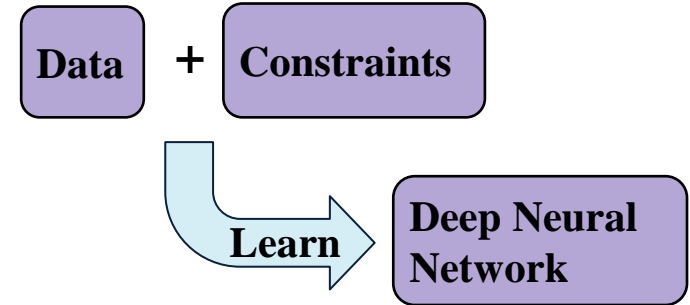
# Learning in Structured Spaces



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

# ***Deep Learning with Logical Constraints***

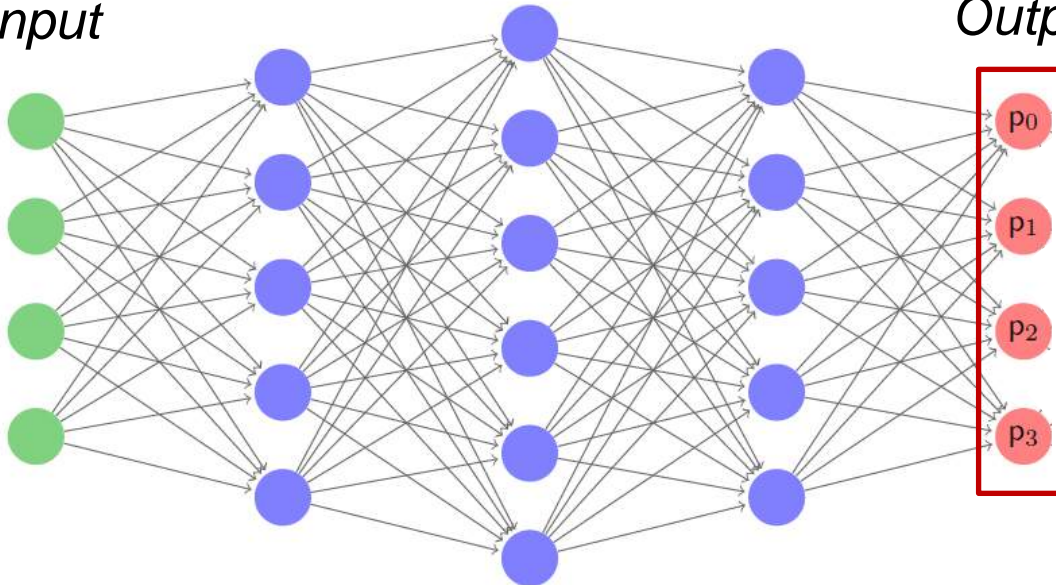
# Deep Learning with Logical Knowledge



*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 Loss!**
  - If  $\mathbf{p}$  is Boolean and satisfies  $\alpha$  then  $L(\alpha, \mathbf{p}) = 0$

# 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:

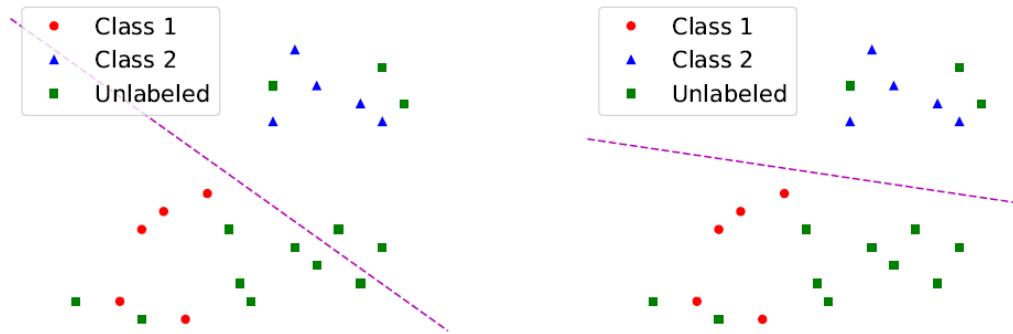
$$L^s(\text{exactly-one}, p) \propto -\log \sum_{i=1}^n p_i \prod_{j=1, j \neq i}^n (1 - p_j)$$

Only  $x_i = 1$  after flipping coins

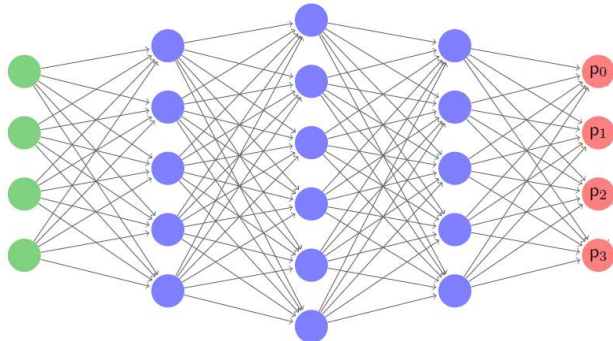
Exactly one true  $x$  after flipping coins

# 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 +  $w \cdot \text{semantic loss}$*



# MNIST Experiment



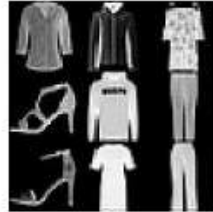
Accuracy % with # of used labels	100	1000	ALL
AtlasRBF (Pitelis et al., 2014)	91.9 ( $\pm 0.95$ )	96.32 ( $\pm 0.12$ )	98.69
Deep Generative (Kingma et al., 2014)	96.67( $\pm 0.14$ )	97.60( $\pm 0.02$ )	99.04
Virtual Adversarial (Miyato et al., 2016)	97.67	98.64	99.36
Ladder Net (Rasmus et al., 2015)	<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

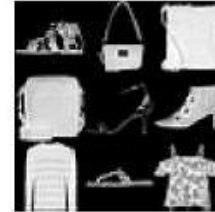
# FASHION Experiment



(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

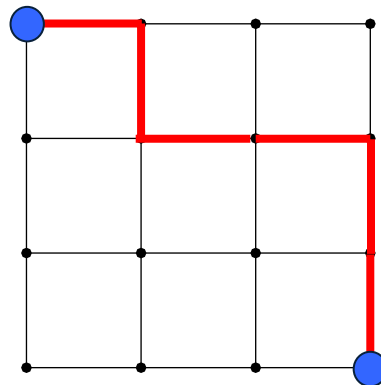
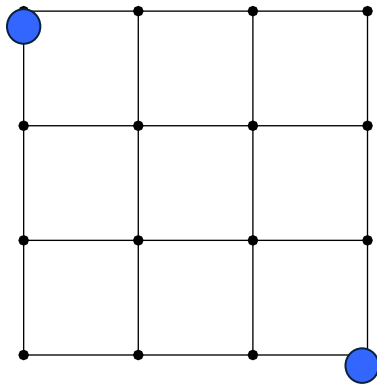
Outperforms Ladder Nets!

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

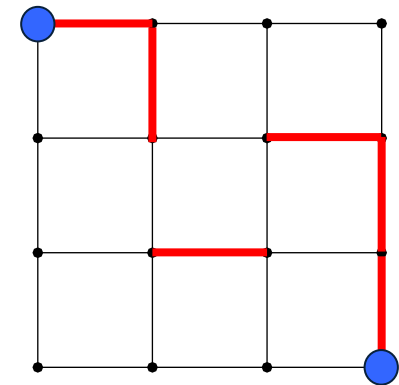
# What about real constraints? Paths

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?

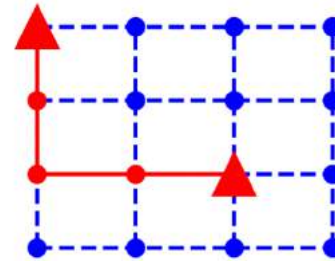
- In general: #P-hard ☹️
- With a logical circuit for  $\alpha$ : Linear!
- Example: exactly-one constraint:

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

- *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?*

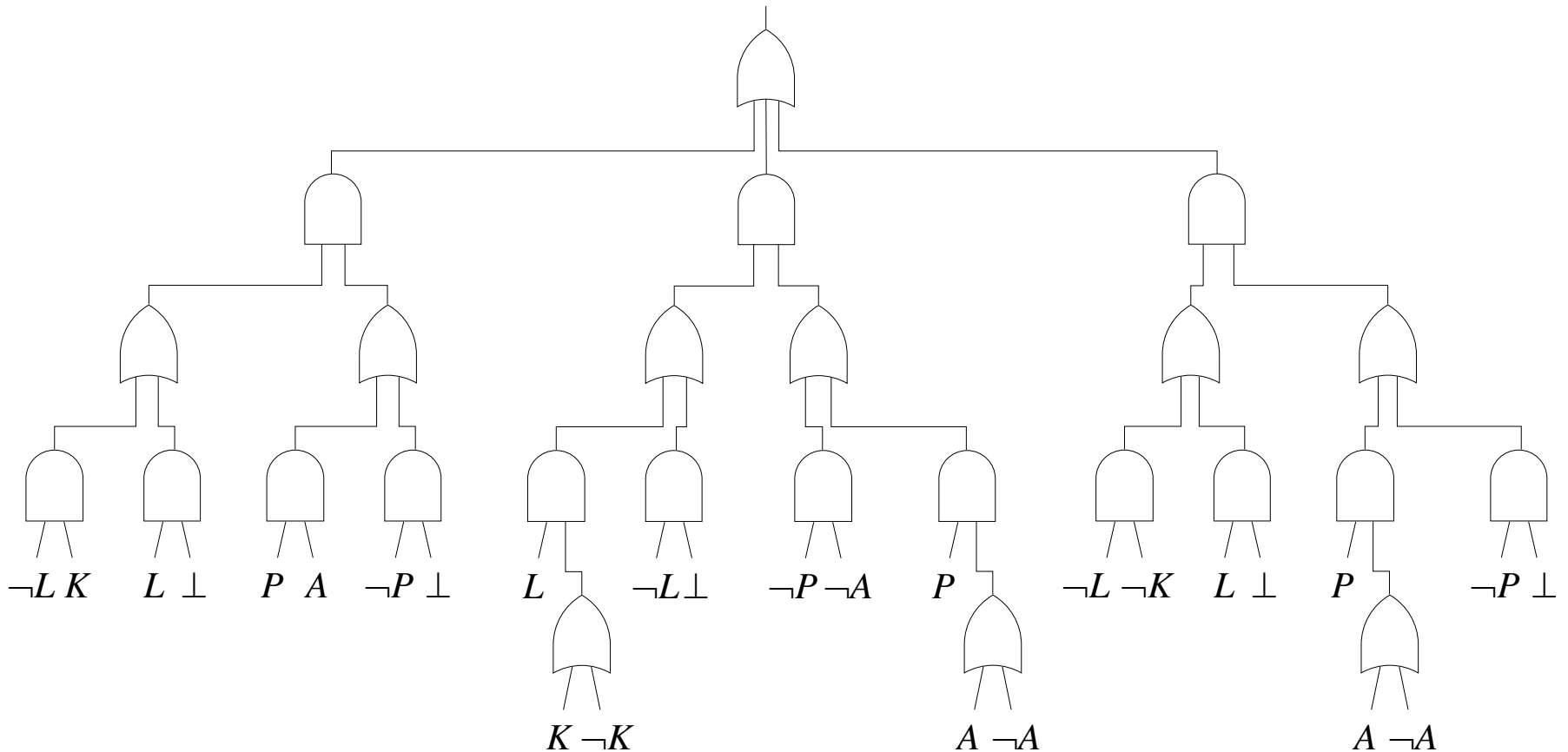
*Is output  
a path?*

(same conclusion for predicting sushi preferences, see paper)

