

Tractable Probabilistic Modeling with Probabilistic Circuits

antonio vergari (he/him)



@tetraduzione

20th June 2024 - Nordic Prob AI Summer School - Copenhagen

april

*april is
probably a
recursive
identifier of a
lab*

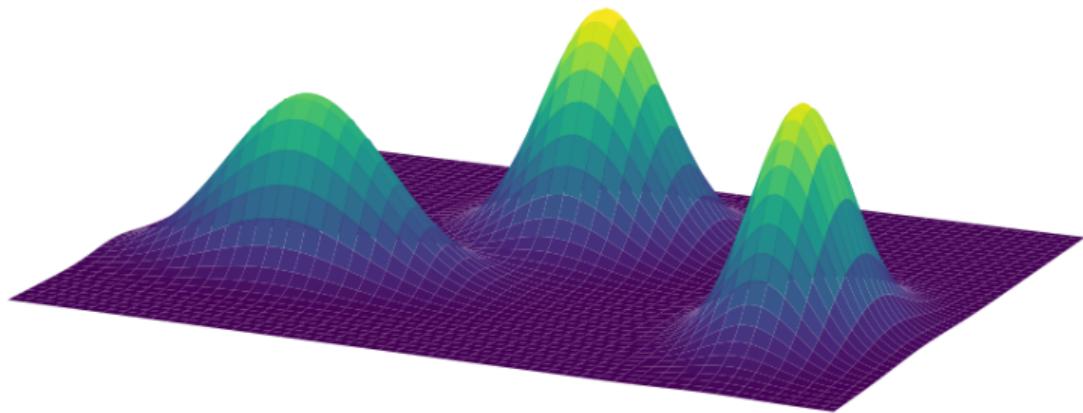
april

*about
probabilities
integrals &
logic*

deep generative models

+

*flexible and reliable
logic & probabilistic reasoning?*



a love letter to mixture models...

Reasoning about ML models



q₁

*"What is the probability of a treatment for a patient with **unavailable records**?"*



q₂

*"How **fair** is the prediction is a certain protected attribute changes?"*



q₃

*"Can we certify no **adversarial examples** exist?"*

Reasoning about ML models



q₁ $\int p(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{X}_m$
(missing values)

q₂ $\mathbb{E}_{\mathbf{x}_c \sim p(\mathbf{X}_c | X_s=0)} [f_0(\mathbf{x}_c)] - \mathbb{E}_{\mathbf{x}_c \sim p(\mathbf{X}_c | X_s=1)} [f_1(\mathbf{x}_c)]$
(fairness)

q₃ $\mathbb{E}_{\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_D)} [f(\mathbf{x} + \mathbf{e})]$
(adversarial robust.)

...in the language of probabilities

more complex reasoning



neuro-symbolic AI



probabilistic programming



*computing uncertainties
(Bayesian inference)*

...and more application scenarios

Reasoning about ML models



q₁ $\int p(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{X}_m$
(missing values)

q₂ $\mathbb{E}_{\mathbf{x}_c \sim p(\mathbf{X}_c | X_s=0)} [f_0(\mathbf{x}_c)] - \mathbb{E}_{\mathbf{x}_c \sim p(\mathbf{X}_c | X_s=1)} [f_1(\mathbf{x}_c)]$
(fairness)

q₃ $\mathbb{E}_{\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_D)} [f(\mathbf{x} + \mathbf{e})]$
(adversarial robust.)

hard to compute in general!

Reasoning about ML models



q₁ $\int p(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{X}_m$
(missing values)

q₂ $\mathbb{E}_{\mathbf{x}_c \sim p(\mathbf{X}_c | X_s=0)} [f_0(\mathbf{x}_c)] - \mathbb{E}_{\mathbf{x}_c \sim p(\mathbf{X}_c | X_s=1)} [f_1(\mathbf{x}_c)]$
(fairness)

q₃ $\mathbb{E}_{\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_D)} [f(\mathbf{x} + \mathbf{e})]$
(adversarial robust.)

it is crucial we compute them exactly and in polytime!

Reasoning about ML models



q₁ $\int p(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{X}_m$
(missing values)

q₂ $\mathbb{E}_{\mathbf{x}_c \sim p(\mathbf{X}_c | X_s=0)} [f_0(\mathbf{x}_c)] - \mathbb{E}_{\mathbf{x}_c \sim p(\mathbf{X}_c | X_s=1)} [f_1(\mathbf{x}_c)]$
(fairness)

q₃ $\mathbb{E}_{\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_D)} [f(\mathbf{x} + \mathbf{e})]$
(adversarial robust.)

it is crucial we compute them tractably!

