

# Tractable Probabilistic Circuits

Guy Van den Broeck

# Outline

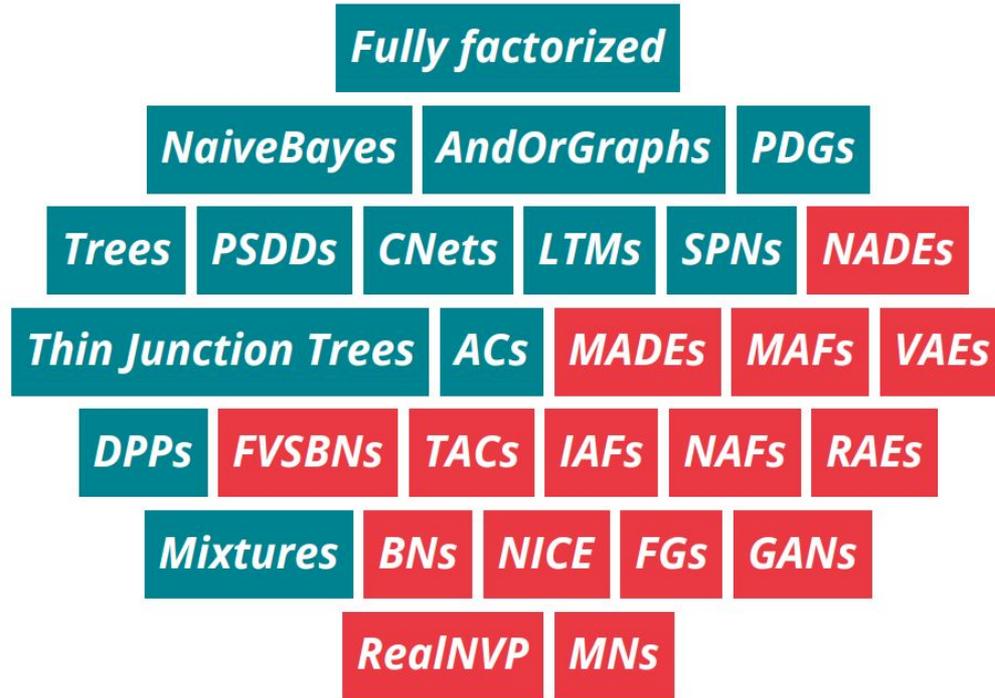


1. What are tractable probabilistic circuits?
2. Are these models any good?
3. How far can we push tractable inference?
4. What is their expressive power?

# Outline

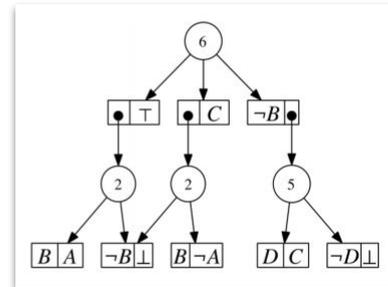
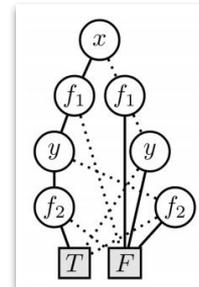
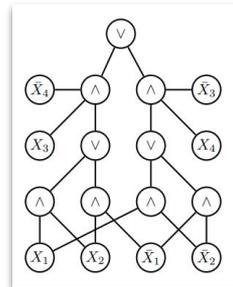
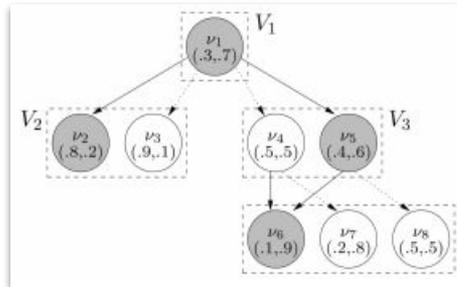
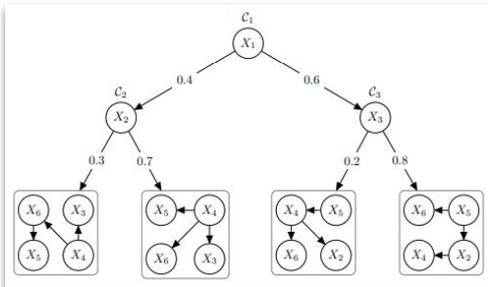
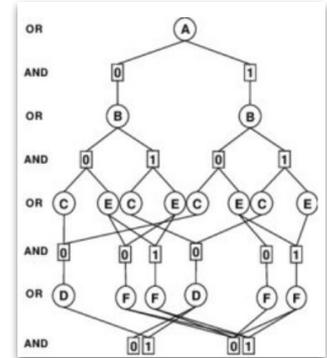
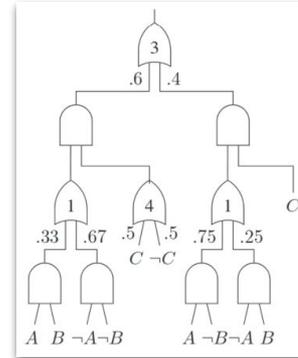
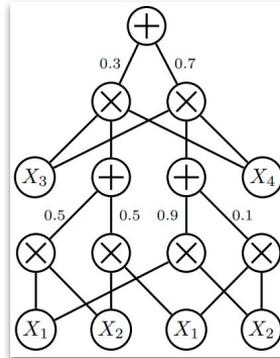
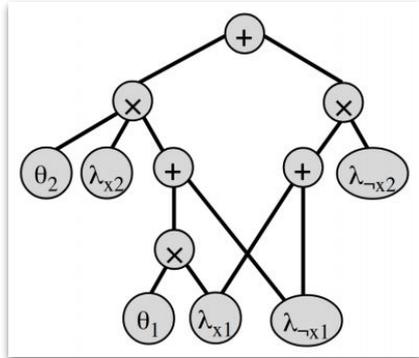
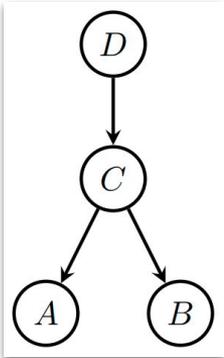


1. **What are tractable probabilistic circuits?**
2. Are these models any good?
3. How far can we push tractable inference?
4. What is their expressive power?

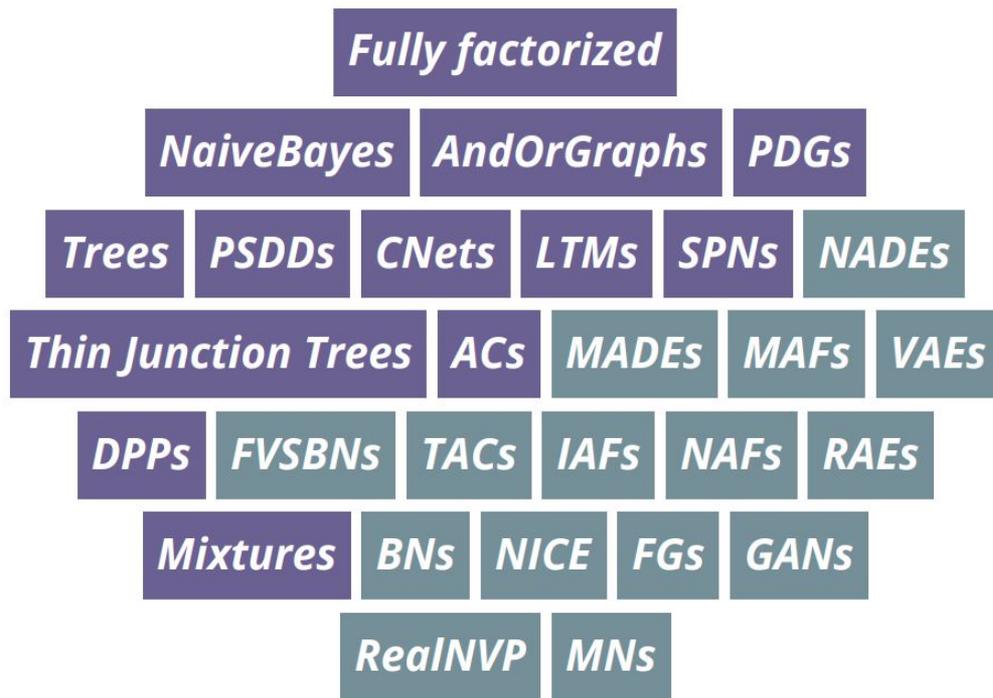


***Intractable*** and ***tractable*** models

# Tractable Probabilistic Models



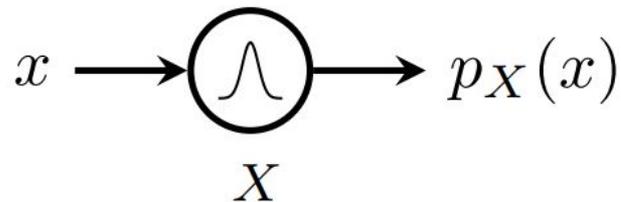
*"Every talk needs a joke and a literature overview slide, not necessarily distinct"*  
 - after Ron Graham



***a unifying framework* for tractable models**

# Probabilistic circuits

*computational graphs* that recursively define distributions



Simple distributions are tractable “black boxes” for:

- EVI: output  $p(\mathbf{x})$  (density or mass)
- MAR: output 1 (normalized) or  $Z$  (unnormalized)
- MAP: output the mode

# Probabilistic circuits

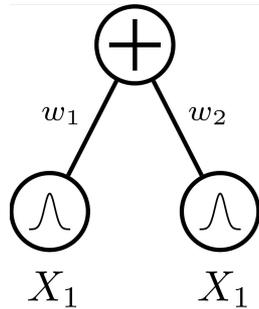
*computational graphs* that recursively define distributions



$\neg X$



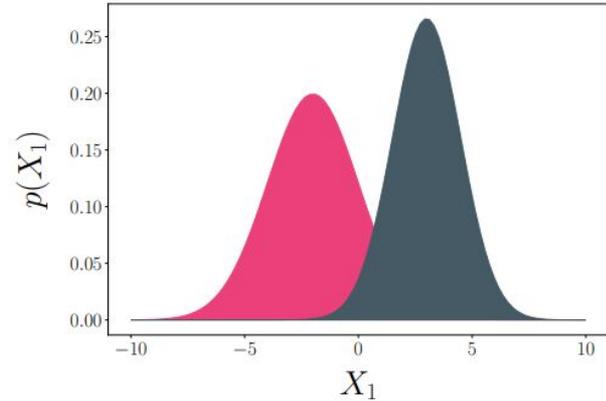
$X_1$



$$p(X_1) = w_1 p_1(X_1) + w_2 p_2(X_1)$$

$\Rightarrow$

*mixtures*



$$p(X) = p(Z = \mathbf{1}) \cdot p_1(X|Z = \mathbf{1}) \\ + p(Z = \mathbf{2}) \cdot p_2(X|Z = \mathbf{2})$$

# Probabilistic circuits

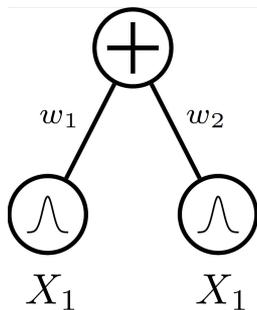
*computational graphs* that recursively define distributions



$\neg X$

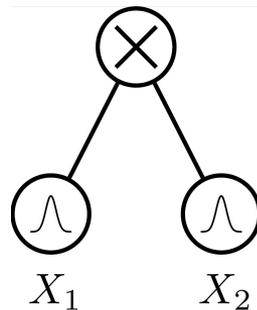


$X_1$



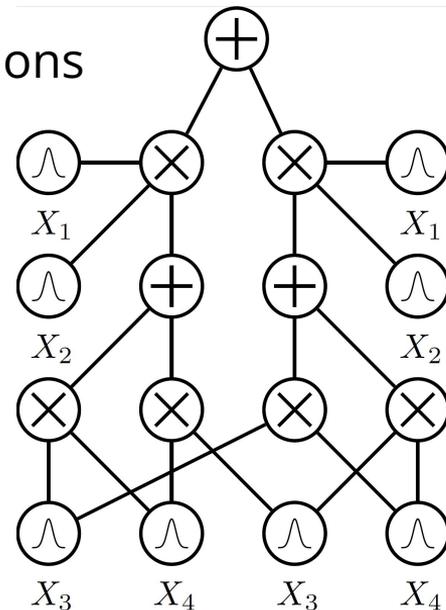
$$p(X_1) = w_1 p_1(X_1) + w_2 p_2(X_1)$$

$\Rightarrow$   
*mixtures*



$$p(X_1, X_2) = p(X_1) \cdot p(X_2)$$

$\Rightarrow$   
*factorizations*



# Tractable Probabilistic Inference

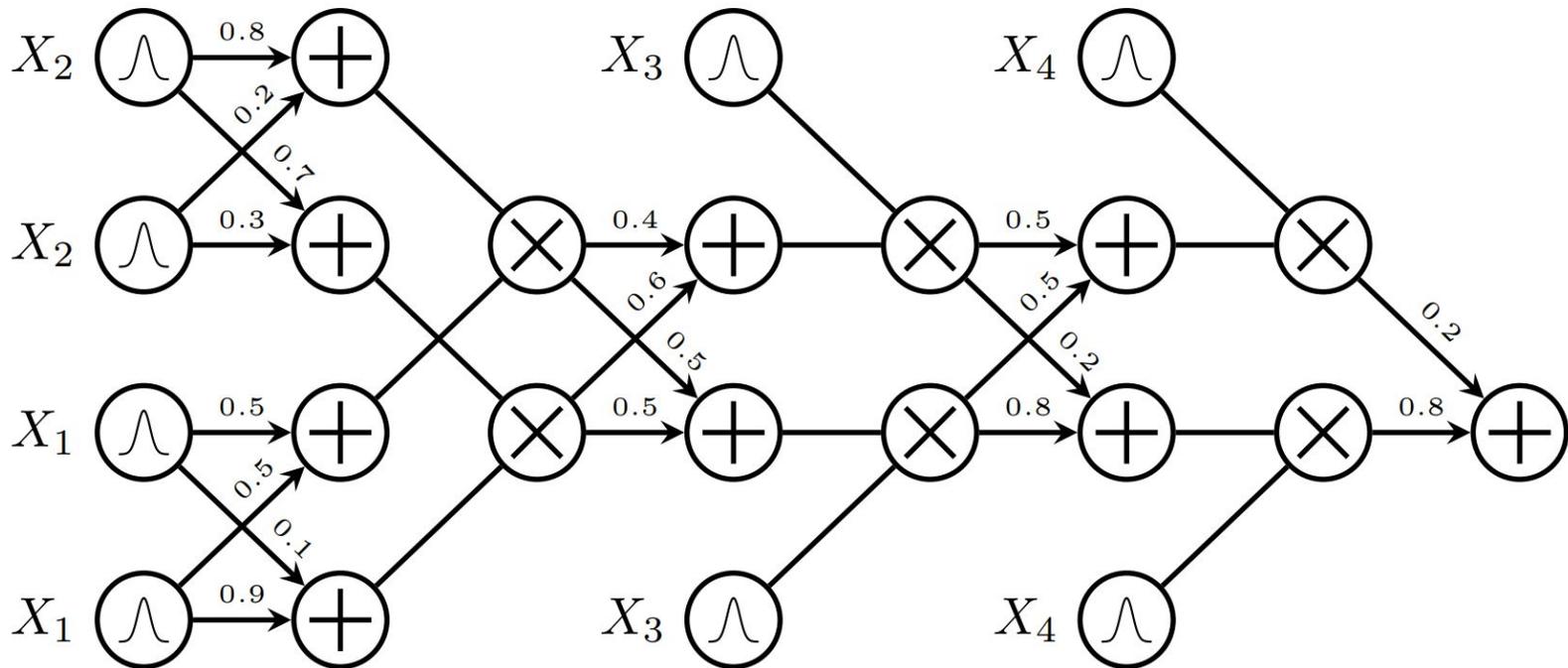
A class of queries  $\mathcal{Q}$  is tractable on a family of probabilistic models  $\mathcal{M}$  iff for any query  $q \in \mathcal{Q}$  and model  $m \in \mathcal{M}$  **exactly** computing  $q(m)$  runs in time  $O(\text{poly}(|m|))$ .

$\Rightarrow$  often poly will in fact be **linear!**

$\Rightarrow$  Note: if  $\mathcal{M}$  is compact in the number of random variables  $\mathbf{X}$ , that is,  $|m| \in O(\text{poly}(|\mathbf{X}|))$ , then query time is  $O(\text{poly}(|\mathbf{X}|))$ .