# ***Probabilistic Circuits***

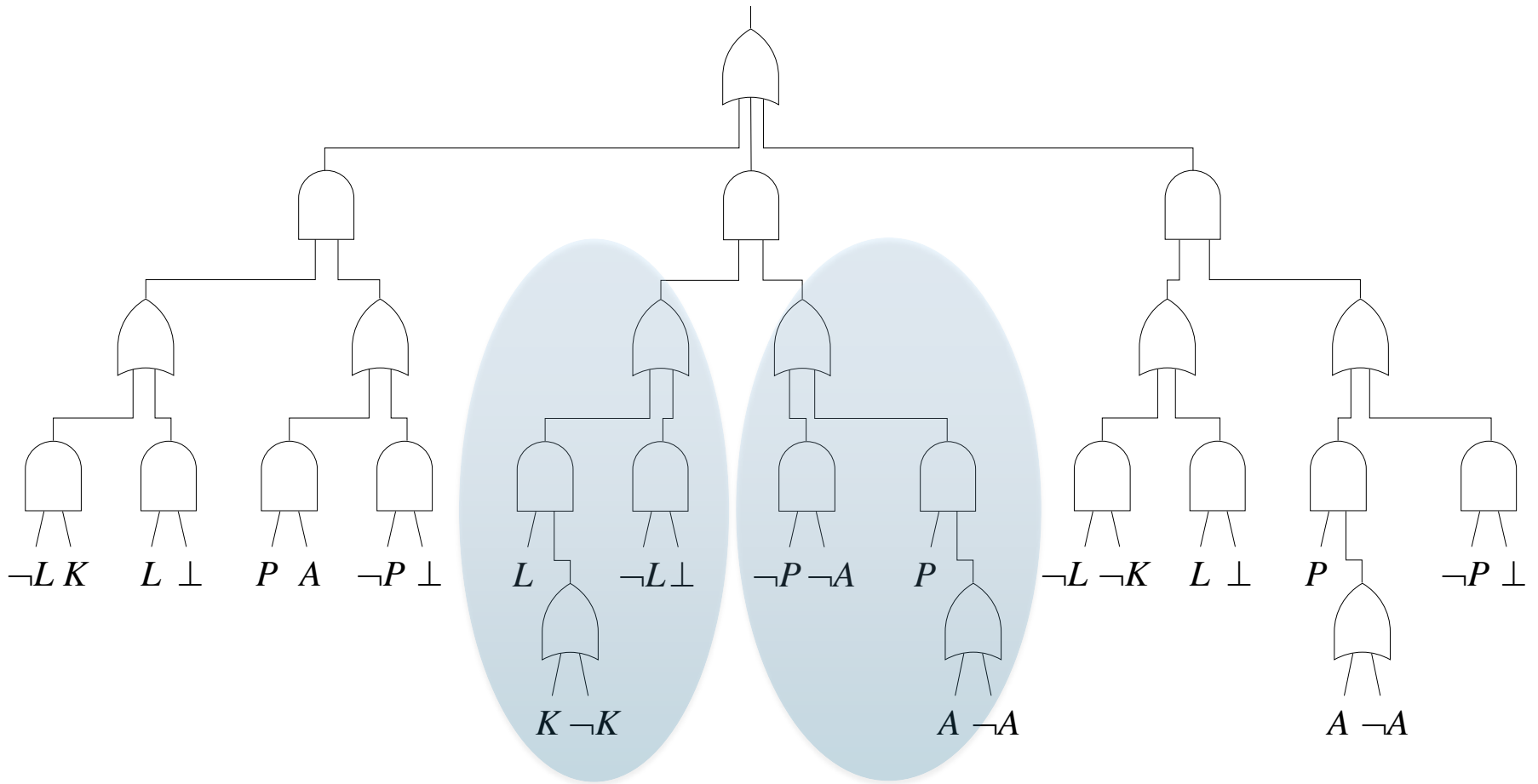
# Logical Circuits

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



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

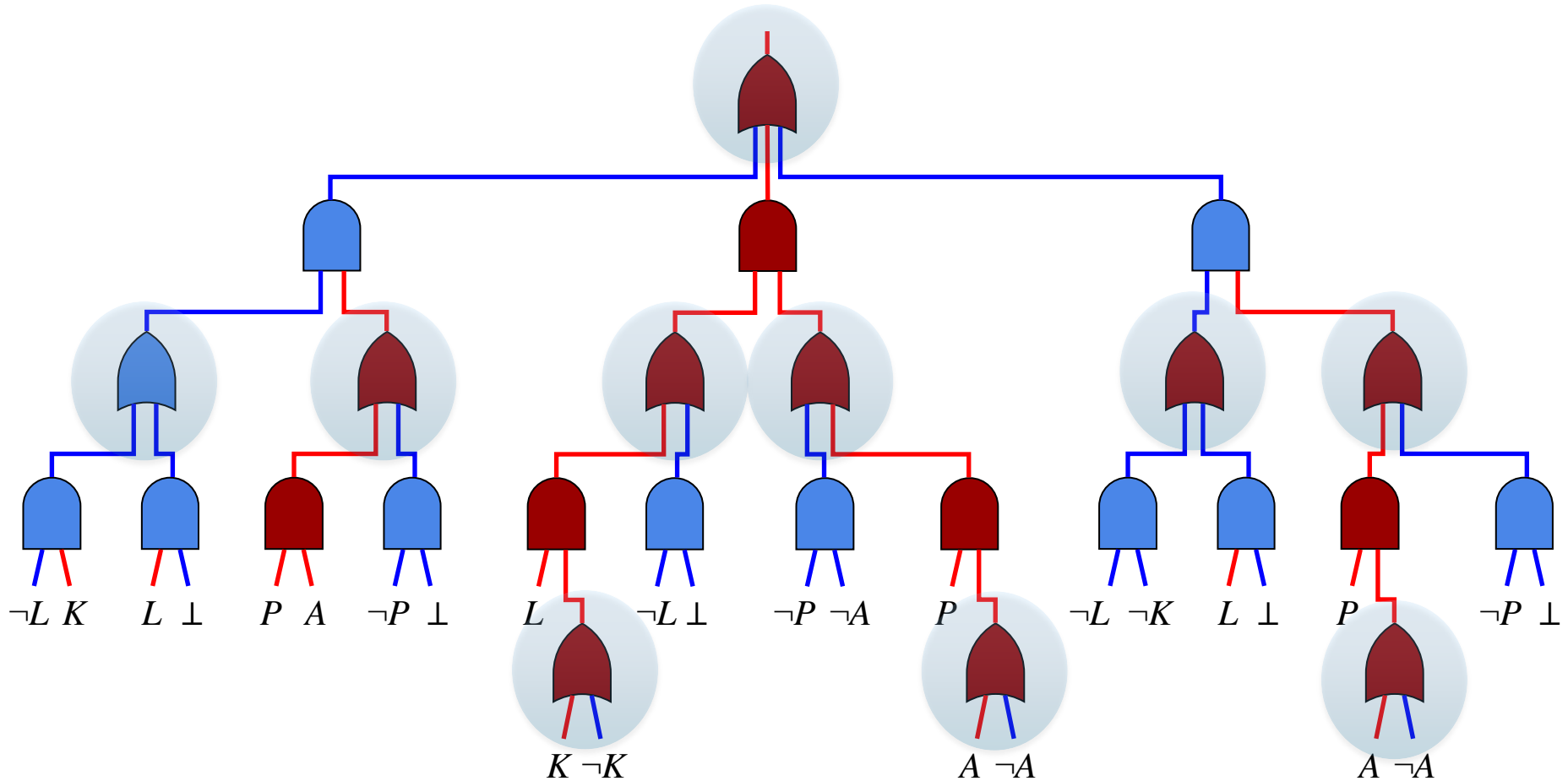
# Recall: Decomposability



AND gates have disjoint input circuits

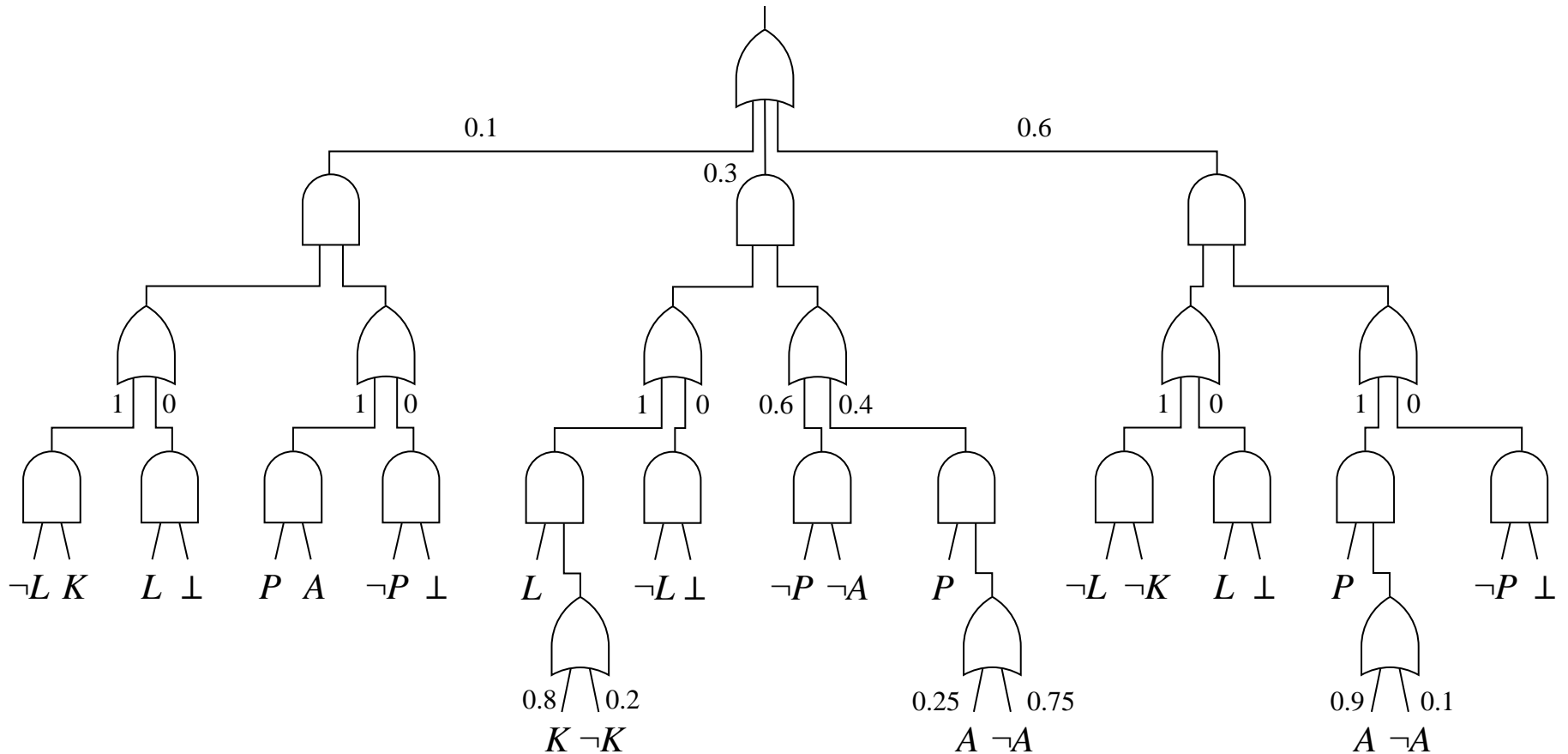


# Recall: Determinism



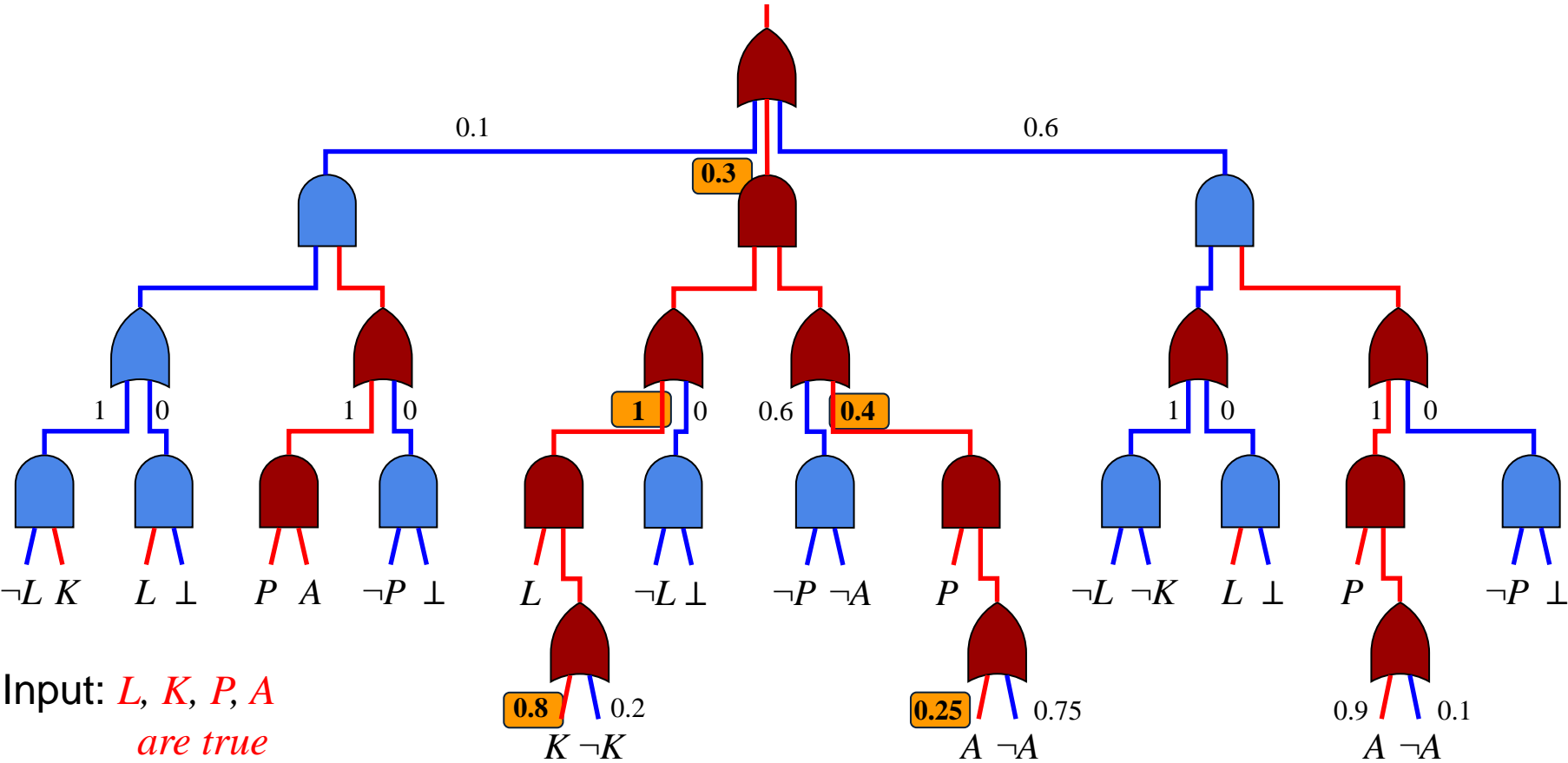
Input:  $L$ ,  $K$ ,  $P$ ,  $A$  are **true** and  $\neg L$ ,  $\neg K$ ,  $\neg P$ ,  $\neg A$  are **false**  
Property: OR gates have at most one **true** input wire