Tractable Probabilistic Inference

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

⇒ **model-centric** definition...

Tractable Probabilistic Inference

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

\Rightarrow **model-centric** definition...

\Rightarrow ...and **query-centric**: Tractability is not a universal property!

Tractable Probabilistic Inference

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

⇒ **model-centric** definition...

⇒ ...and **query-centric**: Tractability is not a universal property!

⇒ often poly will in fact be **linear**!

Tractable Probabilistic Inference

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

⇒ **model-centric** definition...

⇒ ...and **query-centric**: Tractability is not a universal property!

⇒ often poly will in fact be **linear**!

⇒ Note: if $|\mathbf{m}| \in O(\text{poly}(|\mathbf{X}|))$, then query time is $O(\text{poly}(|\mathbf{X}|))$.

Tractable Probabilistic Inference

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

⇒ **model-centric** definition...

⇒ ...and **query-centric**: Tractability is not a universal property!

⇒ often poly will in fact be **linear**!

⇒ Note: if $|\mathbf{m}| \in O(\text{poly}(|\mathbf{X}|))$, then query time is $O(\text{poly}(|\mathbf{X}|))$.

⇒ Why **exactness**? Highest guarantee possible!

why tractable models?

exactness can be crucial in safety-driven applications



guarantee constraint satisfaction

[Ahmed et al. 2022]



estimation error is bounded (0)

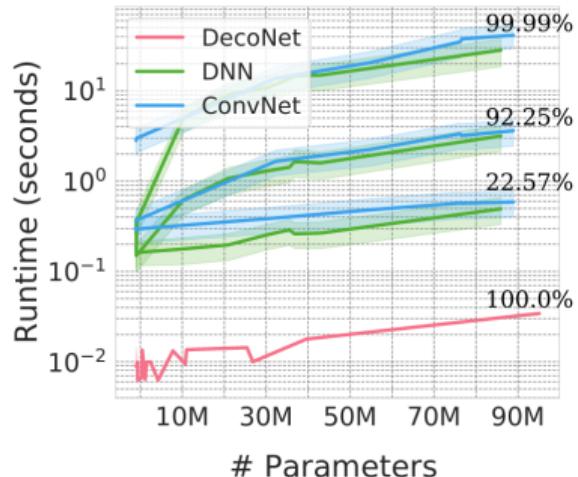
[Choi 2022]

why tractable models?

they can be much faster than intractable ones!

Method	MNIST (10,000 test images)		
	Theoretical bpd	Comp. bpd	En- & decoding time
PC (small)	1.26	1.30	53
PC (large)	1.20	1.24	168
IDF	1.90	1.96	880
BitSwap	1.27	1.31	904

[Liu, Mandt, and Broeck 2022]



[Subramani et al. 2021]

Why?

tractable inference

Why?