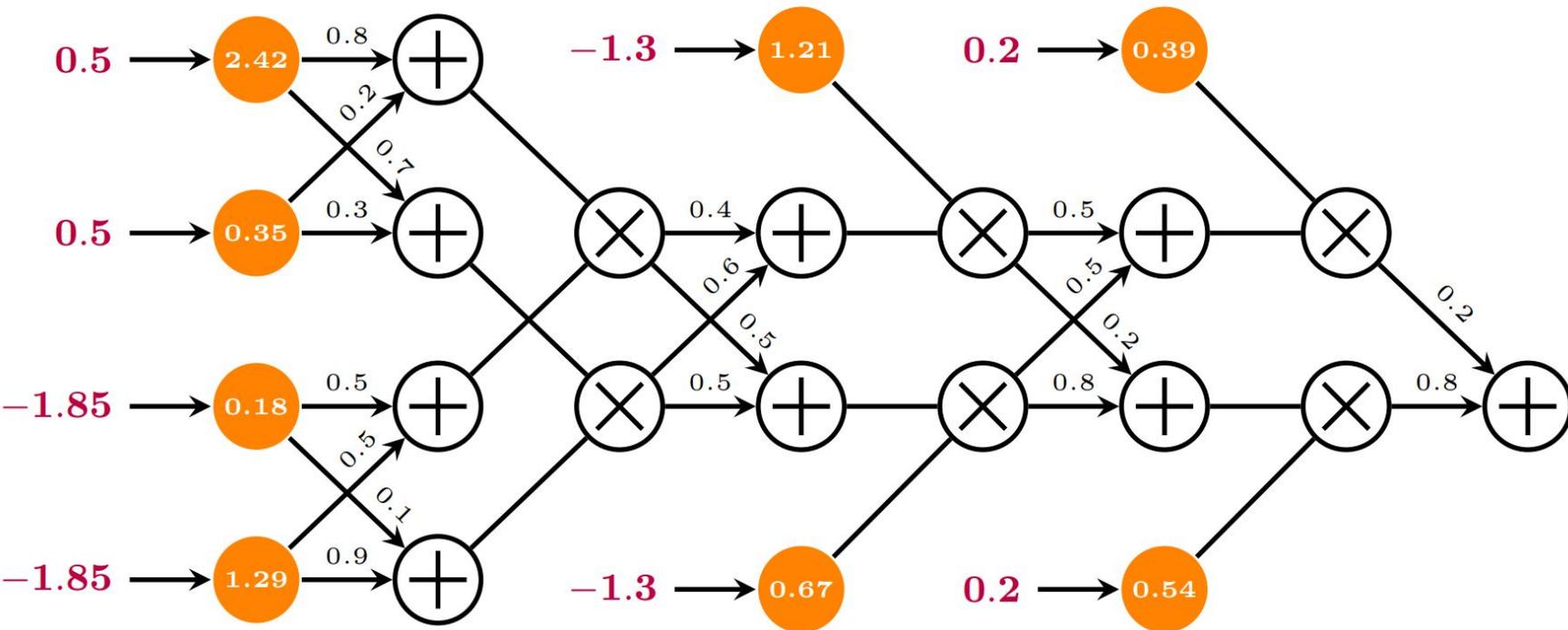
# Likelihood

$$p(X_1 = -1.85, X_2 = 0.5, X_3 = -1.3, X_4 = 0.2)$$



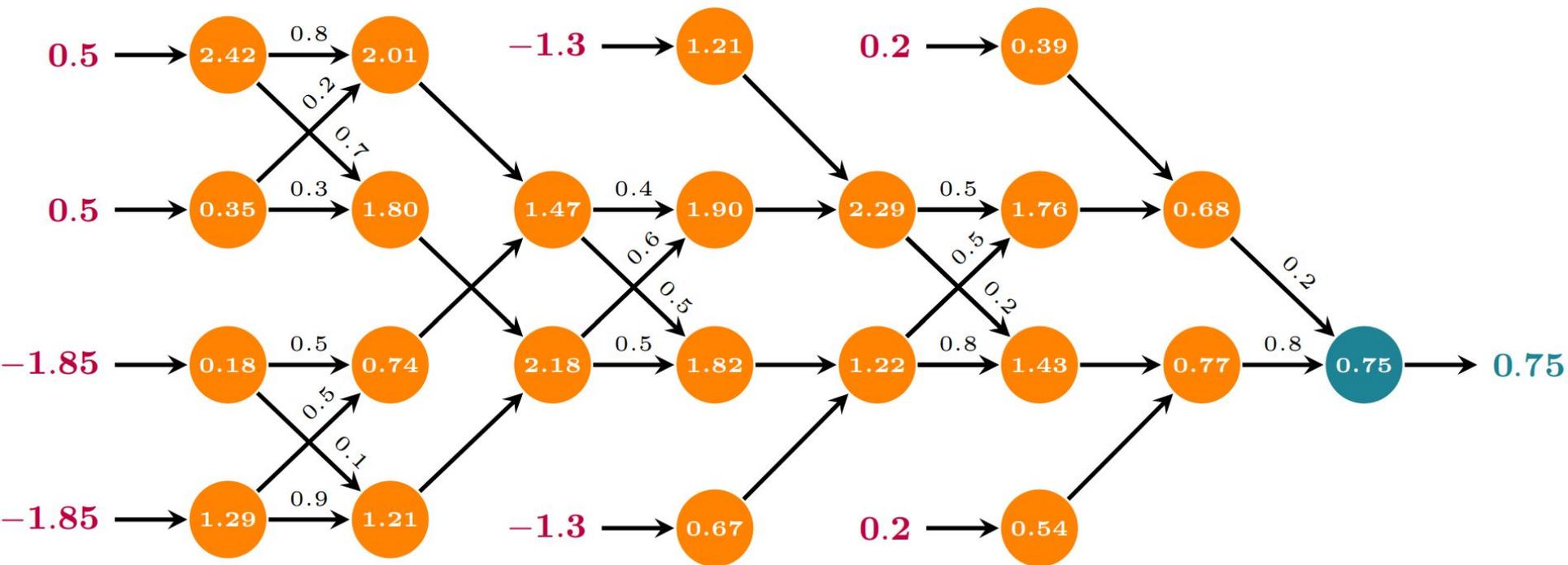
# Likelihood

$$p(X_1 = -1.85, X_2 = 0.5, X_3 = -1.3, X_4 = 0.2)$$

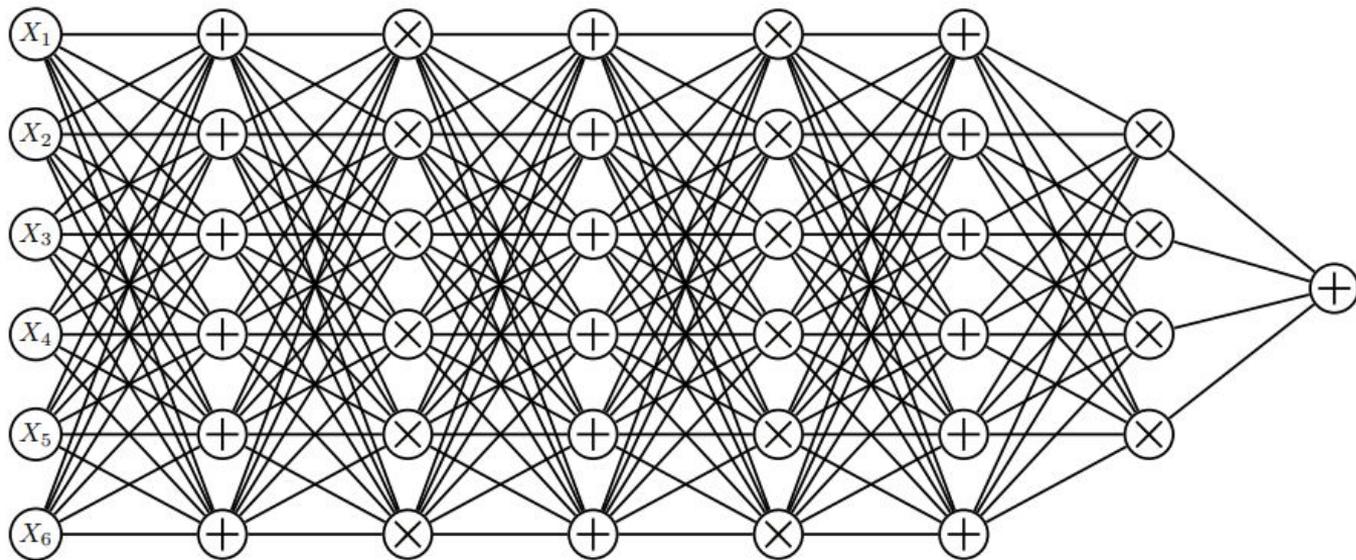


# Likelihood

$$p(X_1 = -1.85, X_2 = 0.5, X_3 = -1.3, X_4 = 0.2)$$

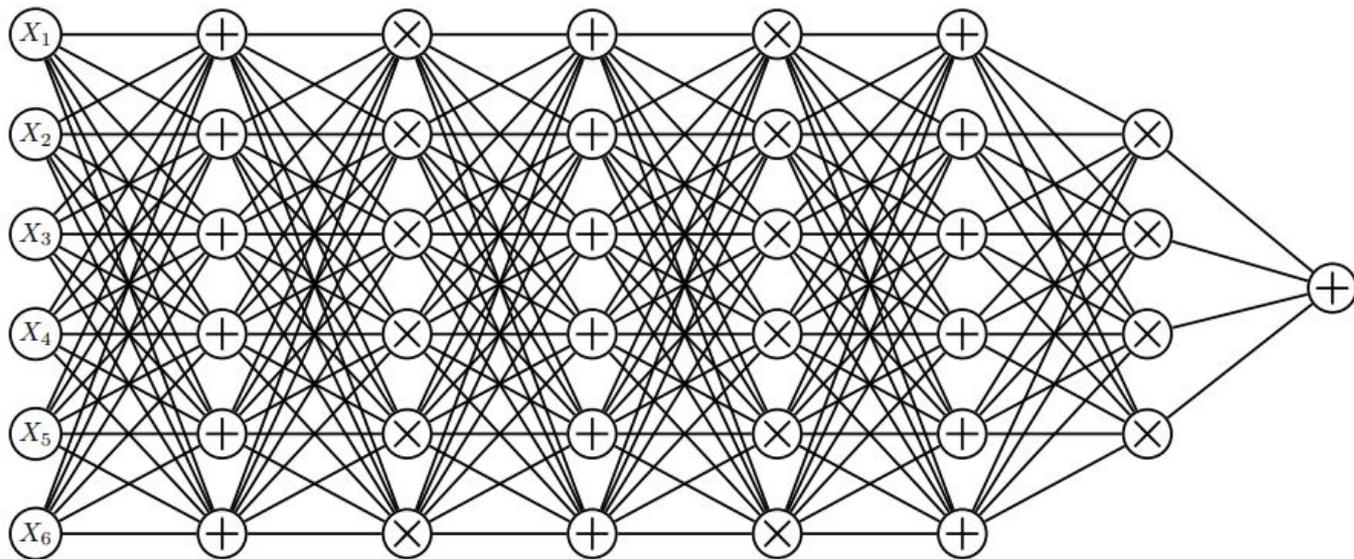


***Just sum, products and distributions?***



***just arbitrarily compose them like a neural network!***

# *Just sum, products and distributions?*



~~**just arbitrarily compose them like a neural network!**~~

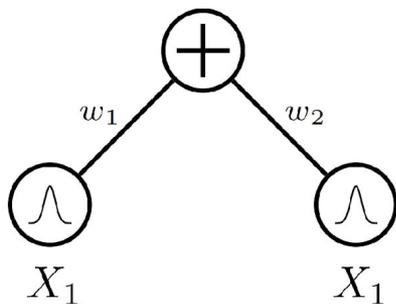


structural constraints needed for tractability

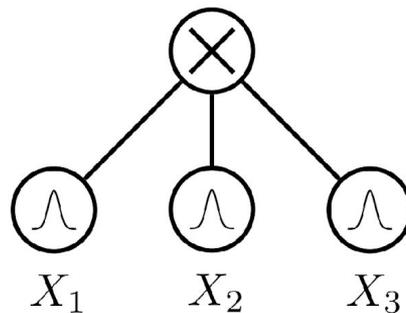
# Tractable marginals

A sum node is *smooth* if its children depend on the same set of variables.

A product node is *decomposable* if its children depend on disjoint sets of variables.



***smooth circuit***



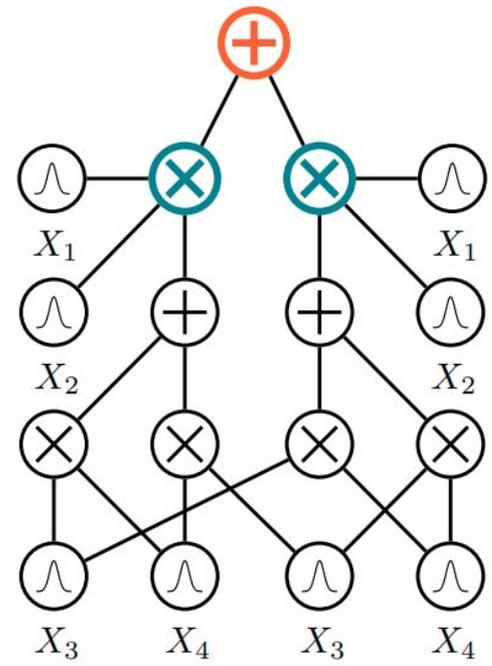
***decomposable circuit***

# Smoothness + decomposability = tractable MAR

If  $p(\mathbf{x}) = \sum_i w_i p_i(\mathbf{x})$ , (**smoothness**):

$$\int p(\mathbf{x}) d\mathbf{x} = \int \sum_i w_i p_i(\mathbf{x}) d\mathbf{x} = \sum_i w_i \int p_i(\mathbf{x}) d\mathbf{x}$$

$\Rightarrow$  integrals are "pushed down" to children

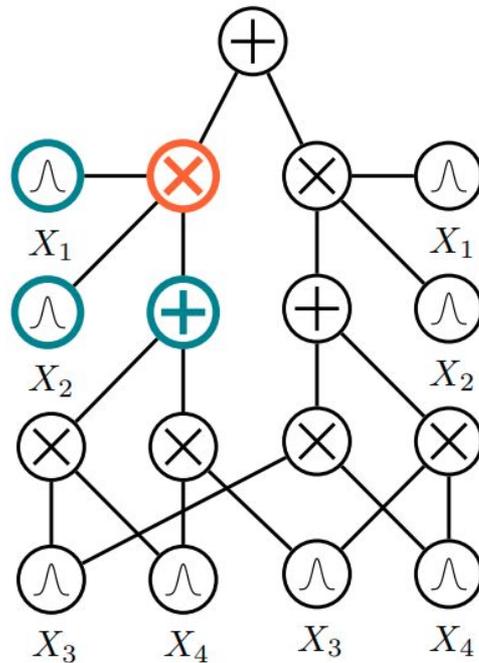


**Smoothness** + **decomposability** = **tractable MAR**

If  $\mathbf{p}(\mathbf{x}, \mathbf{y}, \mathbf{z}) = \mathbf{p}(\mathbf{x})\mathbf{p}(\mathbf{y})\mathbf{p}(\mathbf{z})$ , (**decomposability**):

$$\begin{aligned} & \int \int \int \mathbf{p}(\mathbf{x}, \mathbf{y}, \mathbf{z}) d\mathbf{x} d\mathbf{y} d\mathbf{z} = \\ &= \int \int \int \mathbf{p}(\mathbf{x})\mathbf{p}(\mathbf{y})\mathbf{p}(\mathbf{z}) d\mathbf{x} d\mathbf{y} d\mathbf{z} = \\ &= \int \mathbf{p}(\mathbf{x}) d\mathbf{x} \int \mathbf{p}(\mathbf{y}) d\mathbf{y} \int \mathbf{p}(\mathbf{z}) d\mathbf{z} \end{aligned}$$

$\Rightarrow$  integrals decompose into easier ones





# Smoothness + decomposability = tractable MAR

Forward pass evaluation for MAR

$\Rightarrow$  linear in circuit size!