# PSDD: Probabilistic SDD



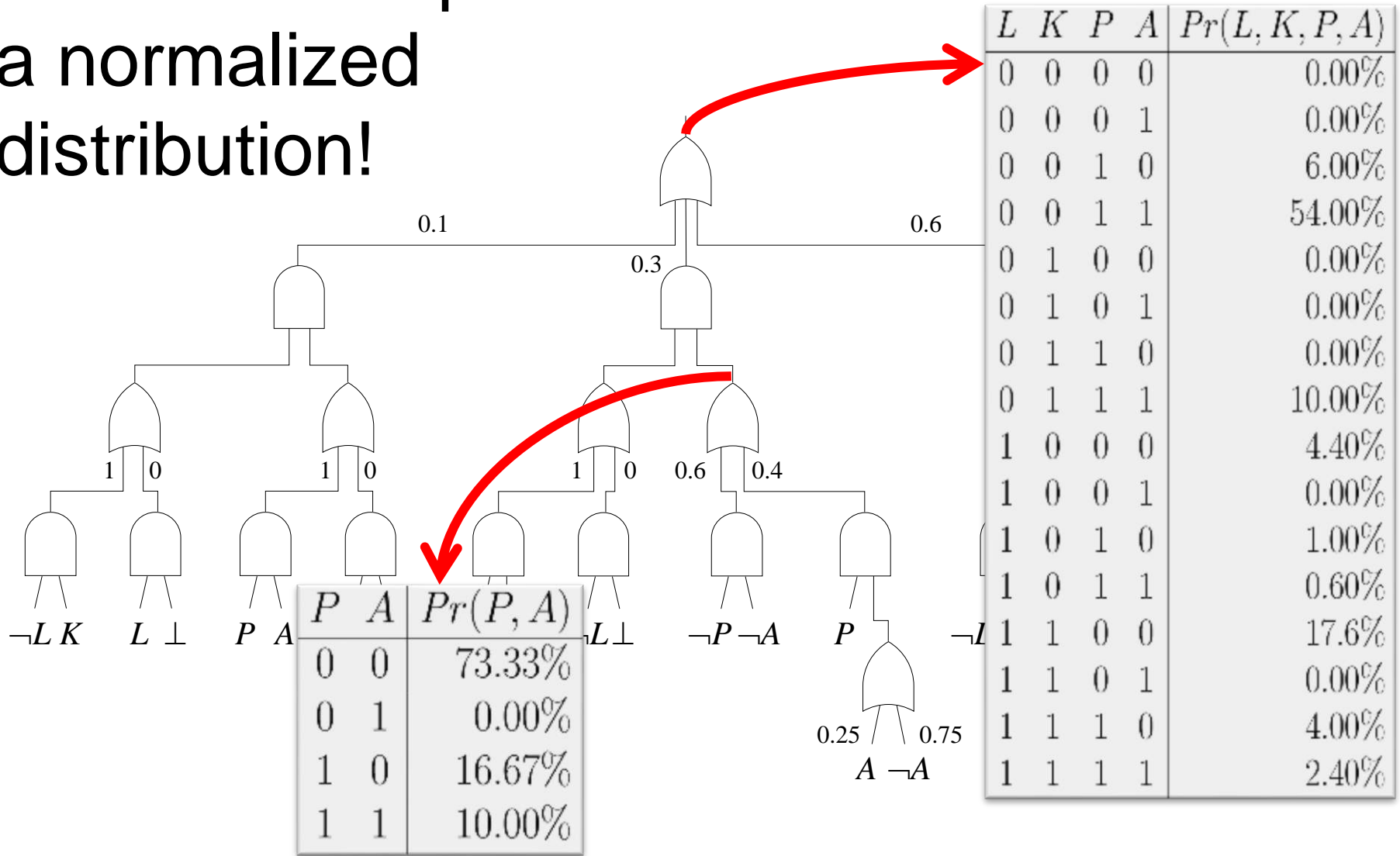
Syntax: assign a normalized probability to each OR gate input

# PSDD: Probabilistic SDD



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

Each node represents  
a normalized  
distribution!

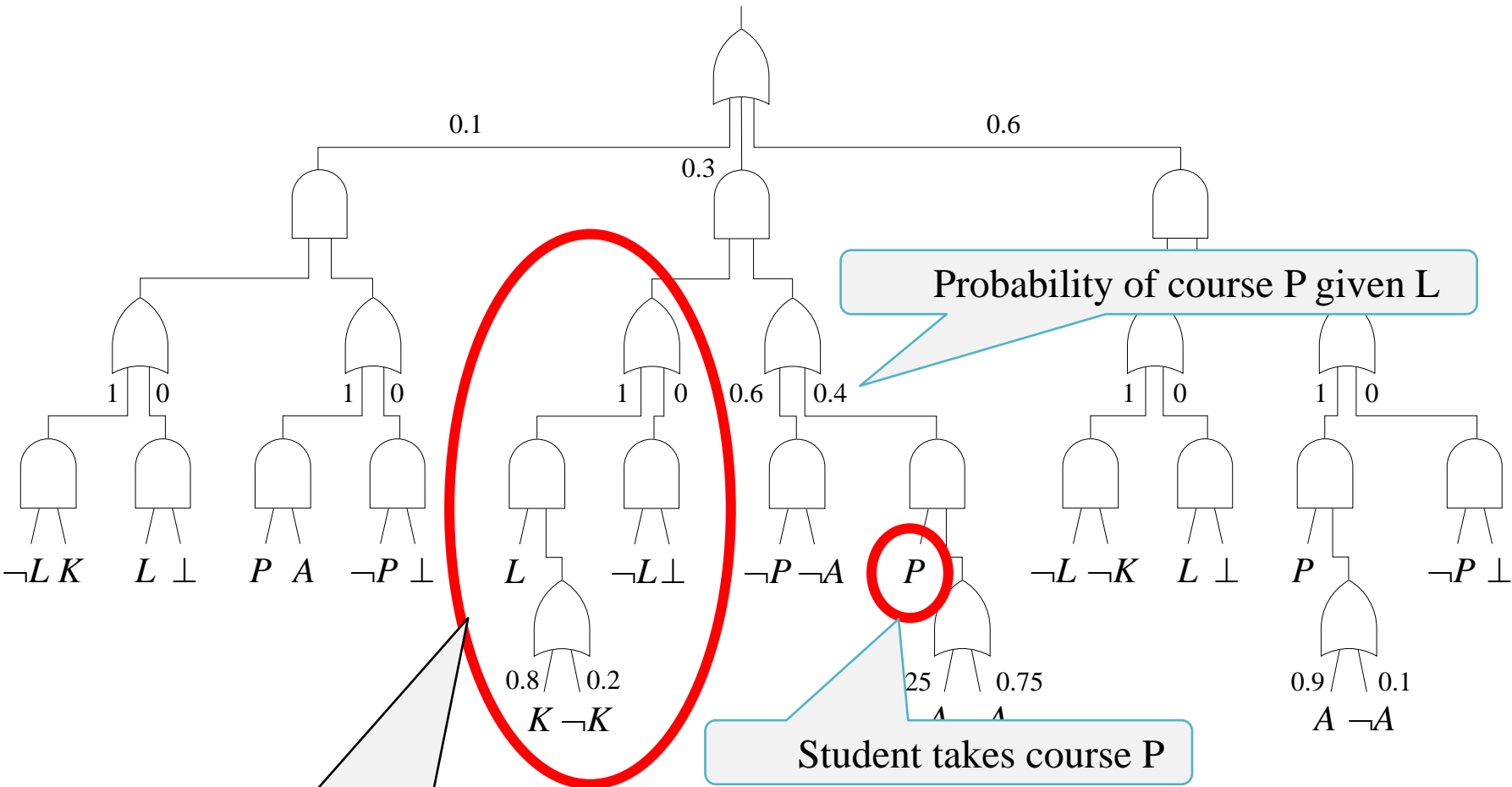


Can read probabilistic independences off the circuit structure

# 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)

# Parameters are Interpretable



Student takes course L

Probability of course P given L

Student takes course P

Explainable AI DARPA Program

***Learning  
Probabilistic Circuit  
Parameters***

# 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! 😊

- Where does the structure come from?  
For now: simply compiled from constraint...



# Combinatorial Objects: Rankings

rank	sushi
1	fatty tuna
2	sea urchin
3	salmon roe
4	shrimp
5	tuna
6	squid
7	tuna roll
8	see eel
9	egg
10	cucumber roll

rank	sushi
1	shrimp
2	sea urchin
3	salmon roe
4	fatty tuna
5	tuna
6	squid
7	tuna roll
8	see eel
9	egg
10	cucumber roll

**10 items:**  
3,628,800  
rankings

**20 items:**  
2,432,902,008,176,640,000  
rankings

# Combinatorial Objects: Rankings

rank	sushi
1	fatty tuna
2	sea urchin
3	salmon roe
4	shrimp
5	tuna
6	squid
7	tuna roll
8	sea eel
9	egg
10	cucumber roll

- Predict Boolean Variables:  
 $A_{ij}$  - item  $i$  at position  $j$
- Constraints:

each item  $i$  assigned to  
a unique position ( $n$  constraints)

$$\bigvee_j A_{ij} \wedge \left( \bigwedge_{k \neq j} \neg A_{ik} \right)$$

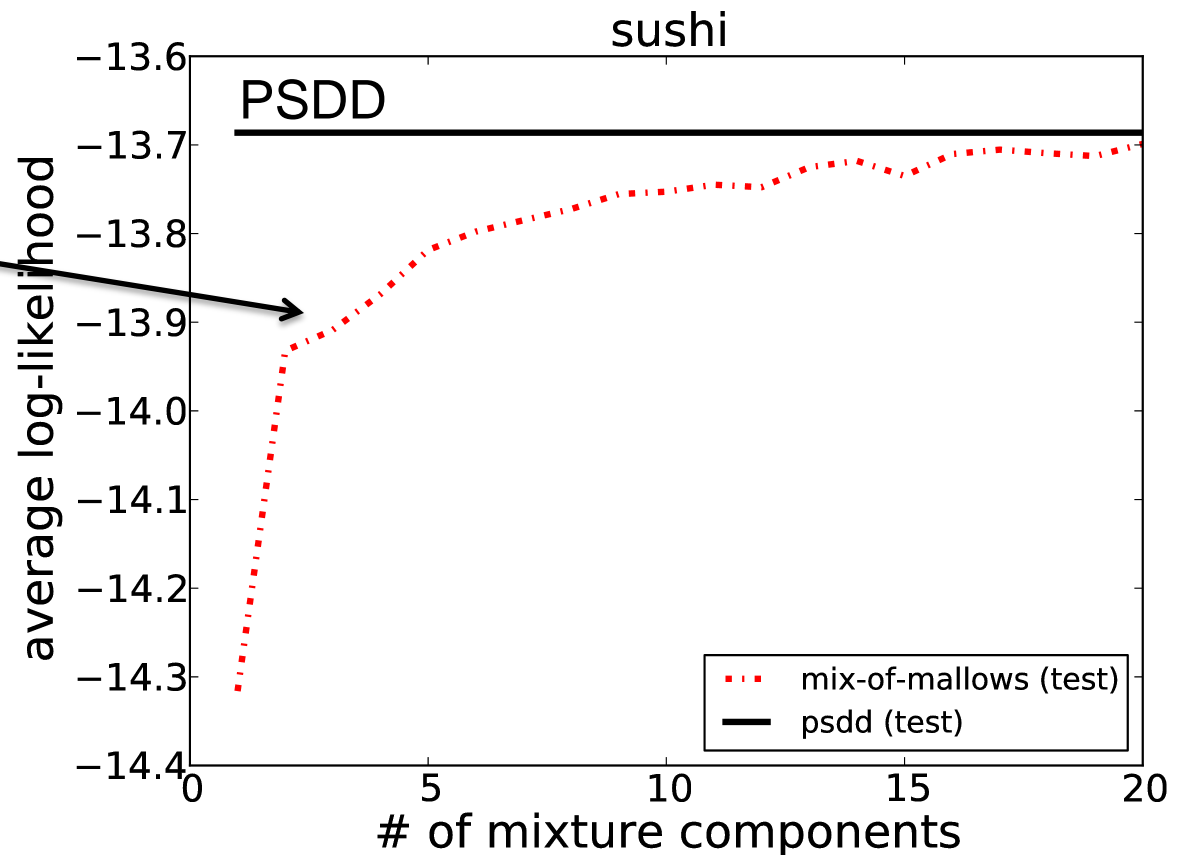
each position  $j$  assigned  
a unique item ( $n$  constraints)