*we always perform
tractable inference
over an approximate model!*

$$\min_{\mathbf{q} \in \mathcal{Q}} \text{KL}(\mathbf{q} || p)$$

we pick a **tractable** variational distribution \mathbf{q}

⇒ e.g., Gaussian, GMM, HMM, flow, etc

VI

$$\min_{\mathbf{q} \in \mathcal{Q}} \text{KL}(\mathbf{q} || p)$$

we pick a **tractable** variational distribution \mathbf{q}

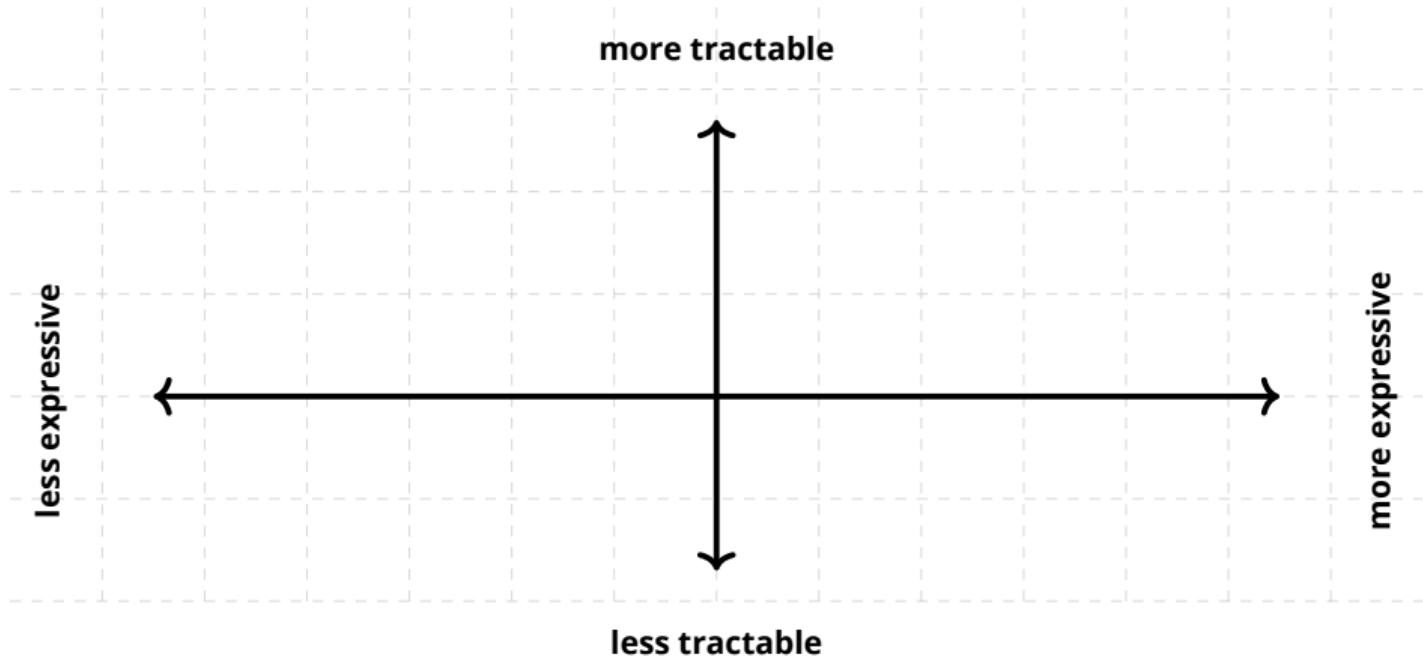
MC

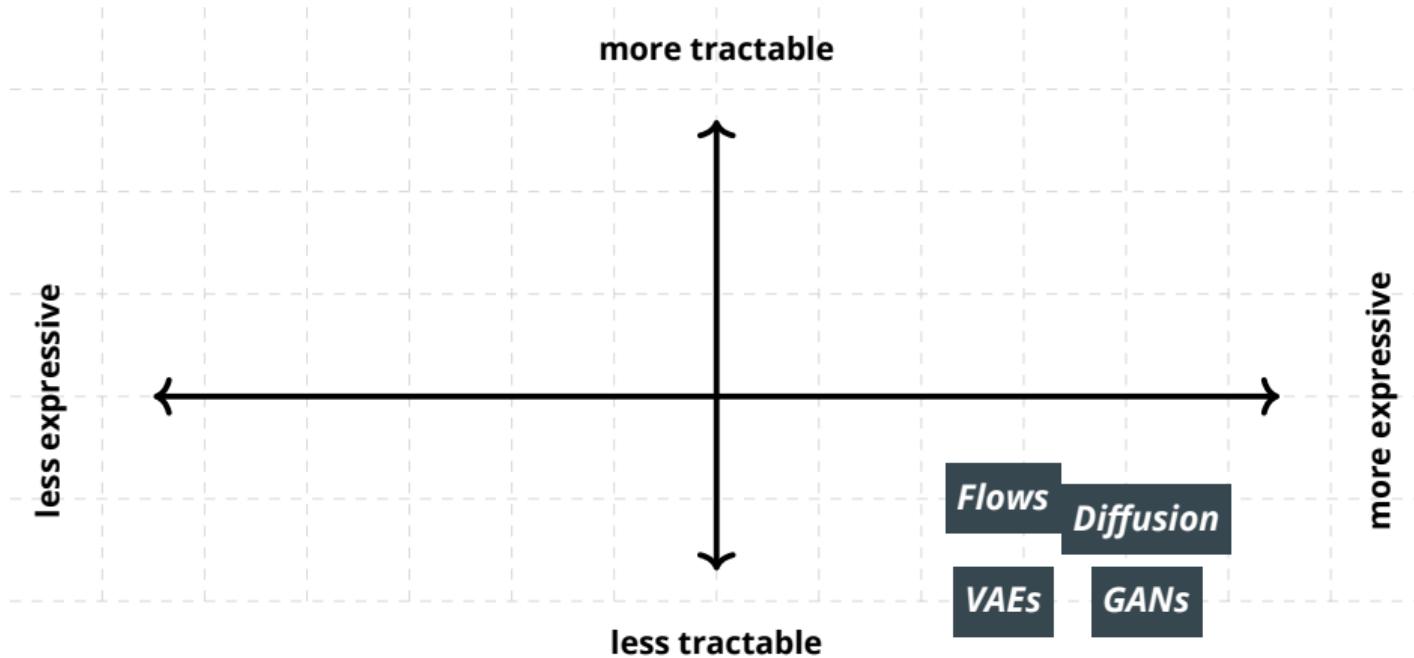
$$\mathbb{E}_{p(\mathbf{x})}[f(\mathbf{x})] \approx \frac{1}{N} \sum_{i=1}^N f(\mathbf{x}^{(i)})$$