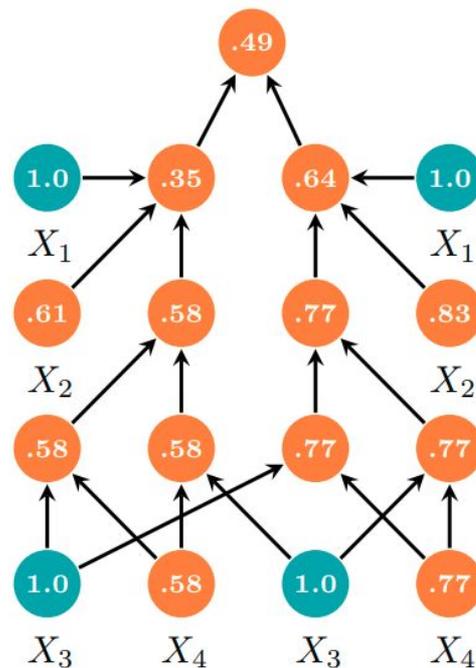
E.g. to compute  $p(x_2, x_4)$ :

leaves over  $X_1$  and  $X_3$  output  $Z_i = \int p(x_i) dx_i$

$\Rightarrow$  for normalized leaf distributions: 1.0

leaves over  $X_2$  and  $X_4$  output **EVI**

feedforward evaluation (bottom-up)

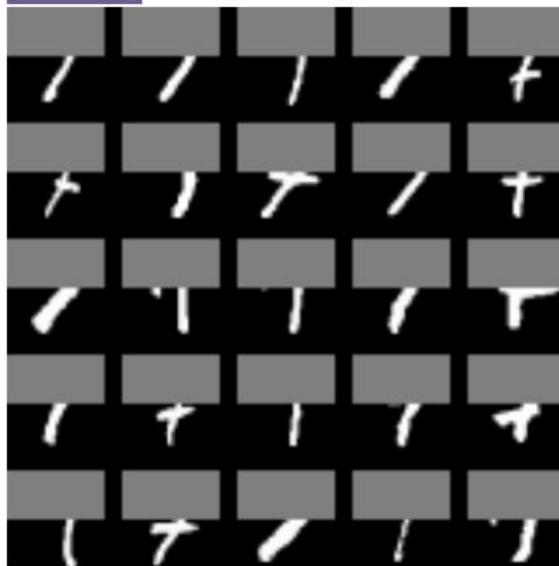


# Tractable MAR on PCs (Einsum Networks)

**EVI** 10,958.72 nats



**MAR** 5,387.55 nats



**Smoothness** + **decomposability** = ~~tractable MAP~~

We **cannot** decompose bottom-up a MAP query:

$$\max_{\mathbf{q}} p(\mathbf{q} \mid \mathbf{e})$$

since for a sum node we are marginalizing out a latent variable

$$\max_{\mathbf{q}} \sum_i w_i p_i(\mathbf{q}, \mathbf{e}) = \max_{\mathbf{q}} \sum_{\mathbf{z}} p(\mathbf{q}, \mathbf{z}, \mathbf{e}) \neq \sum_{\mathbf{z}} \max_{\mathbf{q}} p(\mathbf{q}, \mathbf{z}, \mathbf{e})$$

$\Rightarrow$  MAP for latent variable models is **intractable** [Conaty et al. 2017]

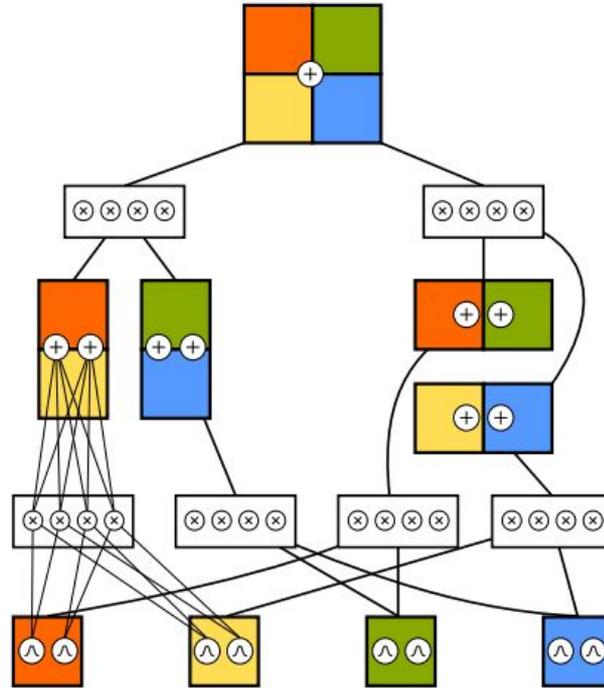
# Outline



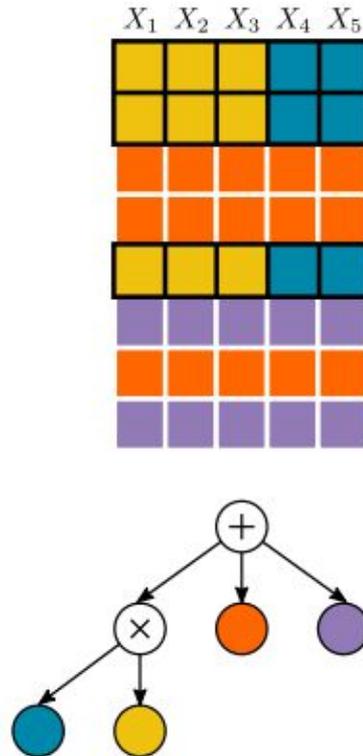
1. What are tractable probabilistic circuits?
2. **Are these models any good?**
3. How far can we push tractable inference?
4. What is their expressive power?

Where do architectures come from?

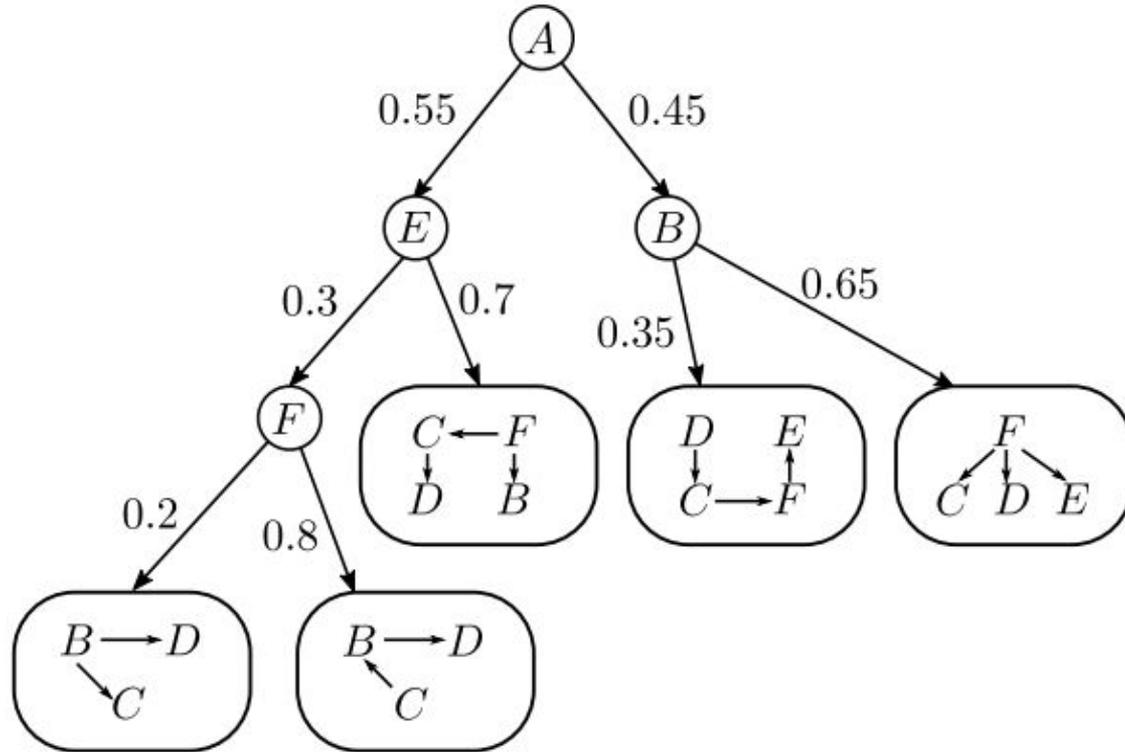
# Where do architectures come from?



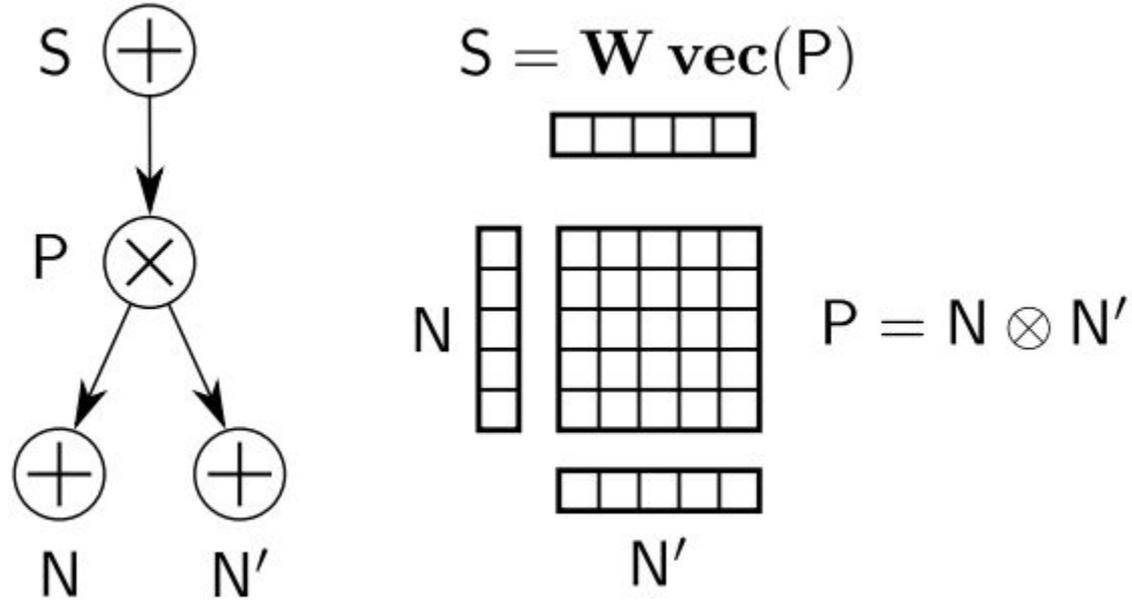
# Where do architectures come from?



# Where do architectures come from?

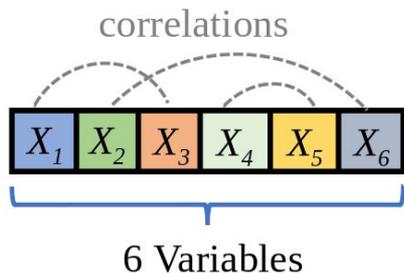


# Where do architectures come from?

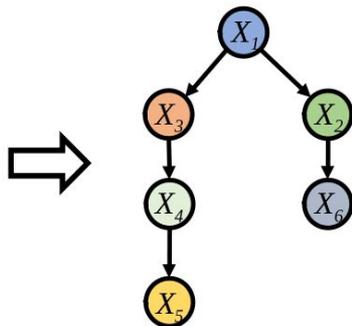


# Learning Expressive Probabilistic Circuits

## Hidden Chow-Liu Trees

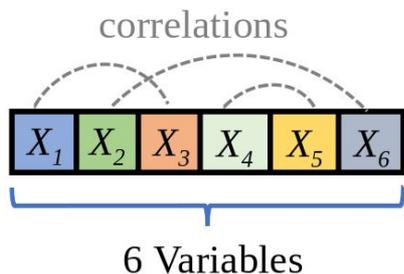


Learned **CLT structure**  
captures strong pairwise  
dependencies

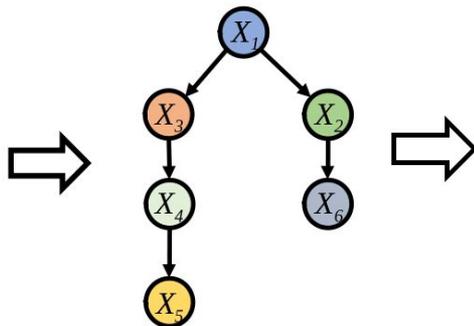


# Learning Expressive Probabilistic Circuits

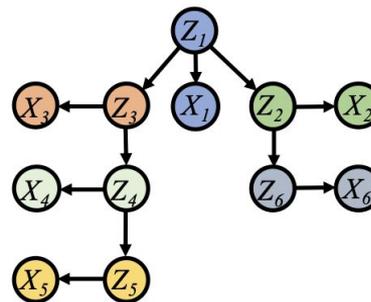
## Hidden Chow-Liu Trees



Learned **CLT** structure captures strong pairwise dependencies



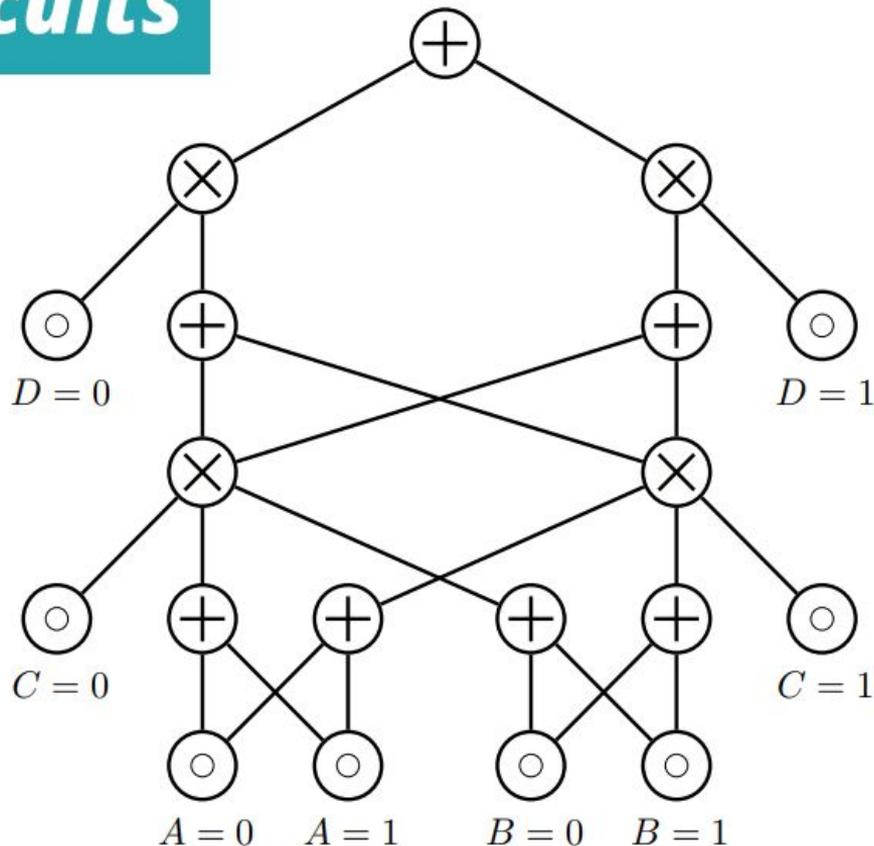
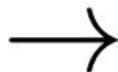
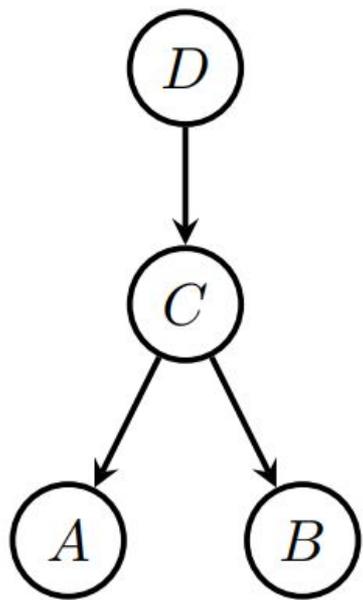
Learned **HCLT** structure