$$\bigvee_i A_{ij} \wedge \left( \bigwedge_{k \neq i} \neg A_{kj} \right)$$

# Learning Preference Distributions

Special-purpose  
distribution:  
Mixture-of-Mallows

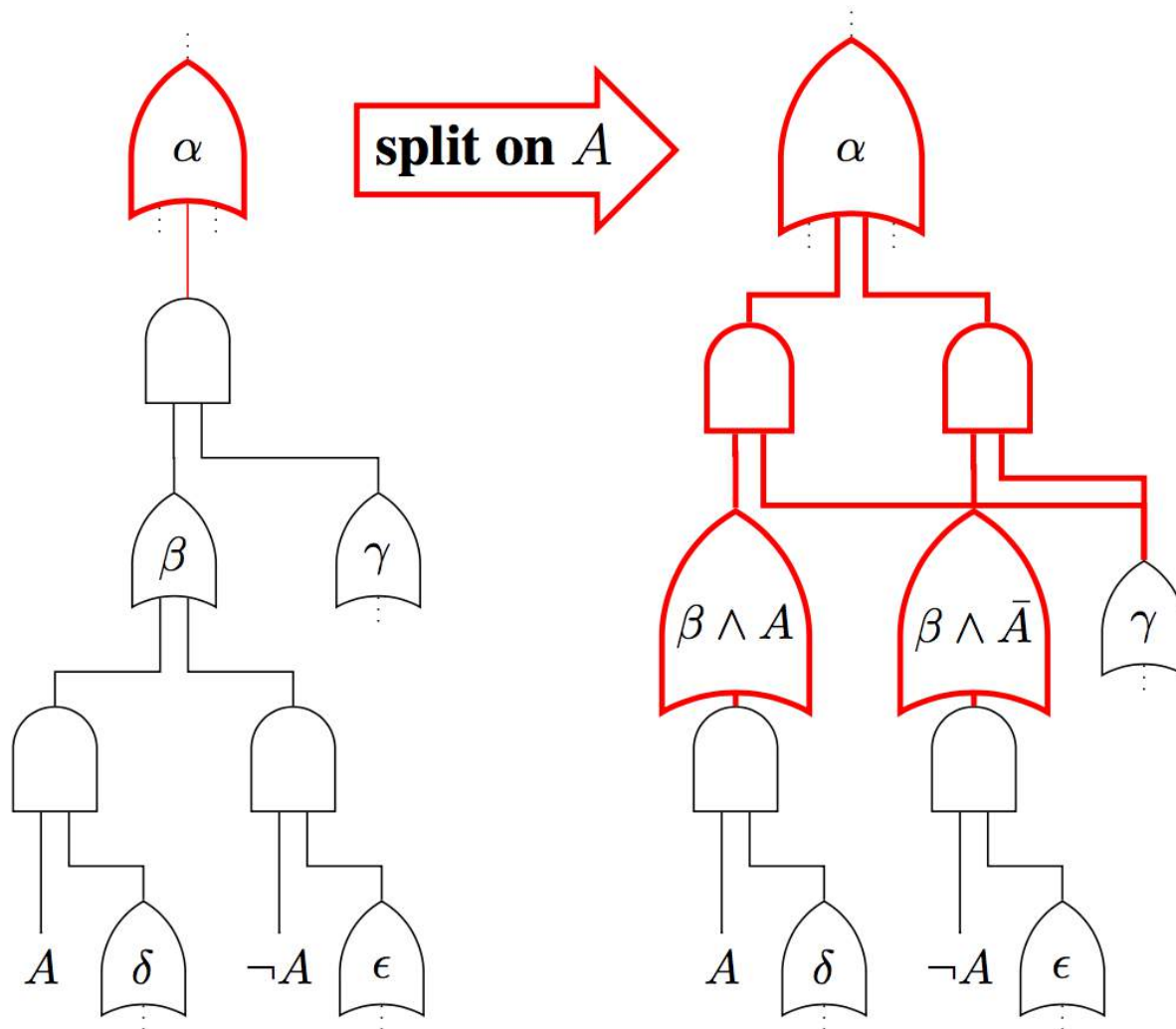
- # of components from 1 to 20
- EM with 10 random seeds
- Implementation of Lu & Boutilier



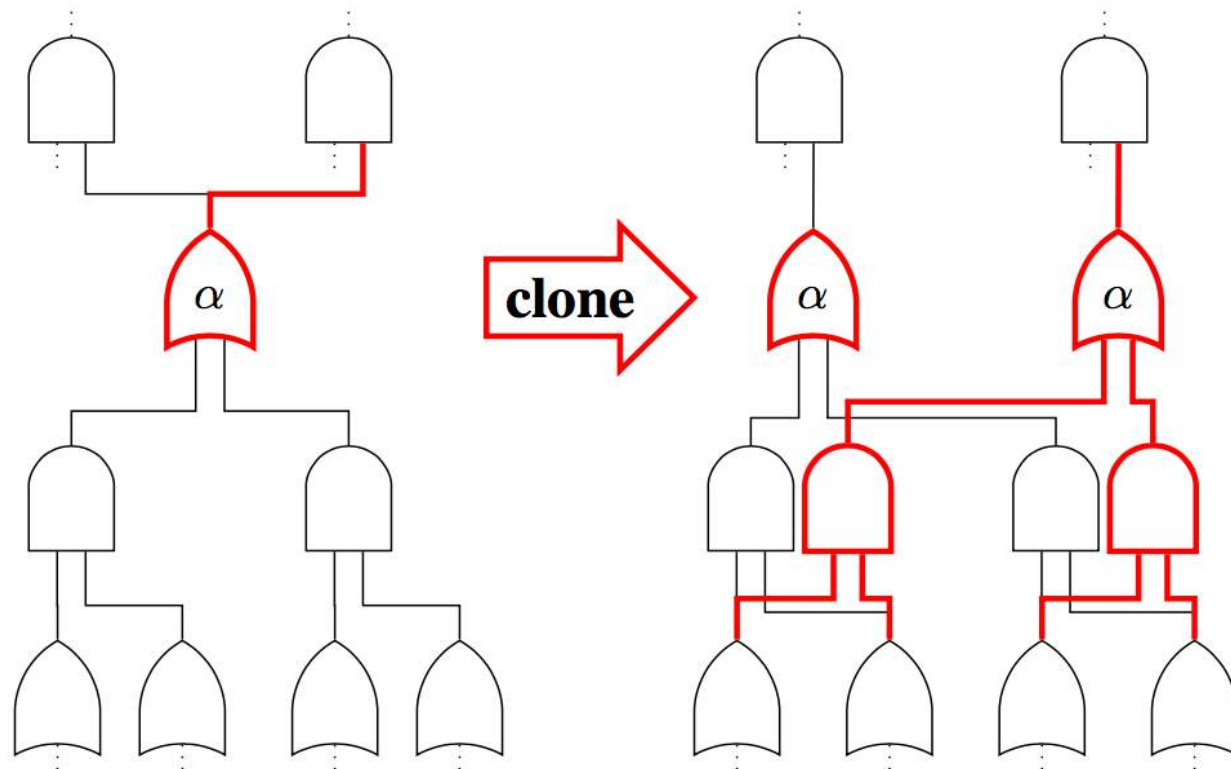
Circuit structure does not even depend on data!

***Learning  
Probabilistic Circuit  
Structure***

# Structure Learning Primitive

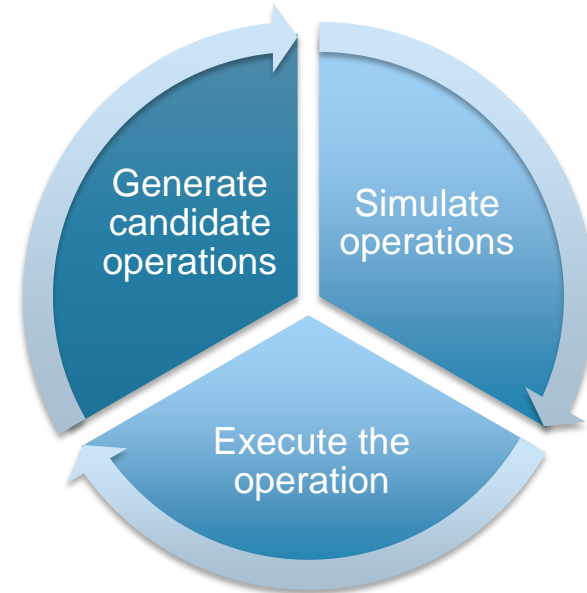
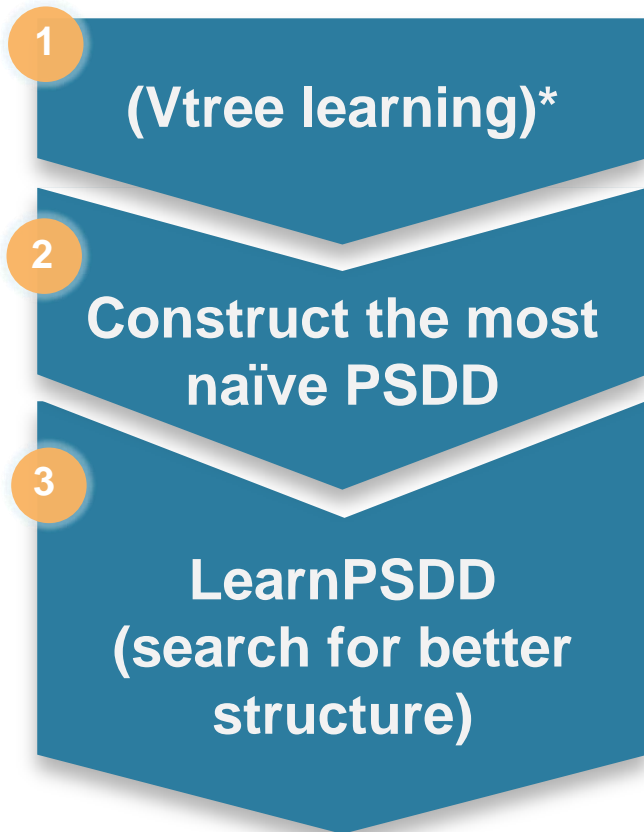


# Structure Learning Primitive



Primitives maintain PSDD properties  
and constraint of root!

# LearnPSDD Algorithm

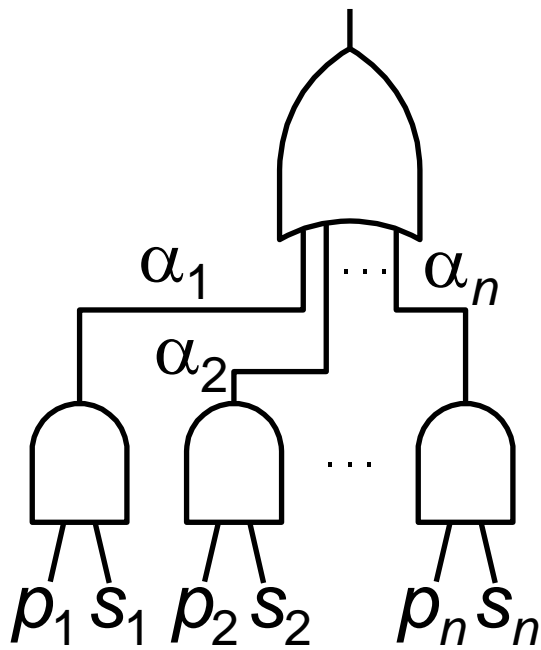


$$\text{score} = \frac{\ln \mathcal{L}(r' | \mathcal{D}) - \ln \mathcal{L}(r | \mathcal{D})}{\text{size}(r') - \text{size}(r)}$$

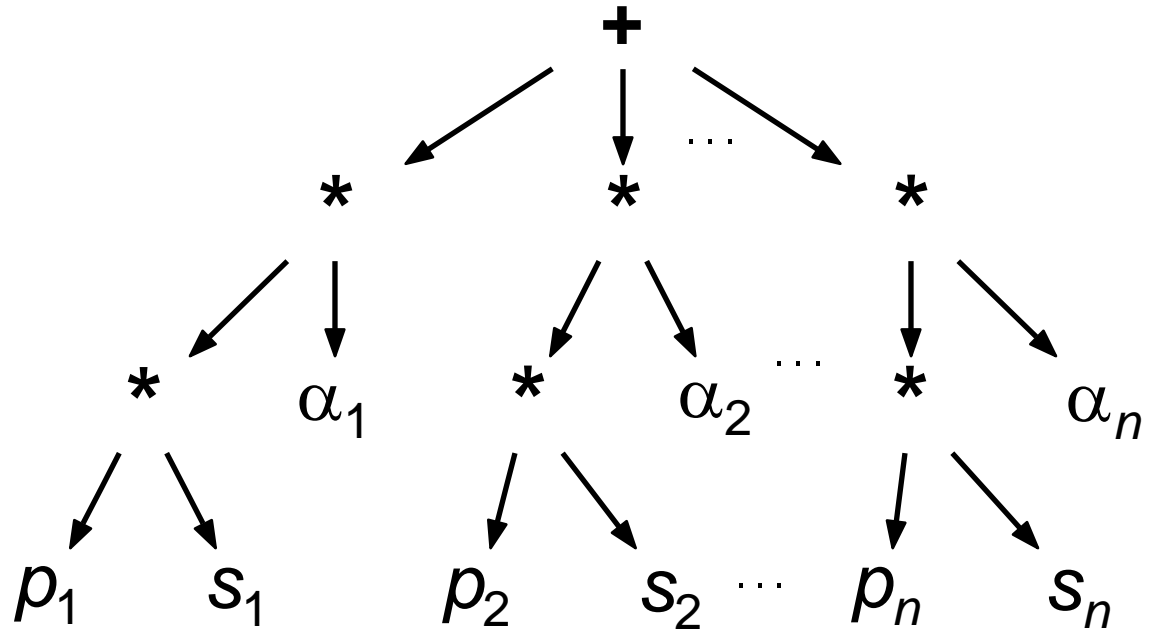
Works with or without logical constraint.