we turn an intractable integral into a **tractable sum**

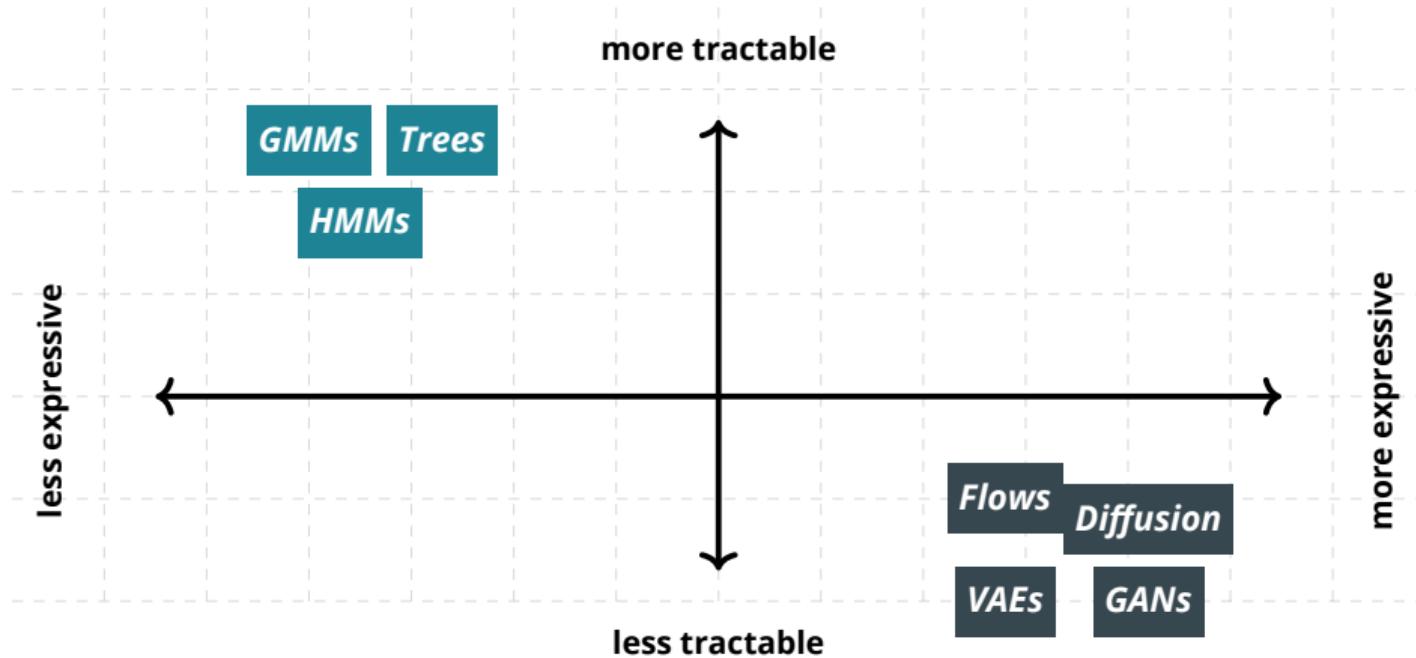
Goal

*“Can we find
a **middle ground**
between