⇒ **Compile into an equivalent PC**

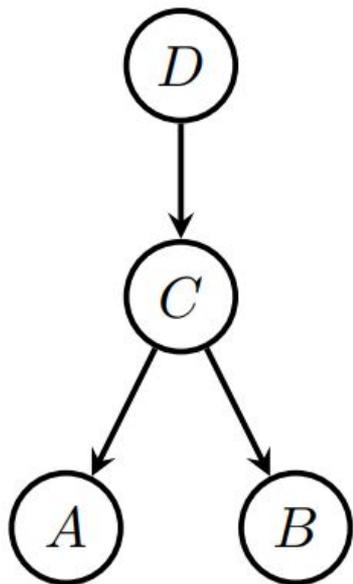
# From BN trees to circuits

*via compilation*



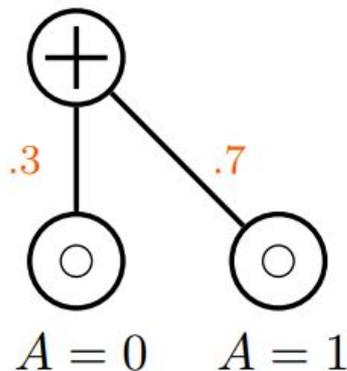
# From BN trees to circuits

*via compilation*



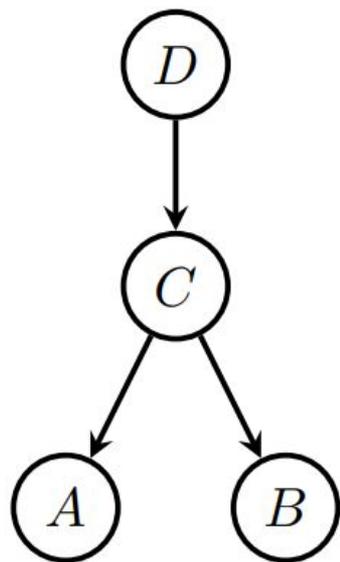
...compile a leaf CPT

$p(A|C = 0)$

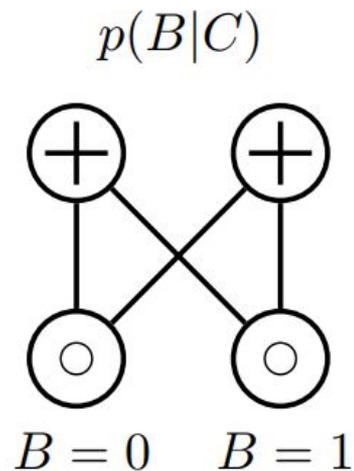
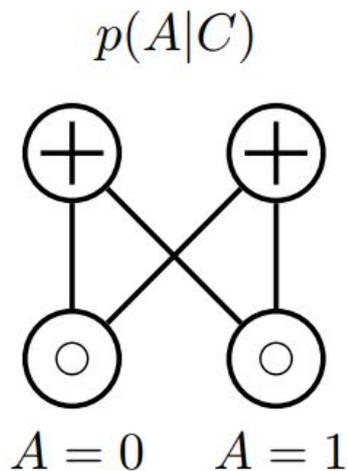


# From BN trees to circuits

*via compilation*



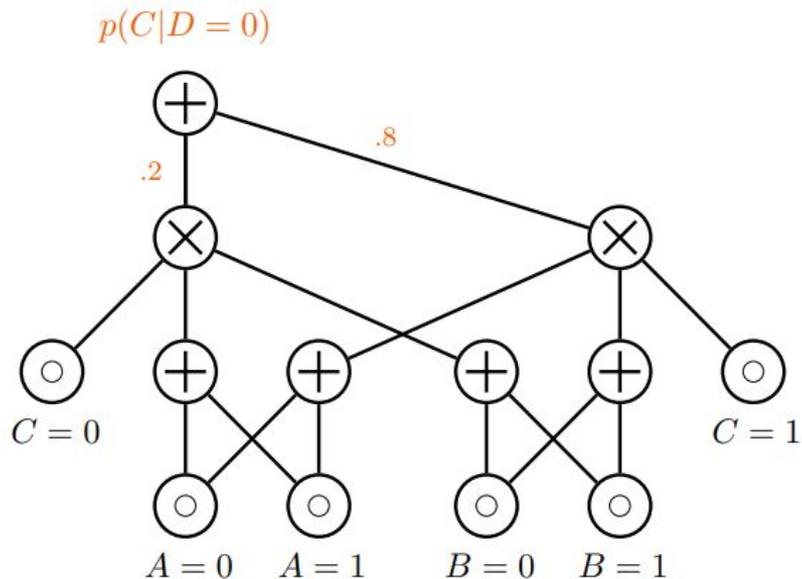
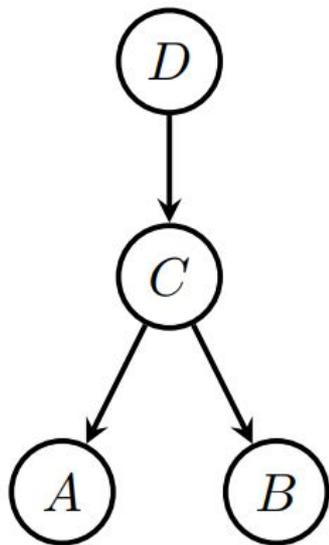
...compile a leaf CPT...for all leaves...



# From BN trees to circuits

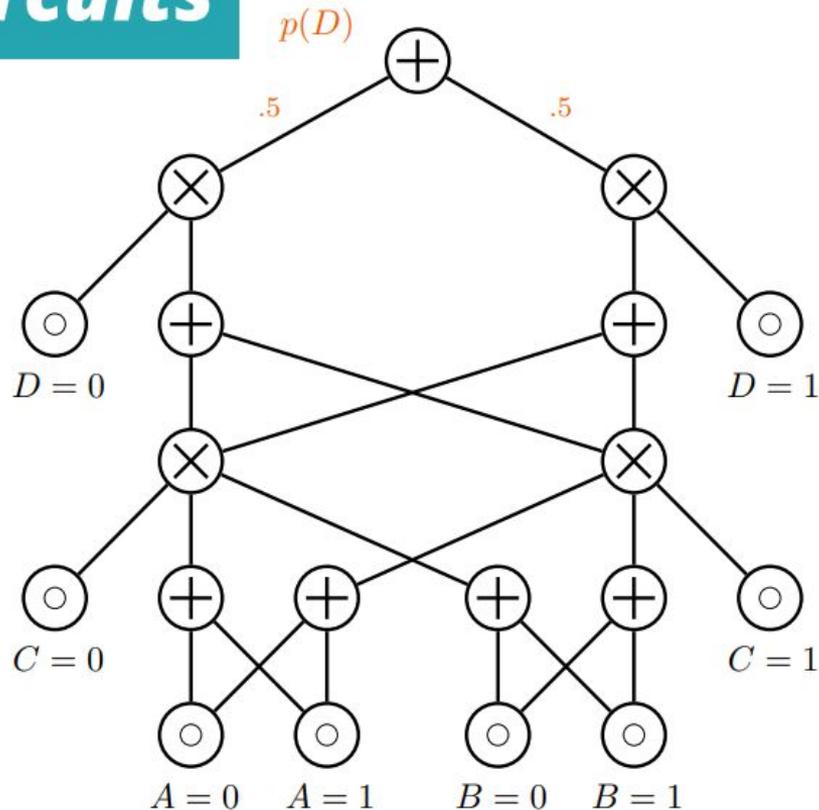
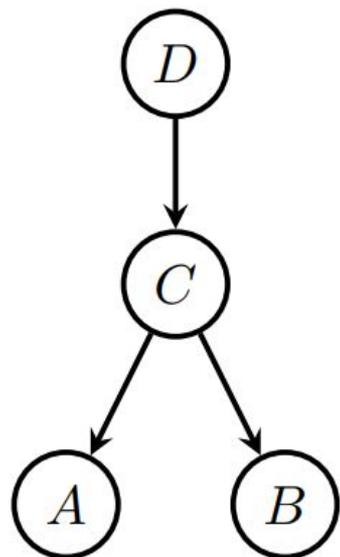
via compilation

...and recurse over parents...



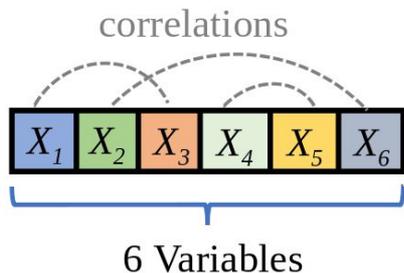
# From BN trees to circuits

via compilation

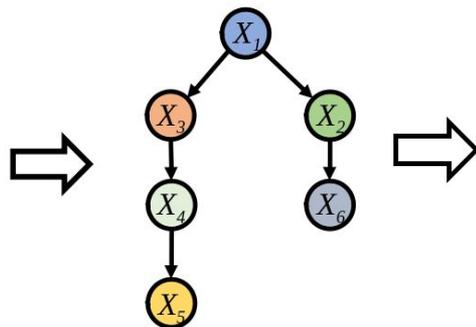


# Learning Expressive Probabilistic Circuits

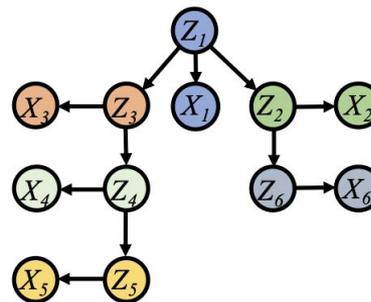
## Hidden Chow-Liu Trees



Learned **CLT structure**  
captures strong pairwise  
dependencies



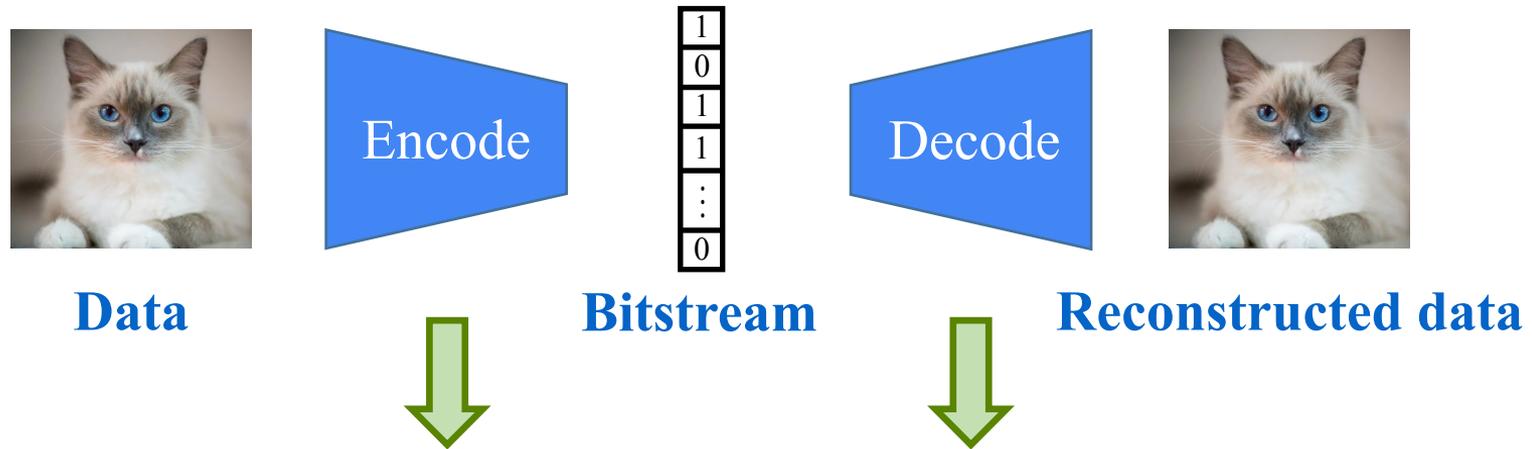
Learned **HCLT structure**



⇒ **Compile into an  
equivalent PC**

⇒ **Mini-batch Stochastic  
Expectation Maximization**

# Lossless Data Compression



Expressive probabilistic model  $p(\mathbf{x})$

+

Efficient coding algorithm



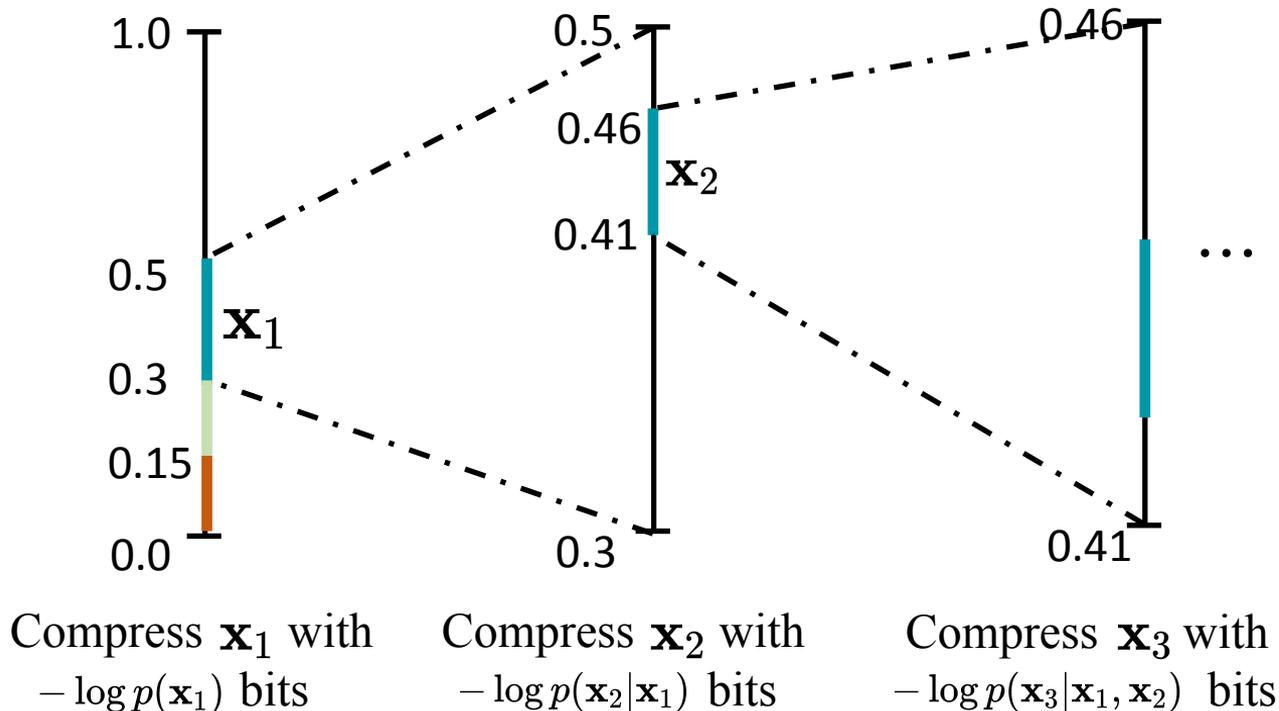
Determines the theoretical limit of compression rate



How close we can approach the theoretical limit

# A Typical Streaming Code – Arithmetic Coding

We want to compress a set of variables (e.g., pixels, letters)  $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\}$



Need to compute

$$p(X_1 < x_1)$$

$$p(X_1 \leq x_1)$$

$$p(X_2 < x_2 | x_1)$$

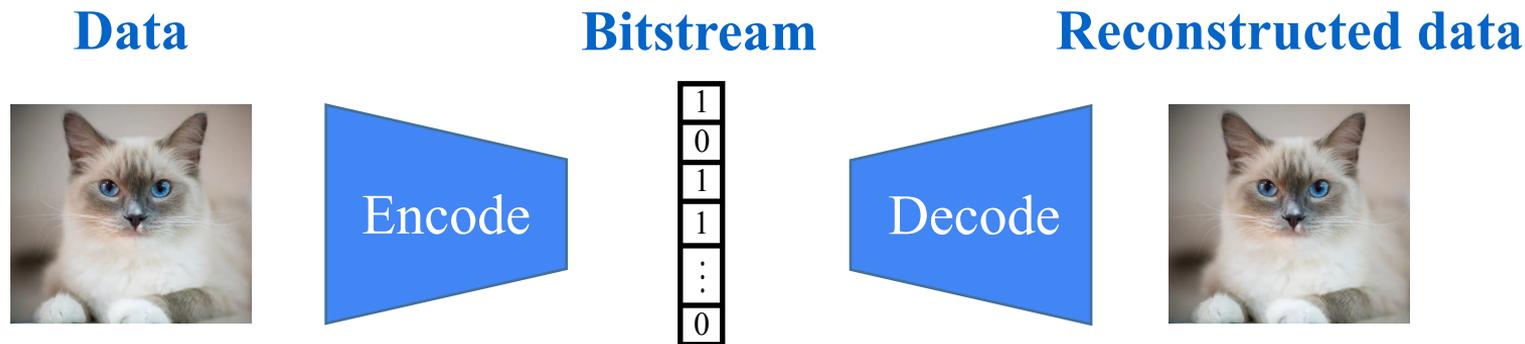
$$p(X_2 \leq x_2 | x_1)$$

$$p(X_3 < x_3 | x_1, x_2)$$

$$p(X_3 \leq x_3 | x_1, x_2)$$

⋮

# Lossless Neural Compression with Probabilistic Circuits



## Probabilistic Circuits

- **Expressive** → SoTA likelihood on MNIST.
- **Fast** → Time complexity of en/decoding is  $\mathbf{O}(|p| \log(\mathbf{D}))$ , where  $\mathbf{D}$  is the # variables and  $|p|$  is the size of the PC.