# PSDDs

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



**PSDD**



**AC**



# Experiments on 20 datasets

Datasets	Var	Train	Valid	Test	LearnPSDD		EM-LearnPSDD		SearchSPN	Merged L-SPN		Merged O-SPN	
					LL	Size	LL	Size	LL	LL	Size	LL	Size
NLTCS	16	16181	2157	3236	-6.03 <sup>†*</sup>	3170	-6.03 <sup>*</sup>	2147	-6.07	-6.04	3988	-6.05	1152
MSNBC	17	291326	38843	58265	-6.05 <sup>†</sup>	8977	-6.04 <sup>*</sup>	3891	-6.06	-6.46	2440	-6.08	9478
KDD	64	1800992	19907	34955	-2.16 <sup>†</sup>	14974	-2.12 <sup>*</sup>	9182	-2.16	-2.14	6670	-2.19	16608
Plants	69	17412	2321	3482	-14.93	13129	-13.79 <sup>*</sup>	13951	-13.12 <sup>†</sup>	-12.69	47802	-13.49	36960
Audio	100	15000	2000	3000	-42.53	13765	-41.98 <sup>*</sup>	9721	-40.13 <sup>†</sup>	-40.02	10804	-42.06	6142
Jester	100	9000	1000	4116	-57.67	11322	-53.47 <sup>*</sup>	7014	-53.08 <sup>†</sup>	-52.97	10002	-55.36	4996
Netflix	100	15000	2000	3000	-58.92	10997	-58.41 <sup>*</sup>	6250	-56.91 <sup>†</sup>	-56.64	11604	-58.64	6142
Accidents	111	12758	1700	2551	-34.13	10489	-33.64 <sup>*</sup>	6752	-30.02 <sup>†</sup>	-30.01	13322	-30.83	6846
Retail	135	22041	2938	4408	-11.13	4091	-10.81 <sup>*</sup>	7251	-10.97 <sup>†</sup>	-10.87	2162	-10.95	3158
Pumsb-Star	163	12262	1635	2452	-34.11	10489	-33.67 <sup>*</sup>	7965	-28.69 <sup>†</sup>	-24.11	17604	-24.34	18338
DNA	180	1600	400	1186	-89.11 <sup>*</sup>	6068	-92.67	14864	-81.76 <sup>†</sup>	-85.51	4320	-87.49	1430
Kosarek	190	33375	4450	6675	-10.99 <sup>†</sup>	11034	-10.81 <sup>*</sup>	10179	-11.00	-10.62	5318	-10.98	6712
MSWeb	294	29441	32750	5000	-10.18 <sup>†</sup>	11389	-9.97 <sup>*</sup>	14512	-10.25	-9.90	16484	-10.06	12770
Book	500	8700	1159	1739	-35.90	15197	-34.97 <sup>*</sup>	11292	-34.91 <sup>†</sup>	-34.76	11998	-37.44	11916
EachMovie	500	4524	1002	591	-56.43 <sup>*</sup>	12483	-58.01	16074	-53.28 <sup>†</sup>	-52.07	15998	-58.05	19846
WebKB	839	2803	558	838	-163.42	10033	-161.09 <sup>*</sup>	18431	-157.88 <sup>†</sup>	-153.55	20134	-161.17	10046
Reuters-52	889	6532	1028	1530	-94.94	10585	-89.61 <sup>*</sup>	9546	-86.38 <sup>†</sup>	-83.90	46232	-87.49	28334
20NewsGrp.	910	11293	3764	3764	-161.41	12222	-161.09 <sup>*</sup>	18431	-153.63 <sup>†</sup>	-154.67	43684	-161.46	29016
BBC	1058	1670	225	330	-260.83	10585	-253.19 <sup>*</sup>	20327	-252.13 <sup>†</sup>	-253.45	21160	-260.59	8454
AD	1556	2461	327	491	-30.49 <sup>*</sup>	9666	-31.78	9521	-16.97 <sup>†</sup>	-16.77	49790	-15.39	31070

Compared to SPN learners, LearnPSDD gives comparable performance yet smaller size

# Learn Mixtures of PSDDs

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

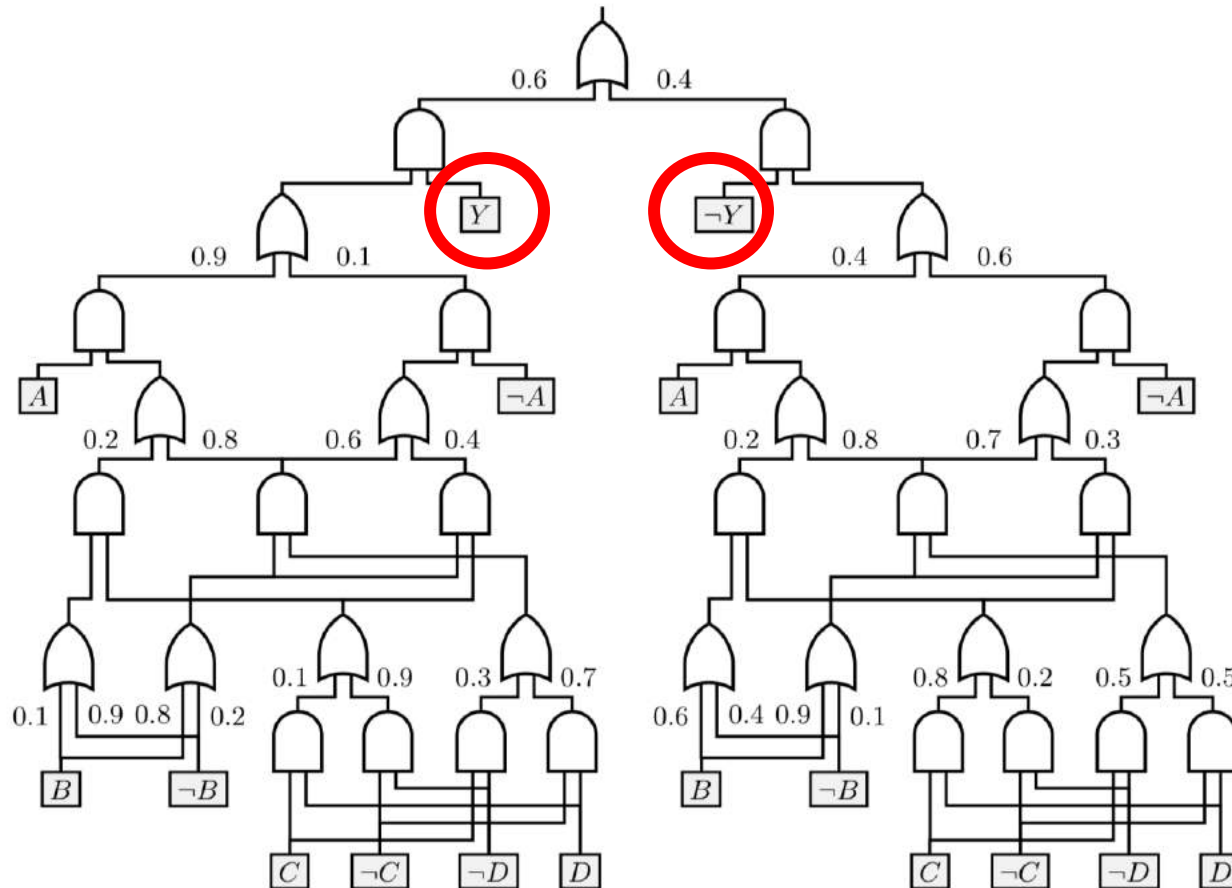
Competitive with

- SPN learners
- Cutset network learners

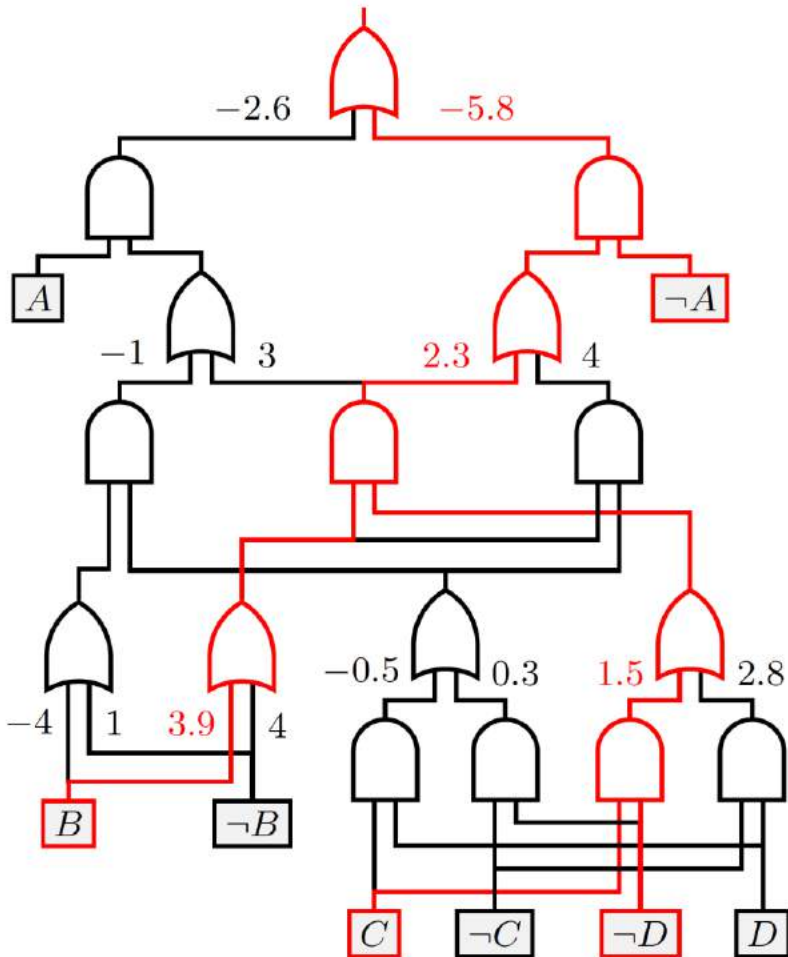
# ***Logistic Circuits***

# What if I only want to classify Y?

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



# Logistic Circuits



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

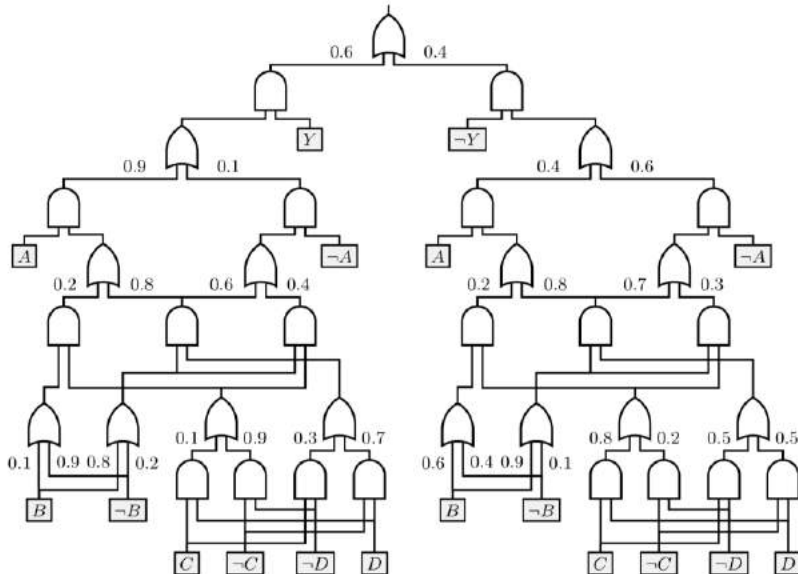
- Take all 'hot' wires
- Sum their weights
- Push through logistic function