tractability and expressiveness?”*

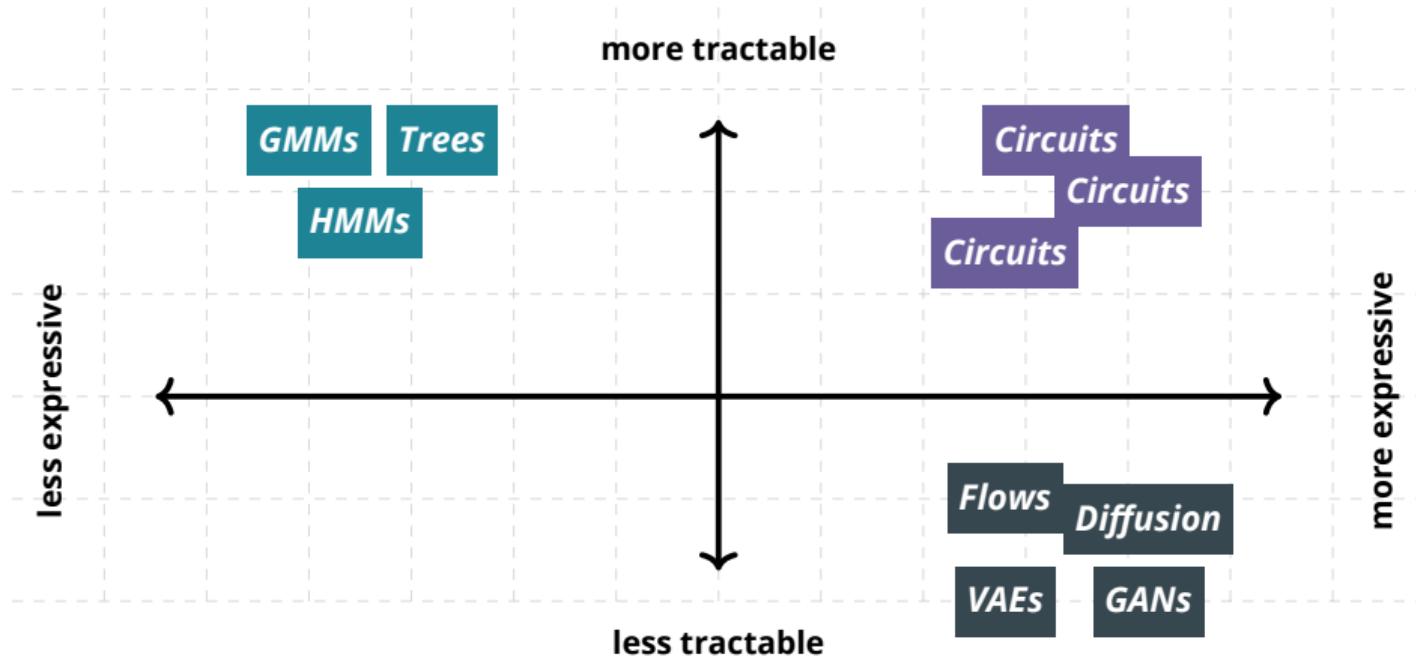




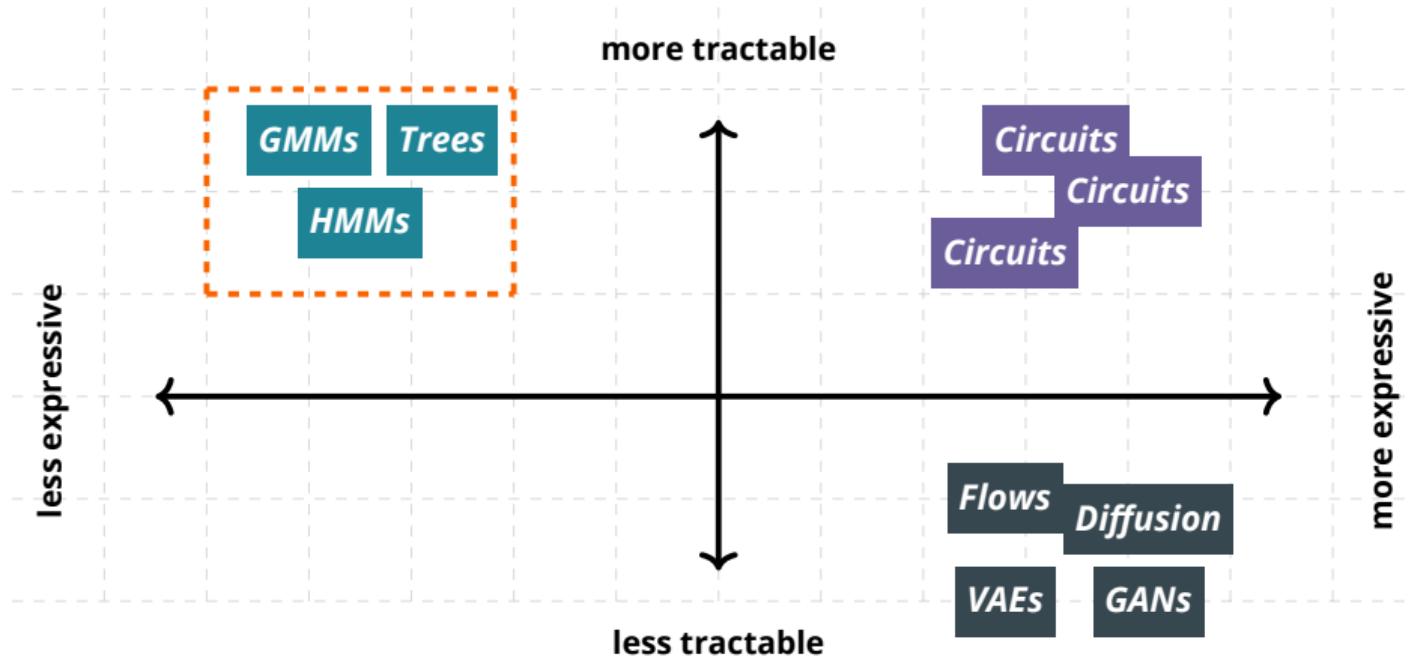
expressive models are not much tractable...