## Arithmetic Coding:

$$\begin{aligned} & p(X_1 < x_1) \\ & p(X_1 \leq x_1) \\ & p(X_2 < x_2 | x_1) \\ & p(X_2 \leq x_2 | x_1) \\ & p(X_3 < x_3 | x_1, x_2) \\ & p(X_3 \leq x_3 | x_1, x_2) \\ & \vdots \end{aligned}$$

# Lossless Neural Compression with Probabilistic Circuits

## SoTA compression rates

Dataset	HCLT (ours)	IDF	BitSwap	BB-ANS	JPEG2000	WebP	McBits
MNIST	<b>1.24</b> (1.20)	1.96 (1.90)	1.31 (1.27)	1.42 (1.39)	3.37	2.09	(1.98)
FashionMNIST	3.37 (3.34)	3.50 (3.47)	<b>3.35</b> (3.28)	3.69 (3.66)	3.93	4.62	(3.72)
EMNIST (Letter)	<b>1.84</b> (1.80)	2.02 (1.95)	1.90 (1.84)	2.29 (2.26)	3.62	3.31	(3.12)
EMNIST (ByClass)	<b>1.89</b> (1.85)	2.04 (1.98)	1.91 (1.87)	2.24 (2.23)	3.61	3.34	(3.14)

Compress and decompress 5-40x faster than NN methods with similar bitrates

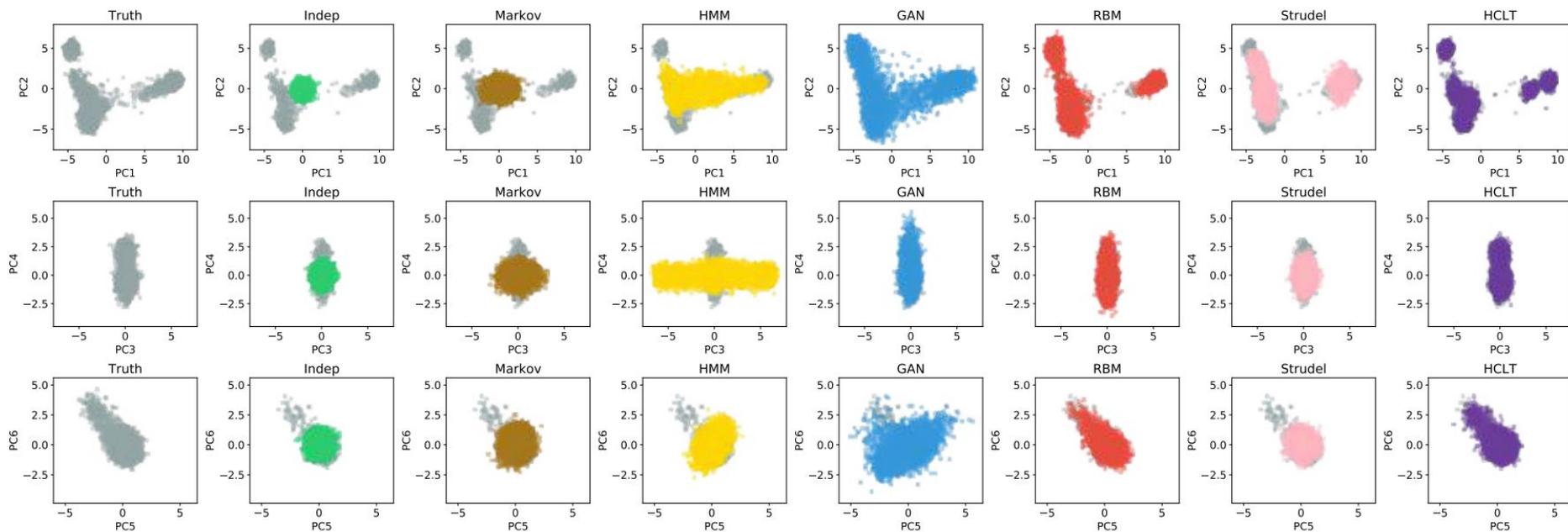
Method	# parameters	Theoretical bpd	Codeword bpd	Comp. time (s)	Decomp. time (s)
PC (HCLT, $M=16$ )	3.3M	1.26	1.30	9	44
PC (HCLT, $M=24$ )	5.1M	1.22	1.26	15	86
PC (HCLT, $M=32$ )	7.0M	1.20	1.24	26	142
IDF	24.1M	1.90	1.96	288	592
BitSwap	2.8M	1.27	1.31	578	326

# Lossless Neural Compression with Probabilistic Circuits

Can be effectively combined with Flow models to achieve better generative performance

Model	CIFAR10	ImageNet32	ImageNet64
RealNVP	3.49	4.28	3.98
Glow	3.35	4.09	3.81
IDF	3.32	4.15	3.90
IDF++	<b>3.24</b>	4.10	3.81
PC+IDF	3.28	<b>3.99</b>	<b>3.71</b>

# Tractable and expressive generative models of genetic variation data



## Text data

Dataset	PC	Bipartite flow	AF/SCF	IAF/SCF
Penn Treebank	<b>1.23</b>	1.38	1.46	1.63

# Training SotA likelihood full MNIST probabilistic circuit model in ~7 minutes on GPU:

[https://github.com/Juice-jl/ProbabilisticCircuits.jl/blob/master/examples/train\\_mnist\\_hclt.ipynb](https://github.com/Juice-jl/ProbabilisticCircuits.jl/blob/master/examples/train_mnist_hclt.ipynb)

Public

<> Code Issues 10 Pull requests Discussions Actions Projects Wiki ...

master ProbabilisticCircuits.jl / examples / train\_mnist\_hclt.ipynb Go to file ...

liuanji update demo notebook ✓ Latest commit c0e062e 2 days ago History

1 contributor

609 lines (609 sloc) | 26.5 KB

Dataset	PC (ours)	IDF	Hierarchical VAE	PixelVAE
MNIST	<b>1.20</b>	2.90	1.27	1.39
FashionMNIST	3.34	3.47	<b>3.28</b>	3.66
EMNIST (Letter split)	<b>1.80</b>	1.95	1.84	2.26
EMNIST (ByClass split)	<b>1.85</b>	1.98	1.87	2.23

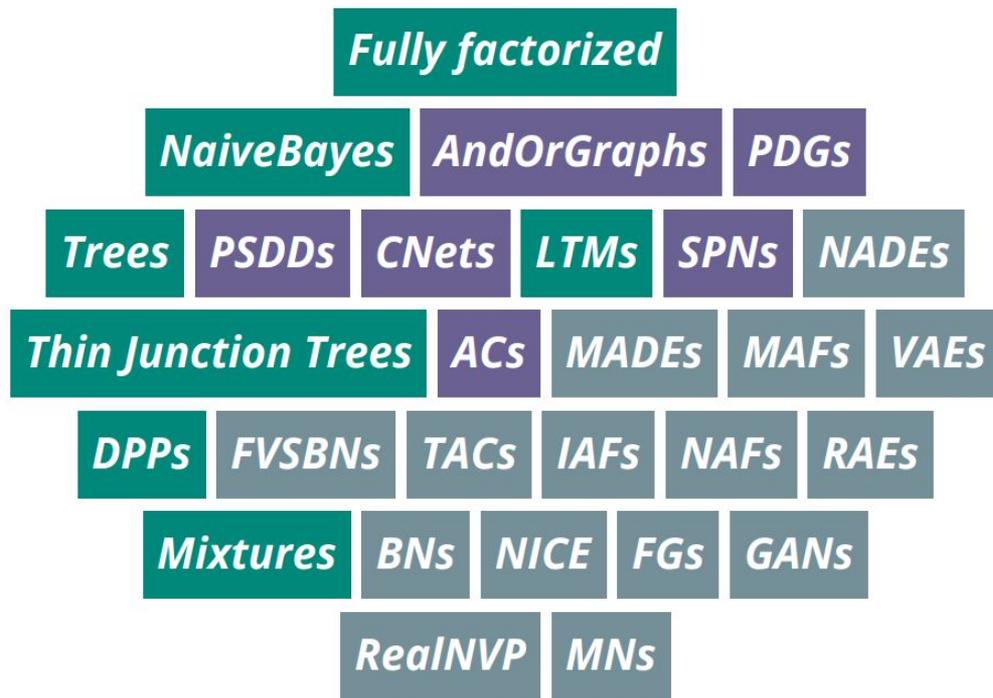
\* Note: all reported numbers are bits-per-dimension (bpd). The results are extracted from [1].

[1] Anji Liu, Stephan Mandt and Guy Van den Broeck. Lossless Compression with Probabilistic Circuits, In International Conference on Learning Representations (ICLR), 2022.

We start by importing ProbabilisticCircuits.jl and other required packages:

```
In [1]: using ProbabilisticCircuits
using MLDatasets
using CUDA
```

We first load the MNIST dataset from MLDatasets.jl and move them to GPU:



***Expressive* models without *compromises***

# Outline



1. What are tractable probabilistic circuits?
2. Are these models any good?
3. **How far can we push tractable inference?**
4. What is their expressive power?

**Smoothness** + **decomposability** = ~~tractable MAP~~

We **cannot** decompose bottom-up a MAP query:

$$\max_{\mathbf{q}} p(\mathbf{q} \mid \mathbf{e})$$

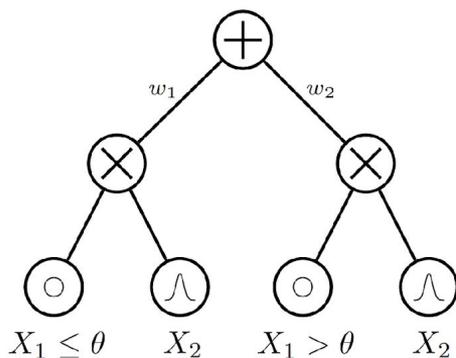
since for a sum node we are marginalizing out a latent variable

$$\max_{\mathbf{q}} \sum_i w_i p_i(\mathbf{q}, \mathbf{e}) = \max_{\mathbf{q}} \sum_{\mathbf{z}} p(\mathbf{q}, \mathbf{z}, \mathbf{e}) \neq \sum_{\mathbf{z}} \max_{\mathbf{q}} p(\mathbf{q}, \mathbf{z}, \mathbf{e})$$

$\Rightarrow$  MAP for latent variable models is **intractable** [Conaty et al. 2017]

# Determinism

A sum node is **deterministic** if only one of its children outputs non-zero for any input



**deterministic circuit**

$\Rightarrow$  allows tractable MAP inference

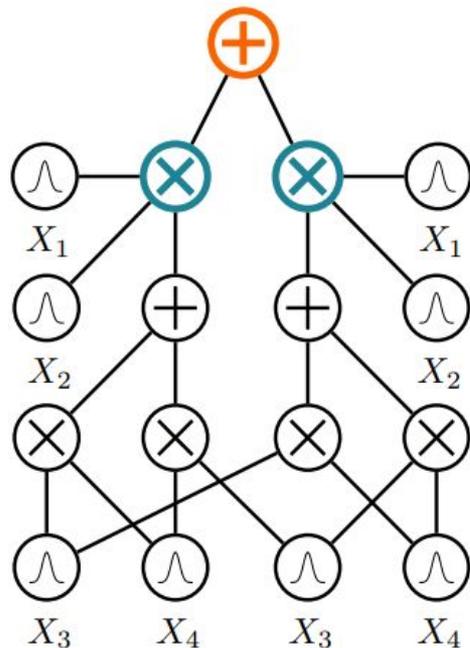
$$\operatorname{argmax}_{\mathbf{x}} p(\mathbf{x})$$

# Determinism + decomposability = tractable MAP

If  $p(\mathbf{q}, \mathbf{e}) = \sum_i w_i p_i(\mathbf{q}, \mathbf{e}) = \max_i w_i p_i(\mathbf{q}, \mathbf{e})$ ,  
 (**deterministic** sum node):

$$\begin{aligned} \max_{\mathbf{q}} p(\mathbf{q}, \mathbf{e}) &= \max_{\mathbf{q}} \sum_i w_i p_i(\mathbf{q}, \mathbf{e}) \\ &= \max_{\mathbf{q}} \max_i w_i p_i(\mathbf{q}, \mathbf{e}) \\ &= \max_i \max_{\mathbf{q}} w_i p_i(\mathbf{q}, \mathbf{e}) \end{aligned}$$

$\Rightarrow$  one non-zero child term, thus sum is max

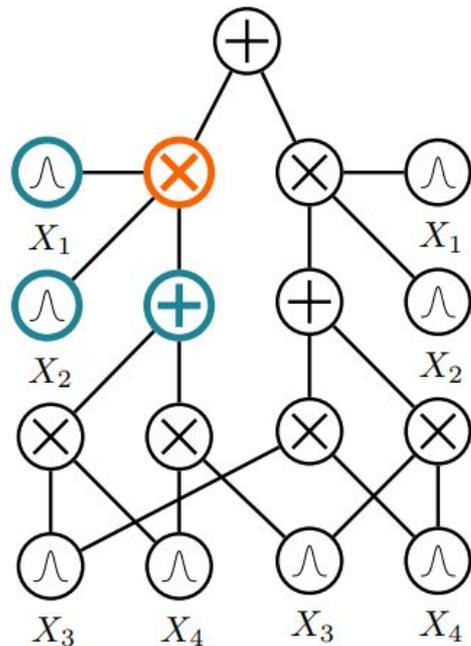


# **Determinism** + **decomposability** = **tractable MAP**

If  $p(\mathbf{q}, \mathbf{e}) = p(\mathbf{q}_x, \mathbf{e}_x, \mathbf{q}_y, \mathbf{e}_y) = p(\mathbf{q}_x, \mathbf{e}_x)p(\mathbf{q}_y, \mathbf{e}_y)$   
(**decomposable** product node):

$$\begin{aligned}\max_{\mathbf{q}} p(\mathbf{q} \mid \mathbf{e}) &= \max_{\mathbf{q}} p(\mathbf{q}, \mathbf{e}) \\ &= \max_{\mathbf{q}_x, \mathbf{q}_y} p(\mathbf{q}_x, \mathbf{e}_x, \mathbf{q}_y, \mathbf{e}_y) \\ &= \max_{\mathbf{q}_x} p(\mathbf{q}_x, \mathbf{e}_x) \cdot \max_{\mathbf{q}_y} p(\mathbf{q}_y, \mathbf{e}_y)\end{aligned}$$

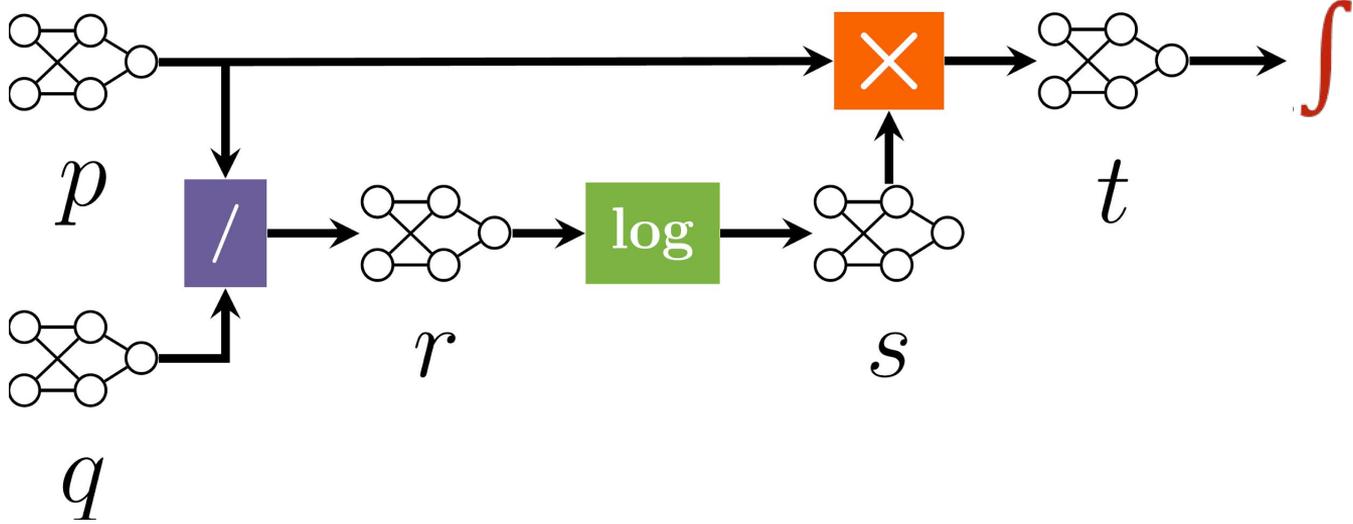
$\Rightarrow$  solving optimization independently