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

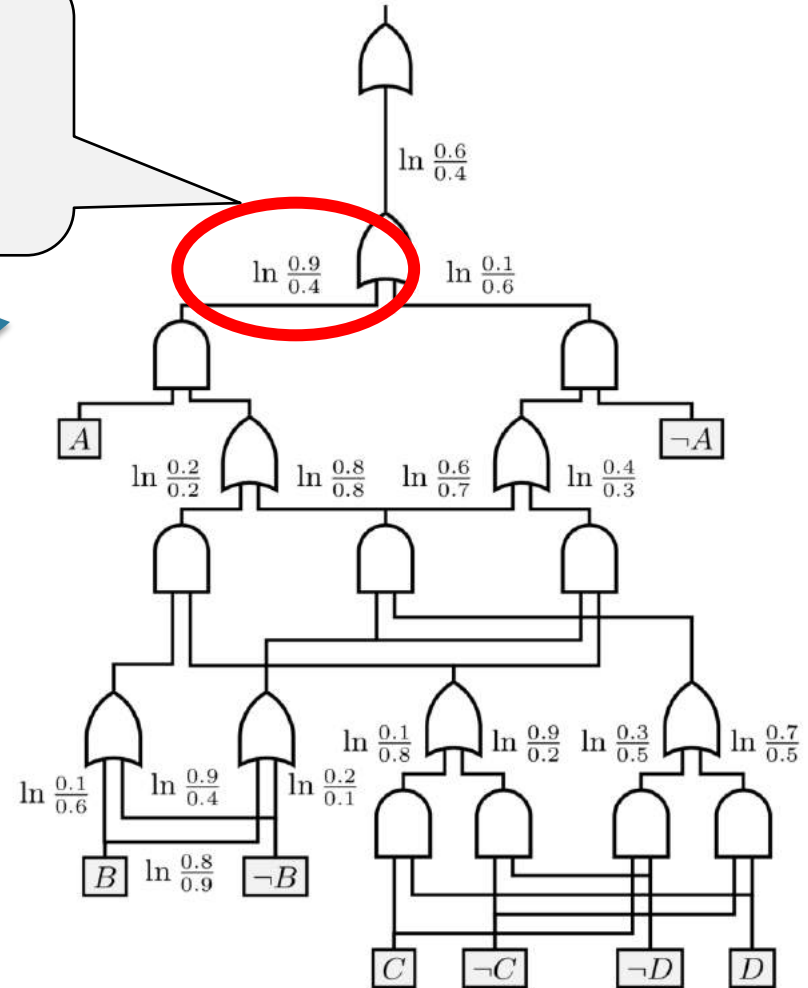
# Logistic vs. Probabilistic Circuits

Probabilities become log-odds

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



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



# Parameter Learning

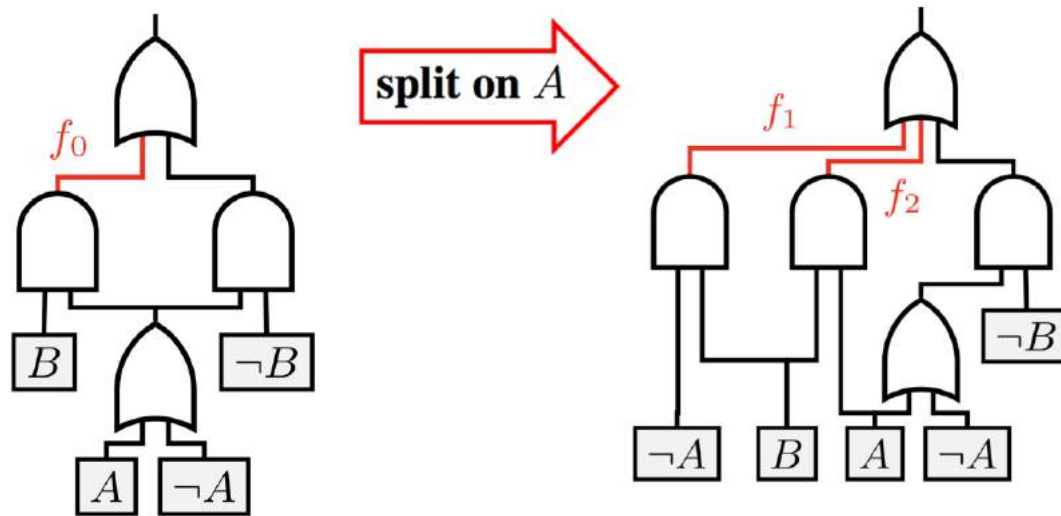
Reduce to logistic regression:

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

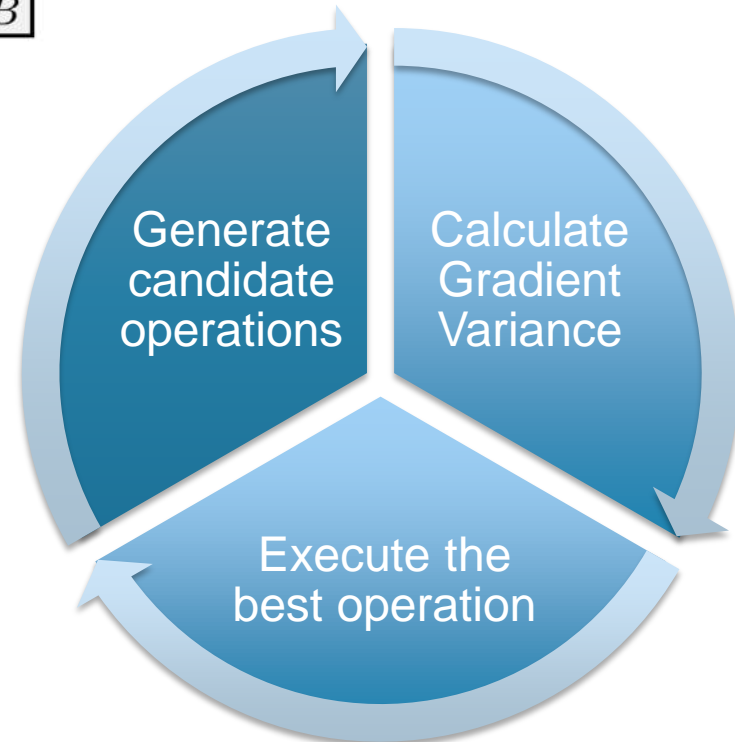
Features associated with each wire  
“Global Circuit Flow” features

Learning parameters  $\theta$  is convex optimization!

# Logistic Circuit Structure Learning



Similar to LearnPSDD  
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

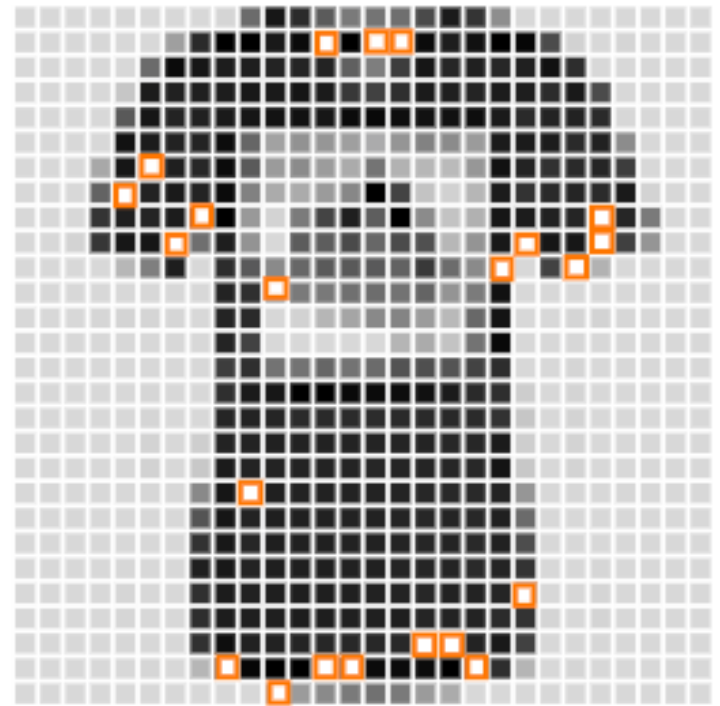
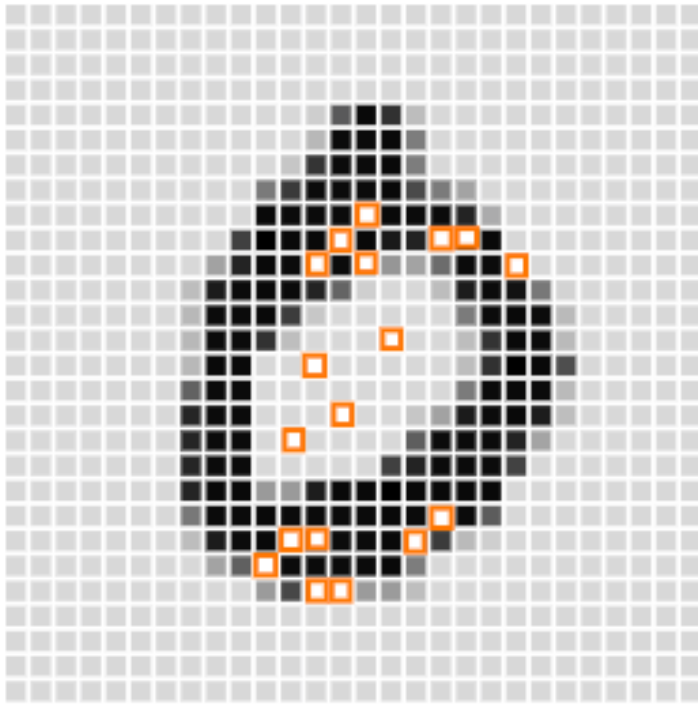
# Significantly Smaller in Size

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

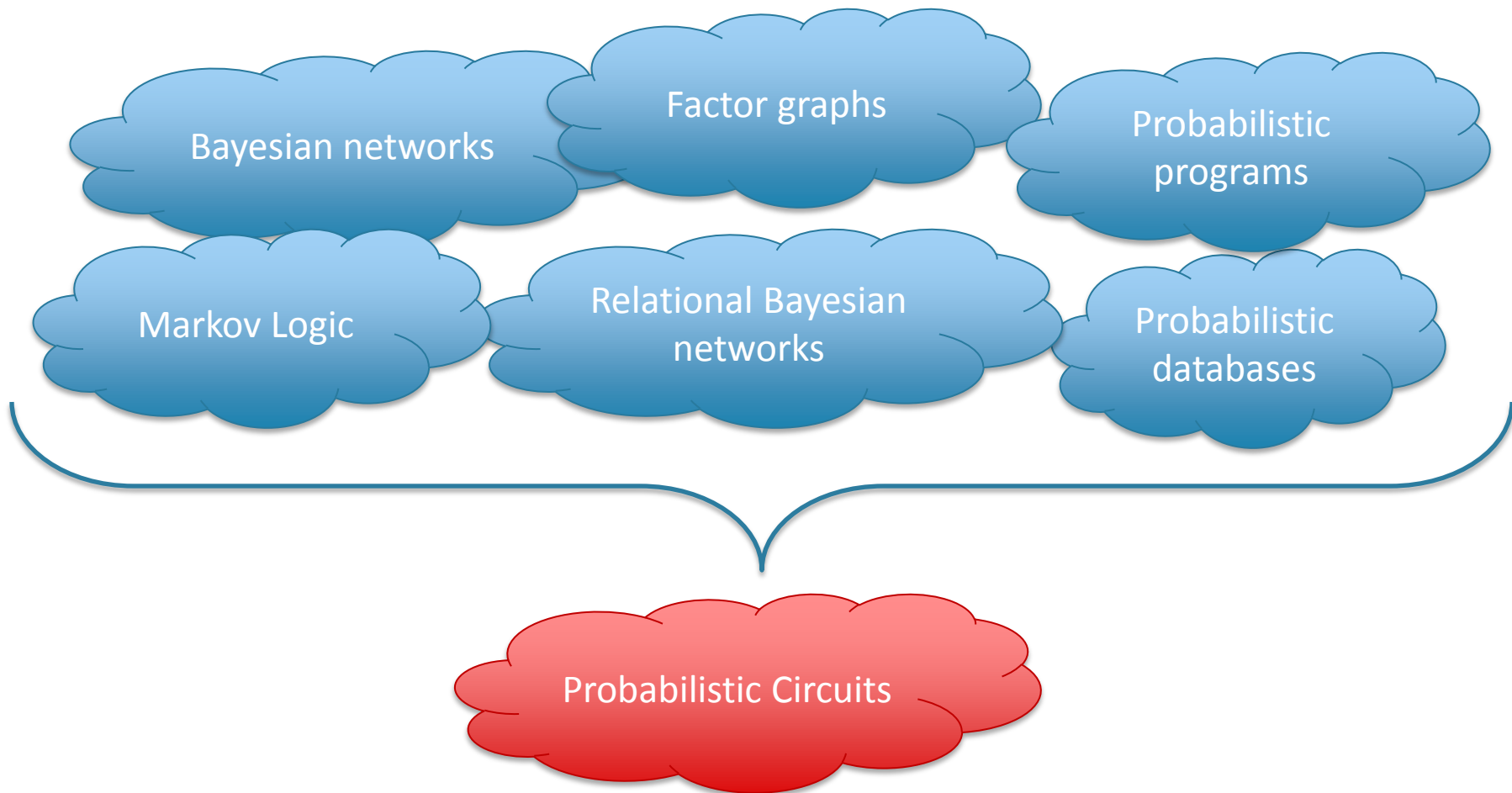
ACCURACY % WITH % OF TRAINING DATA	MNIST			FASHION		
	100%	10%	2%	100%	10%	2%
5-LAYER MLP	99.3	<b>98.2</b>	94.3	89.8	86.5	80.9
CNN WITH 3 CONV LAYERS	99.1	98.1	95.3	90.7	87.6	83.8
LOGISTIC CIRCUIT (BINARY)	97.4	96.9	94.1	87.6	86.7	83.2
LOGISTIC CIRCUIT (REAL-VALUED)	<b>99.4</b>	97.6	<b>96.1</b>	<b>91.3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