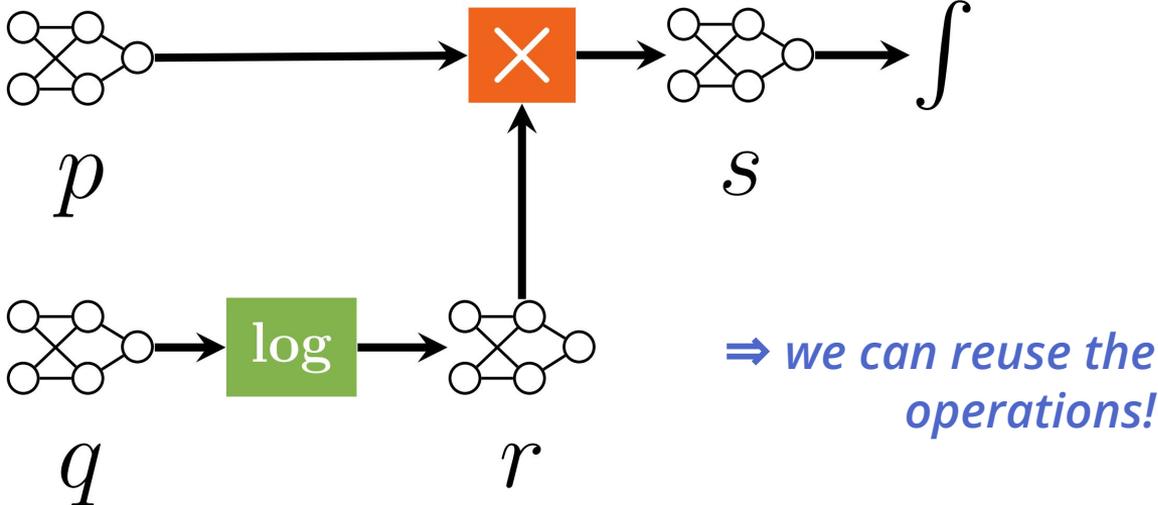
# Queries as pipelines: KLD

$$\text{KLD}(p \parallel q) = \int p(\mathbf{x}) \times \log((p(\mathbf{x})/q(\mathbf{x})))d\mathbf{X}$$

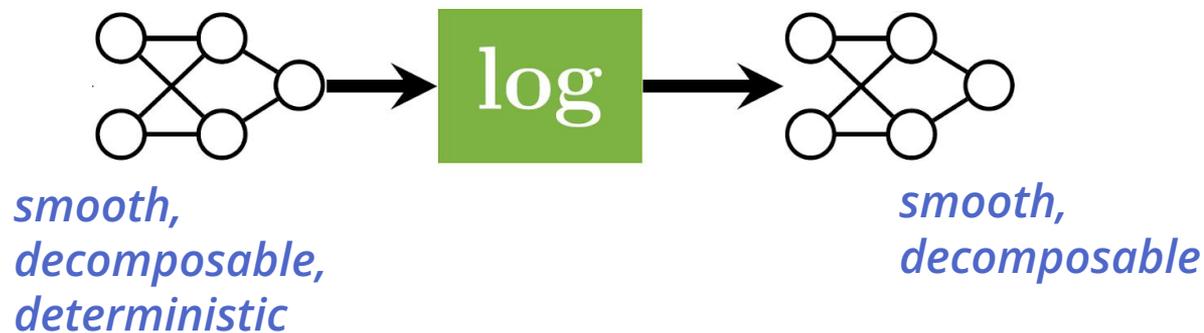


# Queries as pipelines: Cross Entropy

$$H(p, q) = \int p(\mathbf{x}) \times \log(q(\mathbf{x})) d\mathbf{X}$$



Operation	Tractability	
	Input conditions	Output conditions
LOG	Sm, Dec, Det	Sm, Dec



# Tractable circuit operations

Operation		Tractability		Hardness
		Input properties	Output properties	
SUM	$\theta_1 p + \theta_2 q$	(+Cmp)	(+SD)	NP-hard for Det output
PRODUCT	$p \cdot q$	Cmp (+Det, +SD)	Dec (+Det, +SD)	#P-hard w/o Cmp
POWER	$p^n, n \in \mathbb{N}$	SD (+Det)	SD (+Det)	#P-hard w/o SD
	$p^\alpha, \alpha \in \mathbb{R}$	Sm, Dec, Det (+SD)	Sm, Dec, Det (+SD)	#P-hard w/o Det
QUOTIENT	$p/q$	Cmp; $q$ Det (+ $p$ Det,+SD)	Dec (+Det,+SD)	#P-hard w/o Det
LOG	$\log(p)$	Sm, Dec, Det	Sm, Dec	#P-hard w/o Det
EXP	$\exp(p)$	linear	SD	#P-hard

# Inference by tractable operations

*systematically derive* tractable inference algorithm of complex queries

	Query	Tract. Conditions	Hardness
CROSS ENTROPY	$-\int p(\mathbf{x}) \log q(\mathbf{x}) d\mathbf{X}$	Cmp, $q$ Det	#P-hard w/o Det
SHANNON ENTROPY	$-\sum p(\mathbf{x}) \log p(\mathbf{x})$	Sm, Dec, Det	coNP-hard w/o Det
RÉNYI ENTROPY	$(1 - \alpha)^{-1} \log \int p^\alpha(\mathbf{x}) d\mathbf{X}, \alpha \in \mathbb{N}$	SD	#P-hard w/o SD
MUTUAL INFORMATION	$(1 - \alpha)^{-1} \log \int p^\alpha(\mathbf{x}) d\mathbf{X}, \alpha \in \mathbb{R}_+$	Sm, Dec, Det	#P-hard w/o Det
KULLBACK-LEIBLER DIV.	$\int p(\mathbf{x}, \mathbf{y}) \log(p(\mathbf{x}, \mathbf{y}) / (p(\mathbf{x})p(\mathbf{y})))$	Sm, SD, Det*	coNP-hard w/o SD
RÉNYI'S ALPHA DIV.	$\int p(\mathbf{x}) \log(p(\mathbf{x}) / q(\mathbf{x})) d\mathbf{X}$	Cmp, Det	#P-hard w/o Det
ITAKURA-SAITO DIV.	$(1 - \alpha)^{-1} \log \int p^\alpha(\mathbf{x}) q^{1-\alpha}(\mathbf{x}) d\mathbf{X}, \alpha \in \mathbb{N}$	Cmp, $q$ Det	#P-hard w/o Det
CAUCHY-SCHWARZ DIV.	$(1 - \alpha)^{-1} \log \int p^\alpha(\mathbf{x}) q^{1-\alpha}(\mathbf{x}) d\mathbf{X}, \alpha \in \mathbb{R}$	Cmp, Det	#P-hard w/o Det
SQUARED LOSS	$\int [p(\mathbf{x}) / q(\mathbf{x}) - \log(p(\mathbf{x}) / q(\mathbf{x})) - 1] d\mathbf{X}$	Cmp, Det	#P-hard w/o Det
	$-\log \frac{\int p(\mathbf{x}) q(\mathbf{x}) d\mathbf{X}}{\sqrt{\int p^2(\mathbf{x}) d\mathbf{X} \int q^2(\mathbf{x}) d\mathbf{X}}}$	Cmp	#P-hard w/o Cmp
	$\int (p(\mathbf{x}) - q(\mathbf{x}))^2 d\mathbf{X}$	Cmp	#P-hard w/o Cmp

# Even harder queries

## *Marginal MAP*

Given a set of query variables  $\mathbf{Q} \subset \mathbf{X}$  and evidence  $\mathbf{e}$ ,

find:  $\operatorname{argmax}_{\mathbf{q}} p(\mathbf{q}|\mathbf{e})$

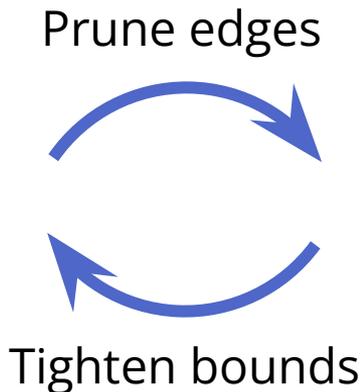
*⇒ i.e. MAP of a marginal distribution on  $\mathbf{Q}$*

**!** *NP<sup>PP</sup>-complete* for PGMs

**!** *NP-hard* even for PCs tractable for marginals, MAP & entropy



# Iterative MMAP solver



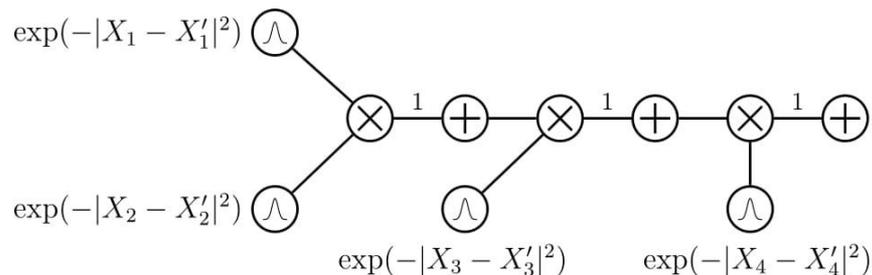
Dataset	runtime (# solved)	
	search	pruning
NLTCS	<b>0.01</b> (10)	0.63 (10)
MSNBC	<b>0.03</b> (10)	0.73 (10)
KDD	<b>0.04</b> (10)	0.68 (10)
Plants	2.95 (10)	<b>2.72</b> (10)
Audio	2041.33 (6)	<b>13.70</b> (10)
Jester	2913.04 (2)	<b>14.74</b> (10)
Netflix	- (0)	<b>47.18</b> (10)
Accidents	109.56 (10)	<b>15.86</b> (10)
Retail	<b>0.06</b> (10)	0.81 (10)
PumSB-star	2208.27 (7)	<b>20.88</b> (10)
DNA	- (0)	<b>505.75</b> (9)
Kosarek	48.74 (10)	<b>3.41</b> (10)
MSWeb	1543.49 (10)	<b>1.28</b> (10)
Book	- (0)	<b>46.50</b> (10)
EachMovie	- (0)	<b>1216.89</b> (8)
WebKB	- (0)	<b>575.68</b> (10)
Reuters-52	- (0)	<b>120.58</b> (10)
20 NewsGrp.	- (0)	<b>504.52</b> (9)
BBC	- (0)	<b>2757.18</b> (3)
Ad	- (0)	<b>1254.37</b> (8)

# Tractable Computation of Expected Kernels

- How to compute the expected kernel given two distributions  $\mathbf{p}$ ,  $\mathbf{q}$ ?

$$\mathbb{E}_{\mathbf{x} \sim \mathbf{p}, \mathbf{x}' \sim \mathbf{q}}[\mathbf{k}(\mathbf{x}, \mathbf{x}')] ]$$

- Circuit representation for kernel functions, e.g.,  $\mathbf{k}(\mathbf{x}, \mathbf{x}') = \exp(-\sum_{i=1}^4 |X_i - X'_i|^2)$



# Tractable Computation of Expected Kernels: Applications

- Reasoning about support vector regression (SVR) with missing features

$$\mathbb{E}_{\mathbf{x}_m \sim \mathbf{p}(\mathbf{X}_m | \mathbf{x}_o)} \left[ \underbrace{\sum_{i=1}^m w_i \mathbf{k}(\mathbf{x}_i, \mathbf{x}) + b}_{\text{SVR model}} \right]$$

missing features

- Collapsed Black-box Importance Sampling: minimize kernelized Stein discrepancy

importance weights  $\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \left\{ \mathbf{w}^\top \mathbf{K}_{p,s} \mathbf{w} \mid \sum_{i=1}^n w_i = 1, w_i \geq 0 \right\}$

↓

expected kernel matrix

# Model-Based Algorithmic Fairness: FairPC

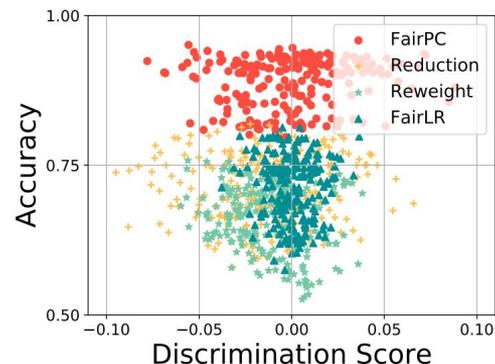
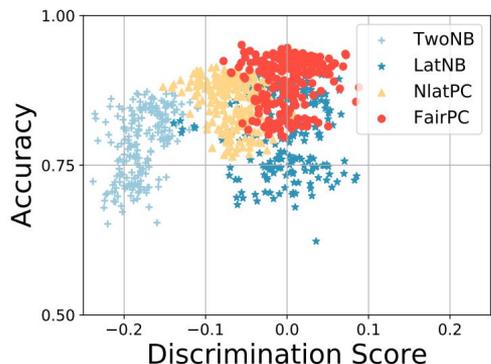
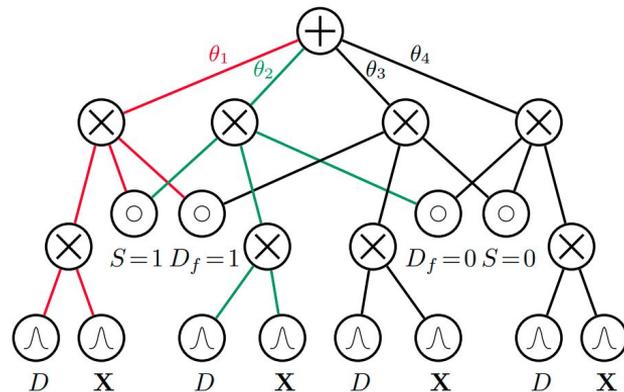
Learn classifier given

- features  $S$  and  $X$
- training labels/decisions  $D$

Group fairness by demographic parity:

*Fair decision  $D_f$  should be independent of the sensitive attribute  $S$*

Discover the **latent fair decision  $D_f$**  by learning a PC.

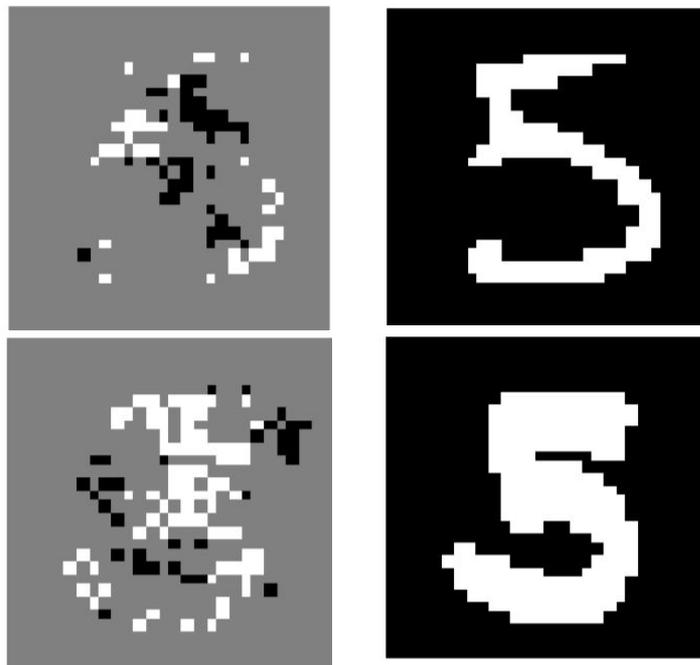


# Probabilistic Sufficient Explanations

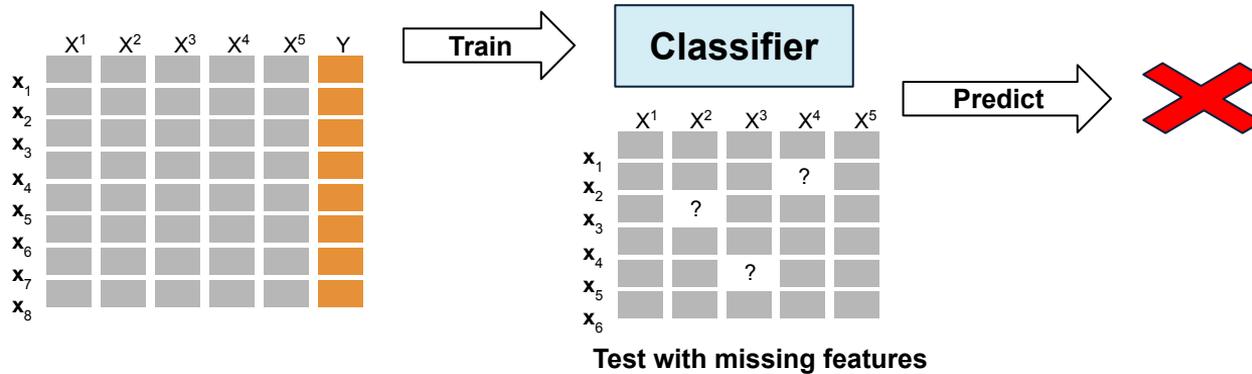
Goal: explain an instance of classification (a specific prediction)

Explanation is a subset of features, s.t.