# Interpretable?

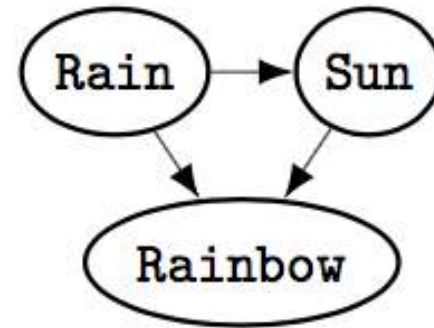
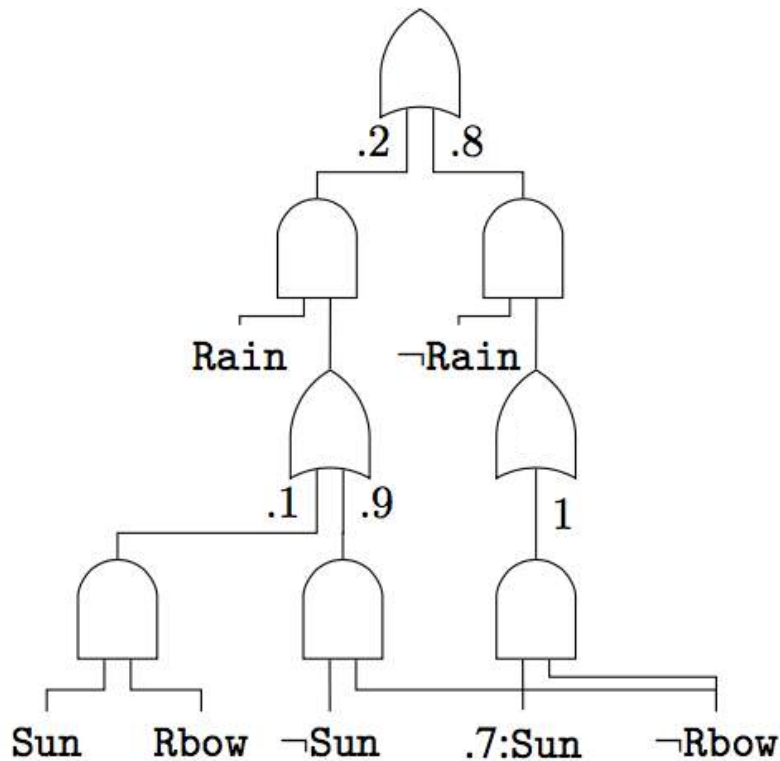


# ***Reasoning with Probabilistic Circuits***

# Compilation target for probabilistic reasoning



# Compilation for Prob. Inference



$$\Pr(\text{Rain}) = 0.2,$$

$$\Pr(\text{Sun} \mid \text{Rain}) = \begin{cases} 0.1 & \text{if Rain} \\ 0.7 & \text{if } \neg\text{Rain} \end{cases}$$

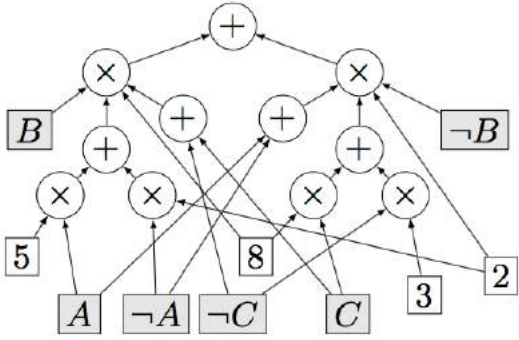
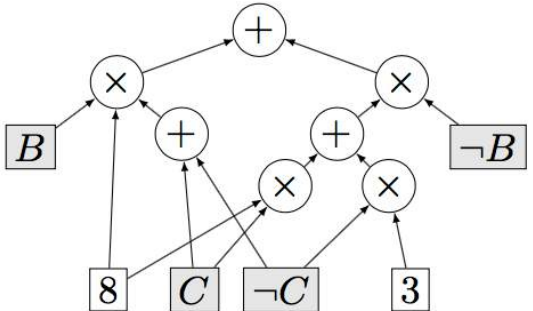
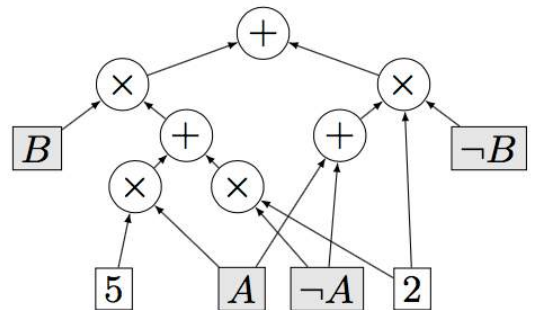
$$\Pr(\text{Rbow} \mid \text{R}, \text{S}) = \begin{cases} 1 & \text{if Rain} \wedge \text{Sun} \\ 0 & \text{otherwise} \end{cases}$$

# Collapsed Compilation

To sample a circuit:

1. Compile bottom up until you reach the size limit
2. Pick a variable you want to sample
3. Sample it according to its marginal distribution in the current circuit
4. Condition on the sampled value
5. (Repeat)

Asymptotically unbiased importance sampler 😊



- 
- 
- 



Circuits +  
importance weights  
approximate any query



# Experiments

Table 2: Hellinger distances across methods with internal treewidth and size bounds

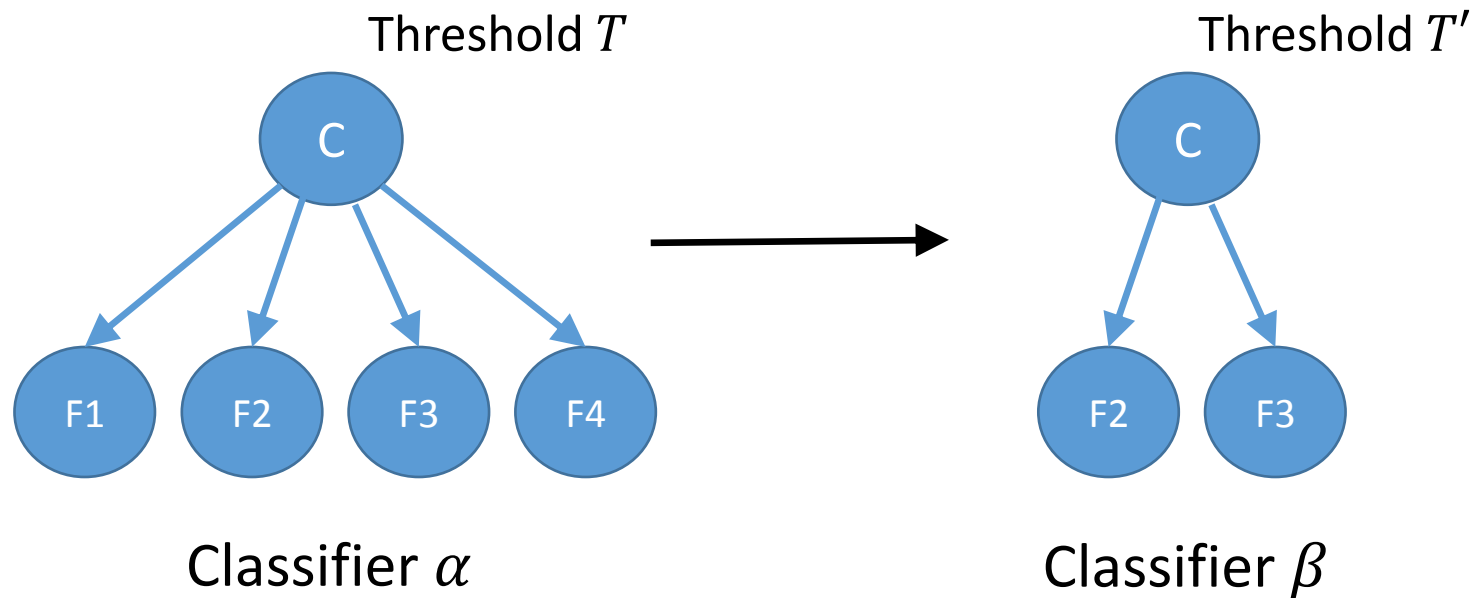
Method	50-20	75-26	DBN	Grids	Segment	linkage	frust
EDBP-100k	$2.19e-3$	$3.17e-5$	$6.39e-1$	$1.24e-3$	$1.63e-6$	$6.54e-8$	$4.73e-3$
EDBP-1m	$7.40e-7$	$2.21e-4$	$6.39e-1$	$1.98e-7$	$1.93e-7$	$5.98e-8$	$4.73e-3$
SS-10	$2.51e-2$	$2.22e-3$	$6.37e-1$	$3.10e-1$	$3.11e-7$	$4.93e-2$	$1.05e-2$
SS-12	$6.96e-3$	$1.02e-3$	$6.27e-1$	$2.48e-1$	$3.11e-7$	$1.10e-3$	$5.27e-4$
SS-15	$9.09e-6$	$1.09e-4$	(Exact)	$8.74e-4$	$3.11e-7$	$4.06e-6$	$6.23e-3$
FD	$9.77e-6$	$1.87e-3$	$1.24e-1$	$1.98e-4$	$6.00e-8$	$5.99e-6$	$5.96e-6$
MinEnt	$1.50e-5$	$3.29e-2$	$1.83e-2$	$3.61e-3$	$3.40e-7$	$6.16e-5$	$3.10e-2$
RBVar	$2.66e-2$	$4.39e-1$	$6.27e-3$	$1.20e-1$	$3.01e-7$	$2.02e-2$	$2.30e-3$

Competitive with state-of-the-art  
approximate inference in graphical models.  
Outperforms it on several benchmarks!