1. The explanation is “probabilistically sufficient”  
*Under the feature distribution, given the explanation, the classifier is likely to make the observed prediction.*
2. It is minimal and “simple”

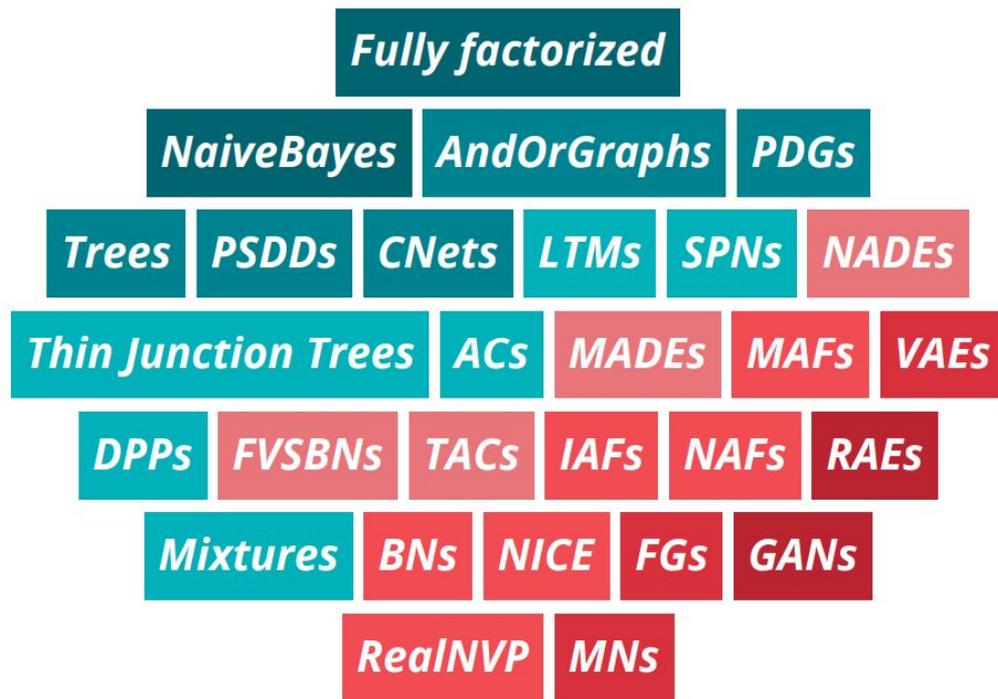


# Prediction with Missing Features

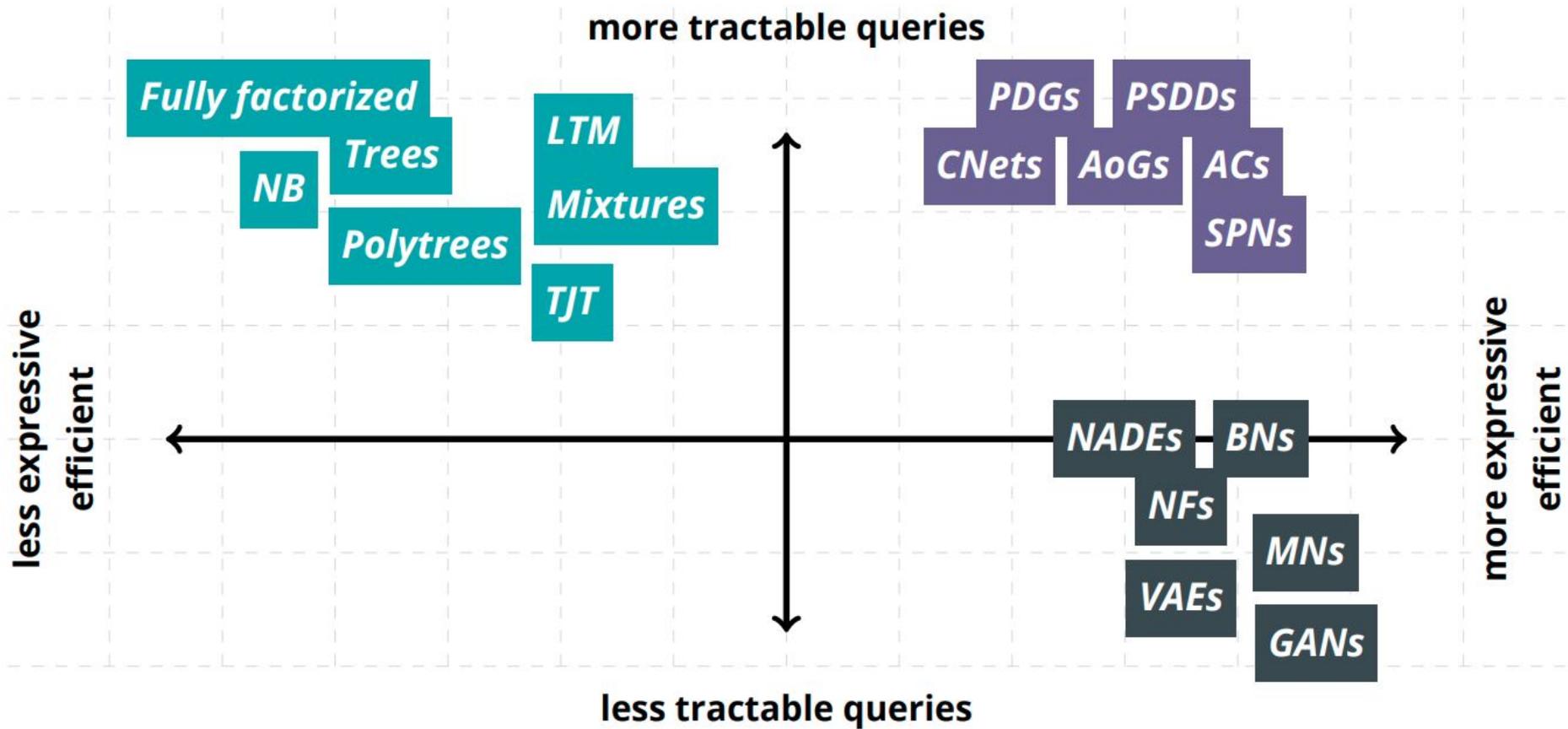


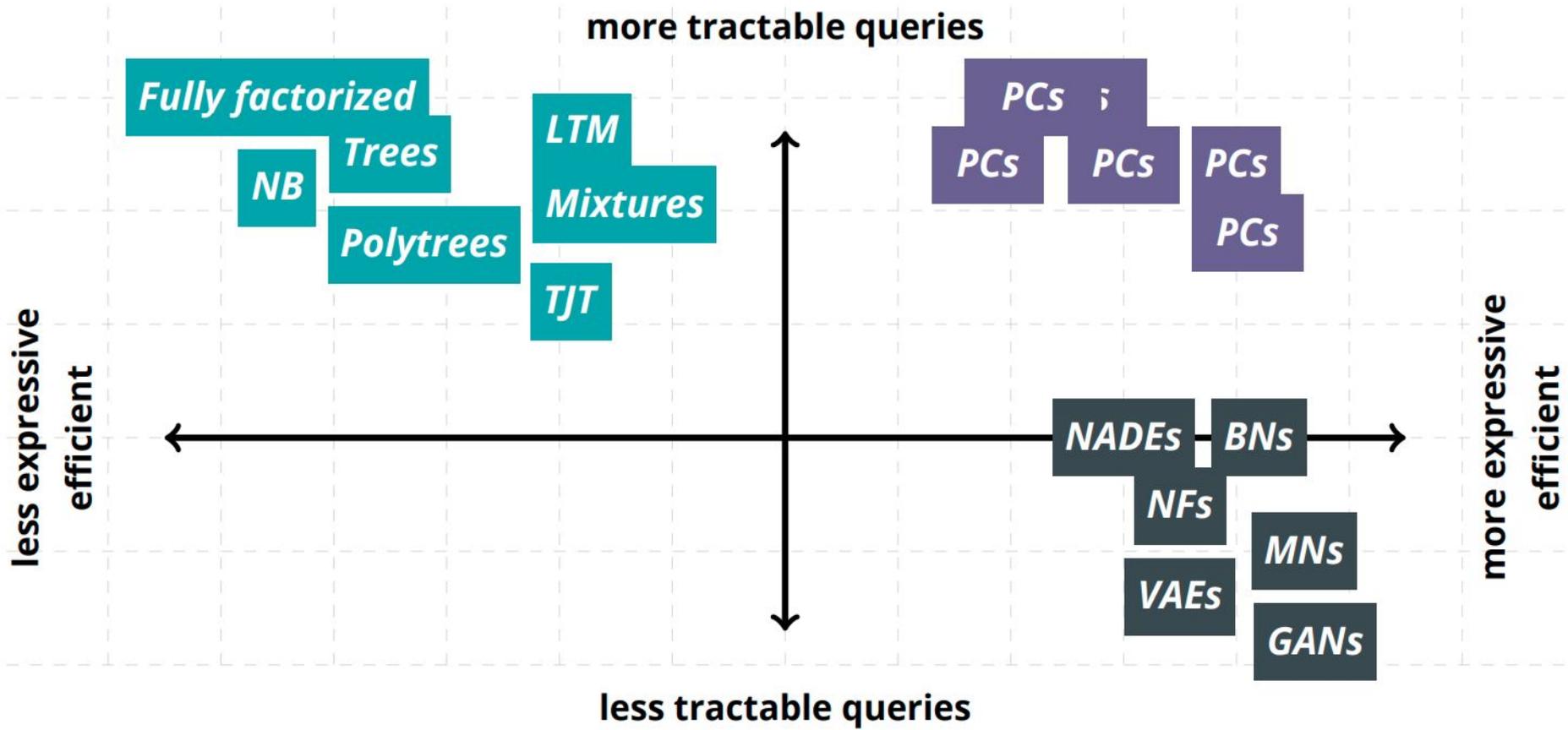
See work on

- Expected predictions / conformant learning [Khosravi et al.]
- Generative forests [Correia et al.]



**tractability is a spectrum**



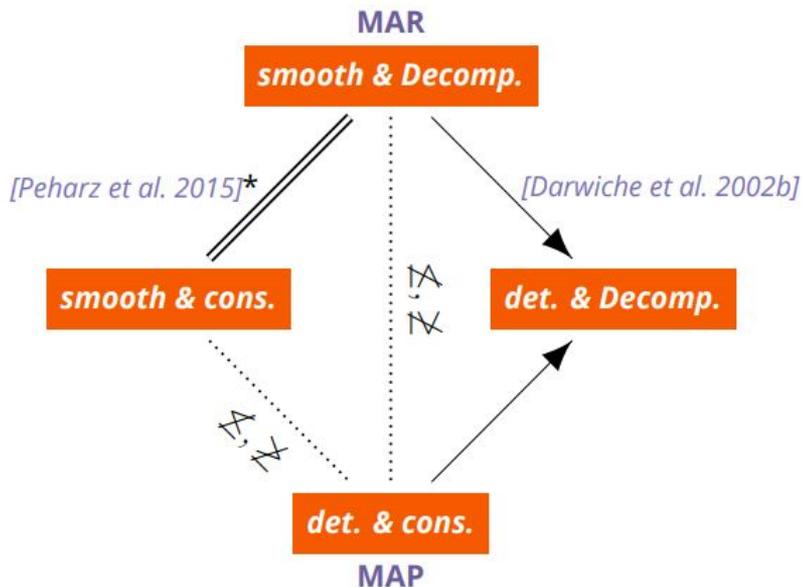


# Outline



1. What are tractable probabilistic circuits?
2. Are these models any good?
3. How far can we push tractable inference?
4. **What is their expressive power?**

# Expressive efficiency of circuits



→ : strictly more succinct  
== : equally succinct

- Neither smooth & decomposable nor deterministic & consistent circuits are more succinct than the other!  
⇒ Choose tractable circuit family based on your query
- More theoretical questions remaining  
⇒ "Complete the map"

ask YooJung Choi and Stefan Mengel

Probabilistic circuits seem awfully general.

*Are all tractable probabilistic models  
probabilistic circuits?*



# Enter: Determinantal Point Processes (DPPs)

DPPs are models where probabilities are specified by (sub)determinants

$$L = \begin{bmatrix} 1 & 0.9 & 0.8 & 0 \\ 0.9 & 0.97 & 0.96 & 0 \\ 0.8 & 0.96 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

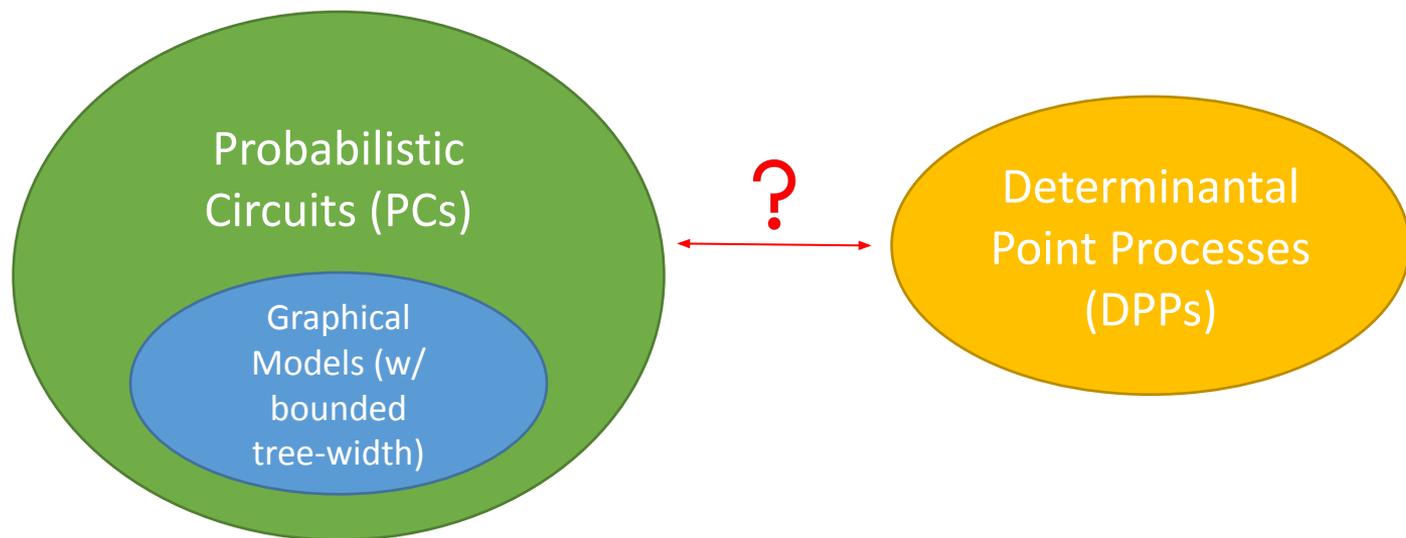
Tractable likelihoods and marginals

Global Negative Dependence

Diversity in recommendation systems

$$\Pr_L(X_1 = 1, X_2 = 0, X_3 = 1, X_4 = 0) = \frac{1}{\det(L + I)} \det(L_{\{1,2\}})$$

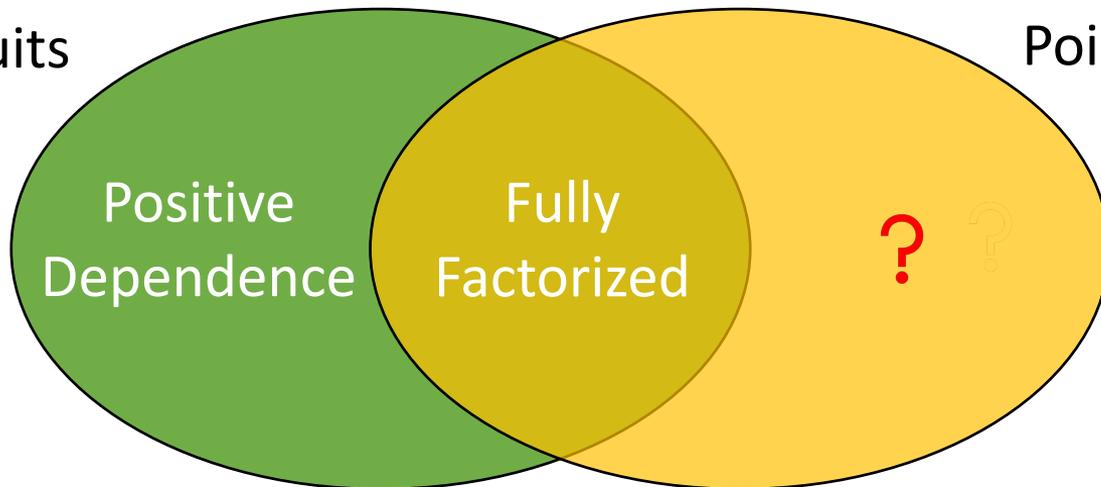
# *Are all tractable probabilistic models probabilistic circuits?*