# ***Reasoning About Classifiers***

# Classifier Trimming

$$C_T(\text{features}) = \mathbb{I}(\Pr(C \mid \text{features}) \geq T)$$



Trim features while maintaining  
classification behavior

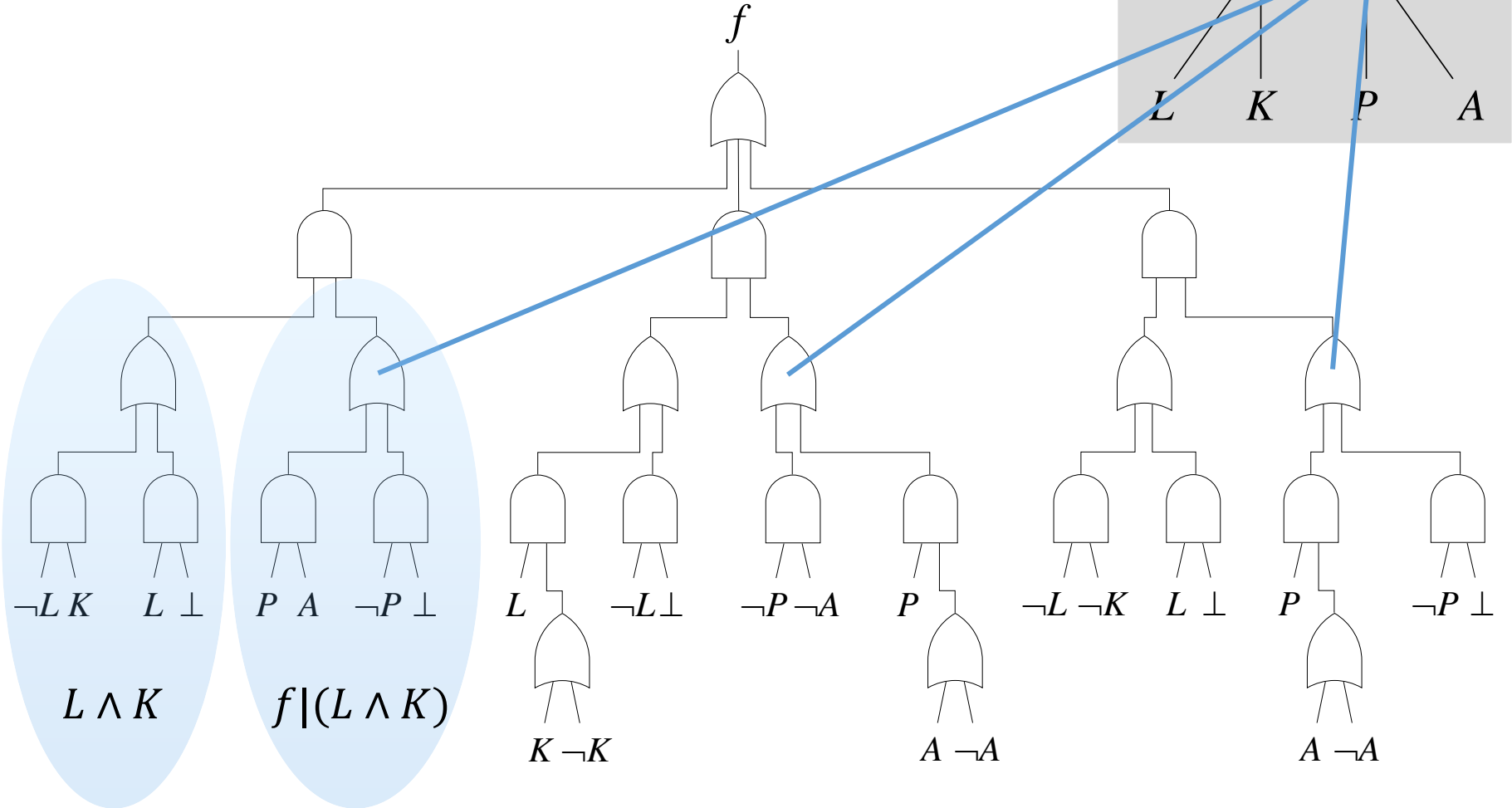
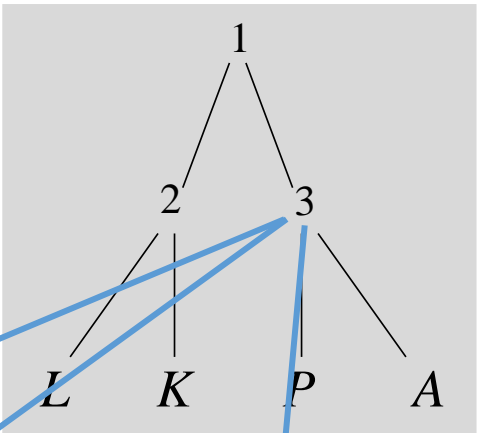
# How to measure Similarity?

“Expected Classification Agreement”

$$\text{ECA}(\alpha, \beta) = \sum_{\mathbf{f}} \mathbb{I}(C_T(\mathbf{f}) = C_{T'}(\mathbf{f}')) \cdot \text{Pr}(\mathbf{f})$$

What is the expected probability that a classifier  $\alpha$  will agree with its trimming  $\beta$ ?

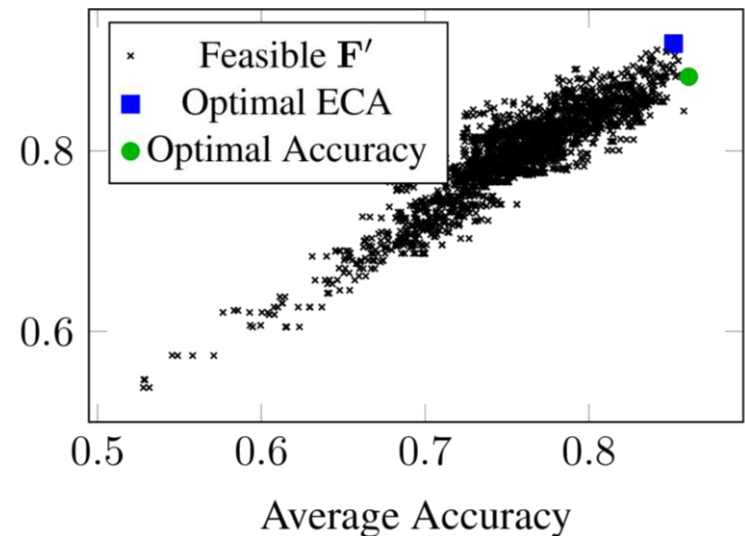
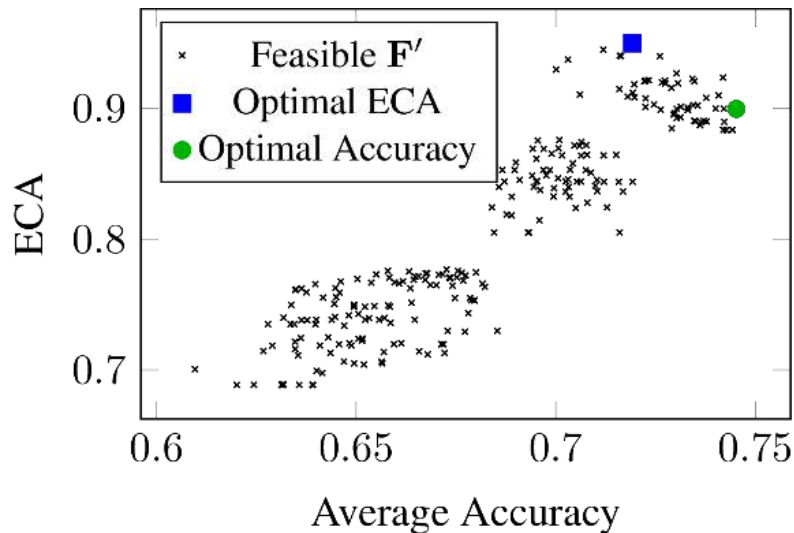
# Solving $PP^{PP}$ problems with constrained SDDs



# SDD method faster than traditional jointree inference

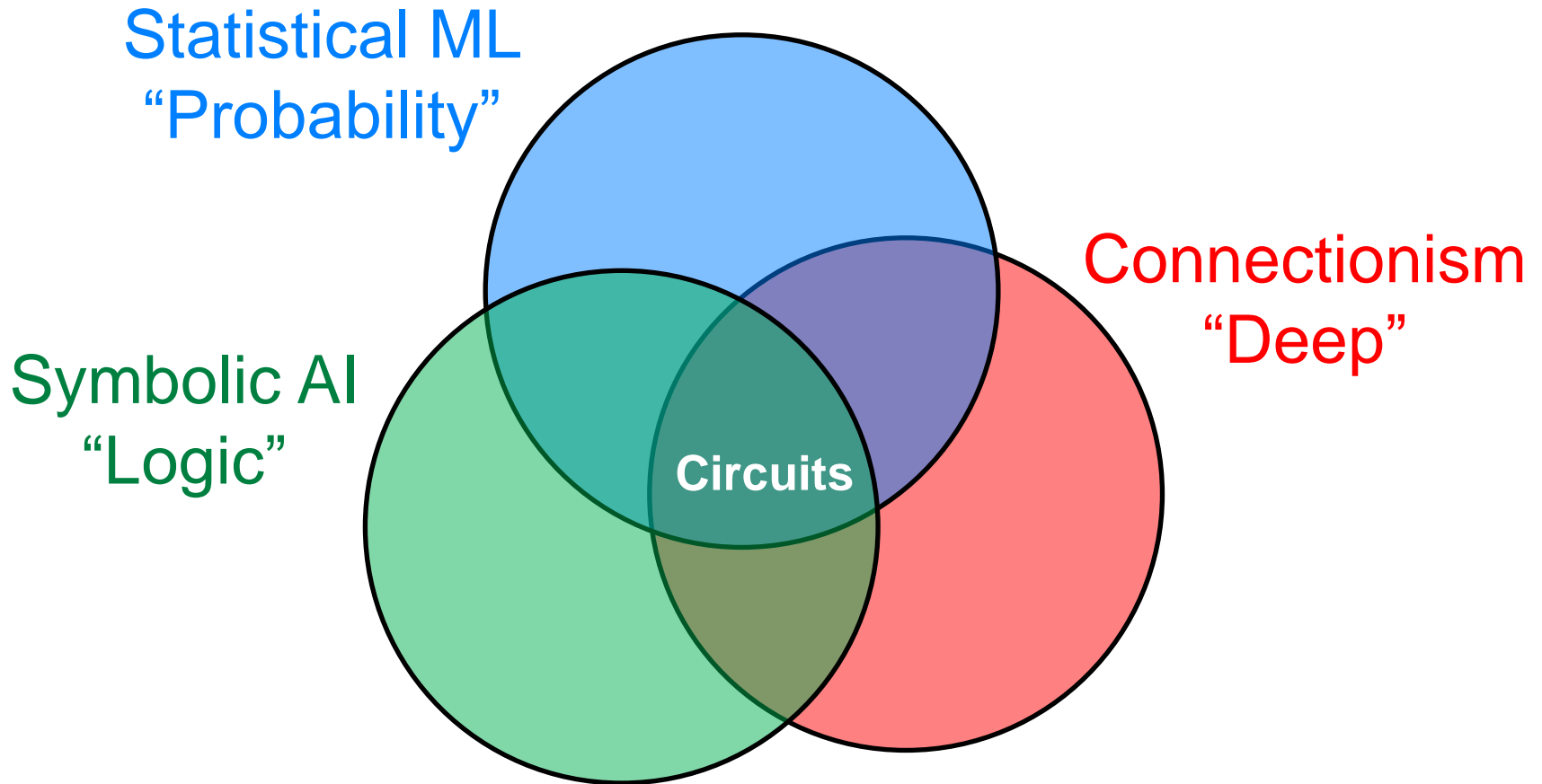
Network	# nodes	naive	FS-SDD
alarm	37	143.920	19.061
win95pts	76	23.581	14.732
tcc4e	98	48.508	2.384
emdec6g	168	28.072	3.688
diagnose	203	105.660	6.667

# Classification agreement and accuracy



Higher agreement tends to get higher accuracy  
Additional dimension for feature selection

# Conclusions





# Questions?



*PSDD with 15,000 nodes*

# References

- Doga Kisa, Guy Van den Broeck, Arthur Choi and Adnan Darwiche. [Probabilistic sentential decision diagrams](#), *In Proceedings of the 14th International Conference on Principles of Knowledge Representation and Reasoning (KR)*, 2014.
- Arthur Choi, Guy Van den Broeck and Adnan Darwiche. [Tractable Learning for Structured Probability Spaces: A Case Study in Learning Preference Distributions](#), *In Proceedings of 24th International Joint Conference on Artificial Intelligence (IJCAI)*, 2015.
- Arthur Choi, Guy Van den Broeck and Adnan Darwiche. [Probability Distributions over Structured Spaces](#), *In Proceedings of the AAAI Spring Symposium on KRR*, 2015.
- Jessa Bekker, Jesse Davis, Arthur Choi, Adnan Darwiche and Guy Van den Broeck. [Tractable Learning for Complex Probability Queries](#), *In Advances in Neural Information Processing Systems 28 (NIPS)*, 2015.
- YooJung Choi, Adnan Darwiche and Guy Van den Broeck. [Optimal Feature Selection for Decision Robustness in Bayesian Networks](#), *In Proceedings of the 26th International Joint Conference on Artificial Intelligence (IJCAI)*, 2017.

# References

- Yitao Liang, Jessa Bekker and Guy Van den Broeck. [Learning the Structure of Probabilistic Sentential Decision Diagrams](#), *In Proceedings of the 33rd Conference on Uncertainty in Artificial Intelligence (UAI)*, 2017.
- Yitao Liang and Guy Van den Broeck. [Towards Compact Interpretable Models: Shrinking of Learned Probabilistic Sentential Decision Diagrams](#), *In IJCAI 2017 Workshop on Explainable Artificial Intelligence (XAI)*, 2017.
- YooJung Choi and Guy Van den Broeck. [On Robust Trimming of Bayesian Network Classifiers](#), *In Proceedings of the 27th International Joint Conference on Artificial Intelligence (IJCAI)*, 2018.
- Jingyi Xu, Zilu Zhang, Tal Friedman, Yitao Liang and Guy Van den Broeck. [A Semantic Loss Function for Deep Learning with Symbolic Knowledge](#), *In Proceedings of the 35th International Conference on Machine Learning (ICML)*, 2018.
- Tal Friedman and Guy Van den Broeck. [Approximate Knowledge Compilation by Online Collapsed Importance Sampling](#), *In Advances in Neural Information Processing Systems 31 (NIPS)*, 2018.
- Yitao Liang and Guy Van den Broeck. [Learning Logistic Circuits](#), *In Proceedings of the 33rd Conference on Artificial Intelligence (AAAI)*, 2019.