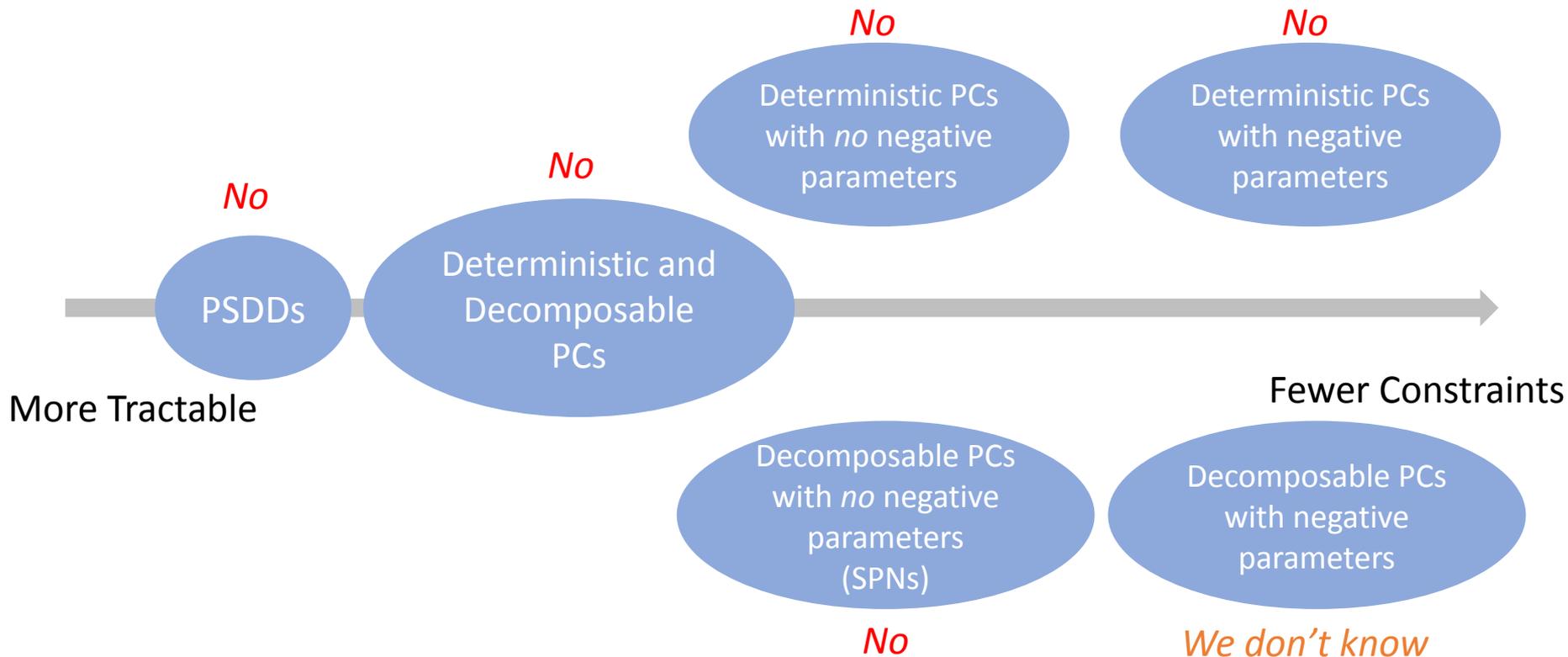
# Relationship between PCs and DPPs

Probabilistic  
Circuits

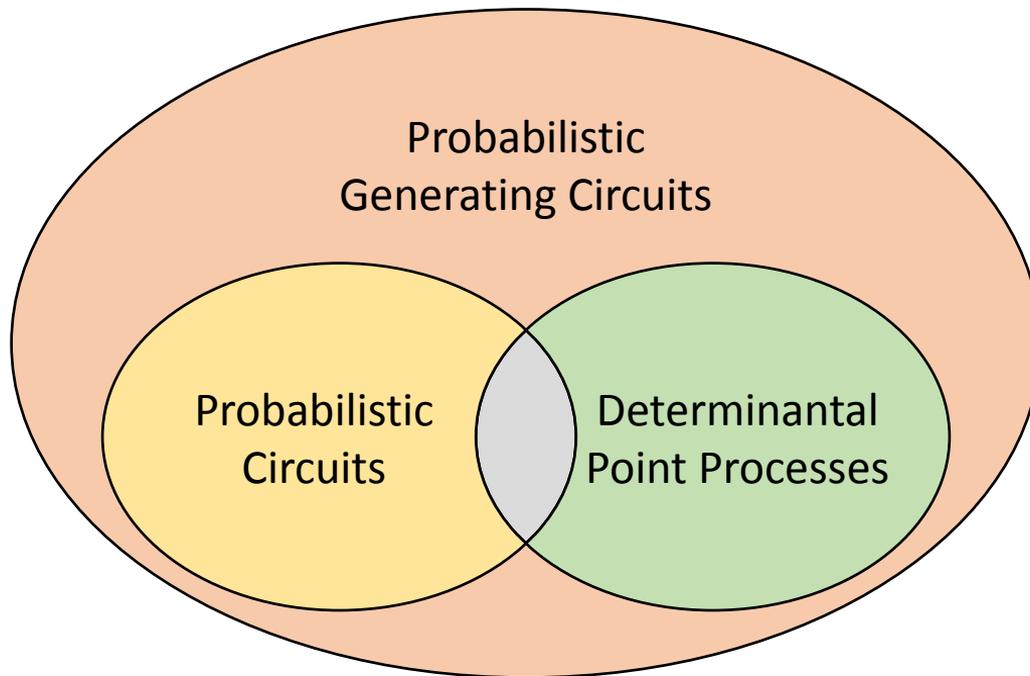
Determinantal  
Point Processes



# We cannot tractably represent DPPs with subclasses of PCs



# Probabilistic Generating Circuits



A Tractable Unifying Framework for PCs and DPPs

# Probability Generating Functions

$X_1$	$X_2$	$X_3$	$\Pr_\beta$
0	0	0	0.02
0	0	1	0.08
0	1	0	0.12
0	1	1	0.48
1	0	0	0.02
1	0	1	0.08
1	1	0	0.04
1	1	1	0.16



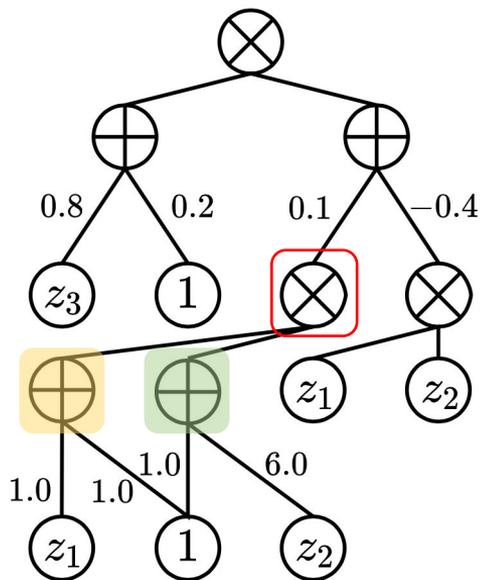
$$g_\beta = 0.16z_1z_2z_3 + 0.04z_1z_2 + 0.08z_1z_3 + 0.02z_1 + 0.48z_2z_3 + 0.12z_2 + 0.08z_3 + 0.02.$$



$$g_\beta = (0.1(z_1 + 1))(6z_2 + 1) - 0.4z_1z_2)(0.8z_3 + 0.2)$$

# Probabilistic Generating Circuits (PGCs)

$$g_{\beta} = (0.1(z_1 + 1)(6z_2 + 1) - 0.4z_1z_2)(0.8z_3 + 0.2)$$



1. Sum nodes  $\oplus$  with weighted edges to children.
2. Product nodes  $\otimes$  with unweighted edges to children.
3. Leaf nodes:  $z_i$  or constant.

# DPPs as PGCs

The generating polynomial for a DPP with kernel  $L$  is given by:

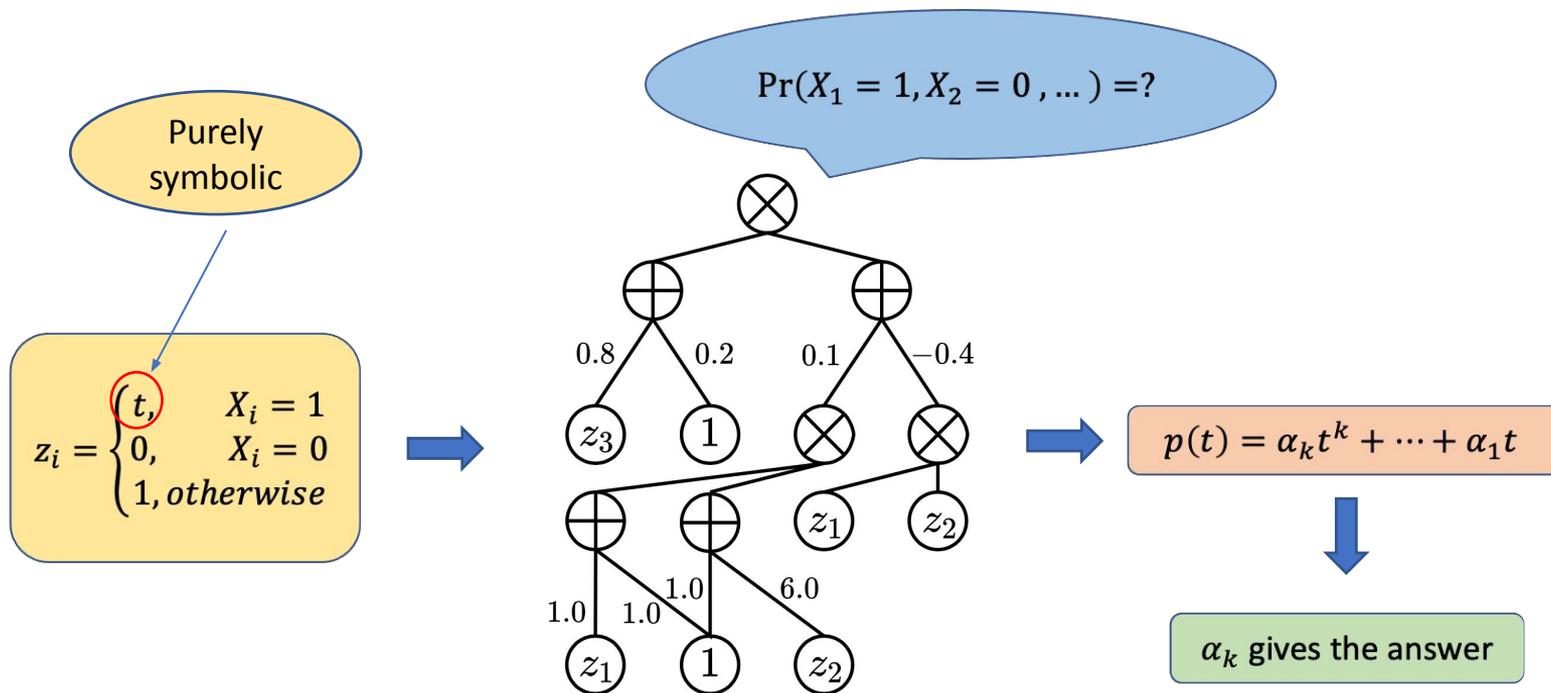
$$g_L = \frac{1}{\det(L + I)} \det(I + L \text{diag}(z_1, \dots, z_n)).$$

Constant

Division-free determinant algorithm  
(Samuelson-Berkowitz algorithm)

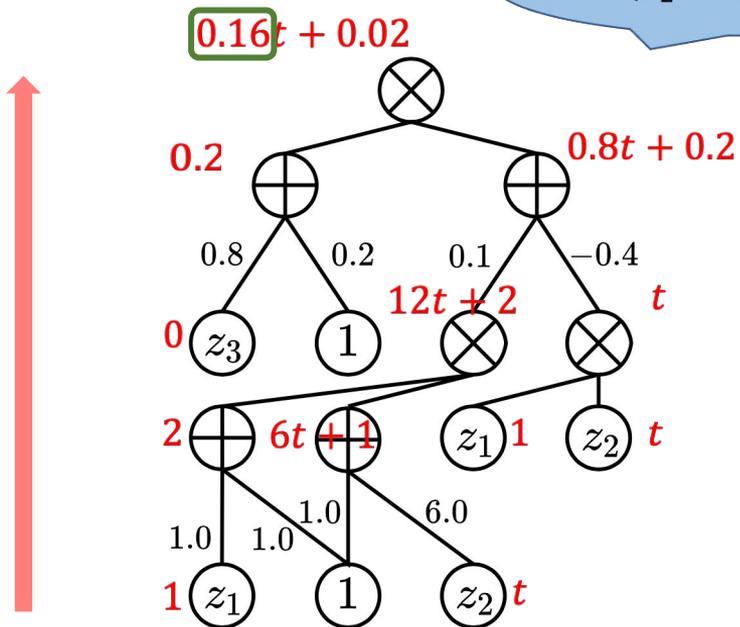
$g_L$  can be represented as a PGC of size  $O(n^4)$

# PGCs Support Tractable Likelihoods/Marginals



# Example

$\Pr(X_2 = 1, X_3 = 0) = ?$



$X_1$	$X_2$	$X_3$	$\Pr_\beta$
0	0	0	0.02
0	0	1	0.08
0	1	0	0.12
0	1	1	0.48
1	0	0	0.02
1	0	1	0.08
1	1	0	0.04
1	1	1	0.16

# Experiment Results: Amazon Baby Registries

	DPP	Strudel	EiNet	MT	SimplePGC
apparel	-9.88	-9.51	-9.24	-9.31	<b>-9.10</b> <sup>*†°</sup>
bath	-8.55	-8.38	-8.49	-8.53	<b>-8.29</b> <sup>*†°</sup>
bedding	-8.65	-8.50	-8.55	-8.59	<b>-8.41</b> <sup>*†°</sup>
carseats	-4.74	-4.79	-4.72	-4.76	<b>-4.64</b> <sup>*†°</sup>
diaper	-10.61	-9.90	-9.86	-9.93	<b>-9.72</b> <sup>*†°</sup>
feeding	-11.86	-11.42	-11.27	-11.30	<b>-11.17</b> <sup>*†°</sup>
furniture	-4.38	-4.39	-4.38	-4.43	<b>-4.34</b> <sup>*†°</sup>
gear	-9.14	-9.15	-9.18	-9.23	<b>-9.04</b> <sup>*†°</sup>
gifts	-3.51	<b>-3.39</b>	-3.42	-3.48	-3.47 <sup>°</sup>
health	-7.40	-7.37	-7.47	-7.49	<b>-7.24</b> <sup>*†°</sup>
media	-8.36	<b>-7.62</b>	-7.82	-7.93	-7.69 <sup>†°</sup>
moms	-3.55	-3.52	<b>-3.48</b>	-3.54	-3.53 <sup>°</sup>
safety	-4.28	-4.43	-4.39	-4.36	<b>-4.28</b> <sup>*†°</sup>
strollers	-5.30	-5.07	-5.07	-5.14	<b>-5.00</b> <sup>*†°</sup>
toys	-8.05	<b>-7.61</b>	-7.84	-7.88	-7.62 <sup>†°</sup>

SimplePGC achieves SOTA  
result on 11/15 datasets

# Conclusion



1. What are tractable probabilistic circuits?
2. Are these models any good?
3. How far can we push tractable inference?
4. What is their expressive power?

# Learn more about probabilistic circuits?



## Tutorial (3h)

**Probabilistic Circuits**

**Inference**  
**Representations**  
**Learning**  
**Theory**

**Antonio Vergari**  
University of California, Los Angeles

**Robert Peharz**  
TU Eindhoven

**YooJung Choi**  
University of California, Los Angeles

**Guy Van den Broeck**  
University of California, Los Angeles

September 14th, 2020 - Ghent, Belgium - ECML-PKDD 2020

<https://youtu.be/2RAG5-L9R70>

## Overview Paper (80p)

### Probabilistic Circuits: A Unifying Framework for Tractable Probabilistic Models\*

YooJung Choi

Antonio Vergari

Guy Van den Broeck

*Computer Science Department*

*University of California*

*Los Angeles, CA, USA*

#### Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Probabilistic Inference: Models, Queries, and Tractability</b>	<b>4</b>
2.1	Probabilistic Models . . . . .	5
2.2	Probabilistic Queries . . . . .	6
2.3	Tractable Probabilistic Inference . . . . .	8
2.4	Properties of Tractable Probabilistic Models . . . . .	9

<http://starai.cs.ucla.edu/papers/ProbCirc20.pdf>