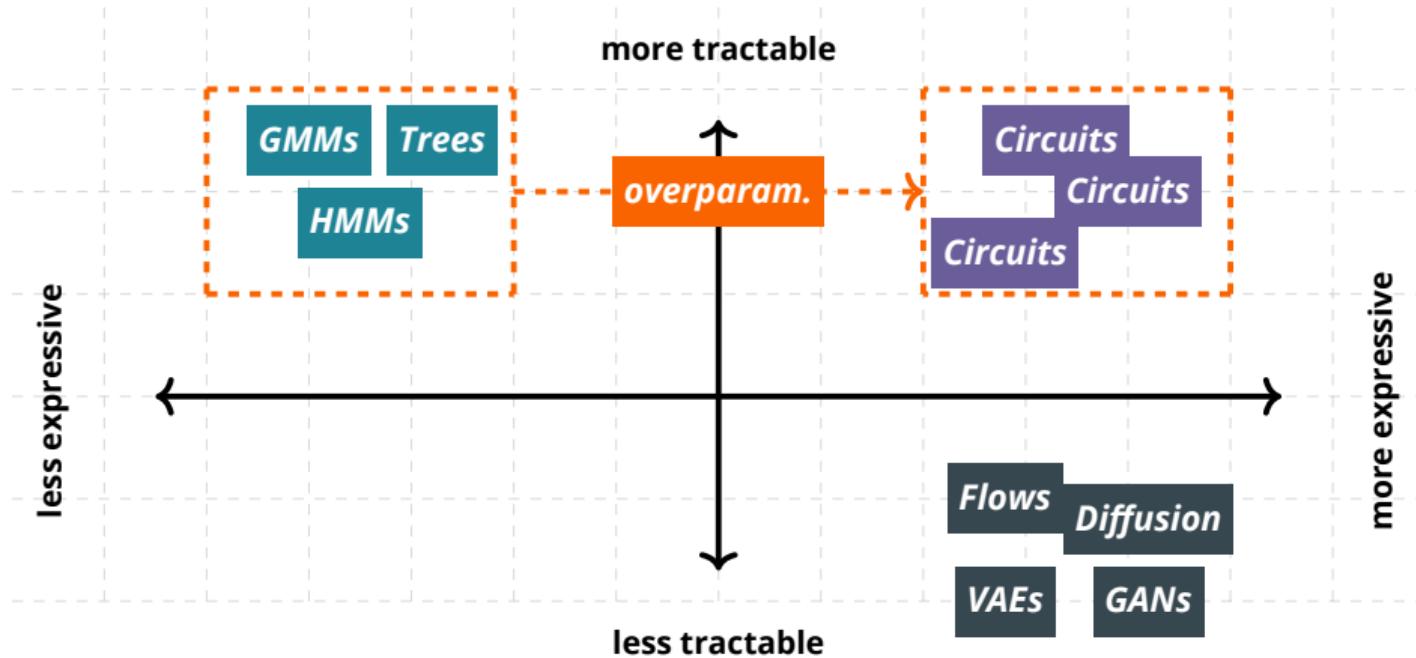
tractable models are not that expressive...



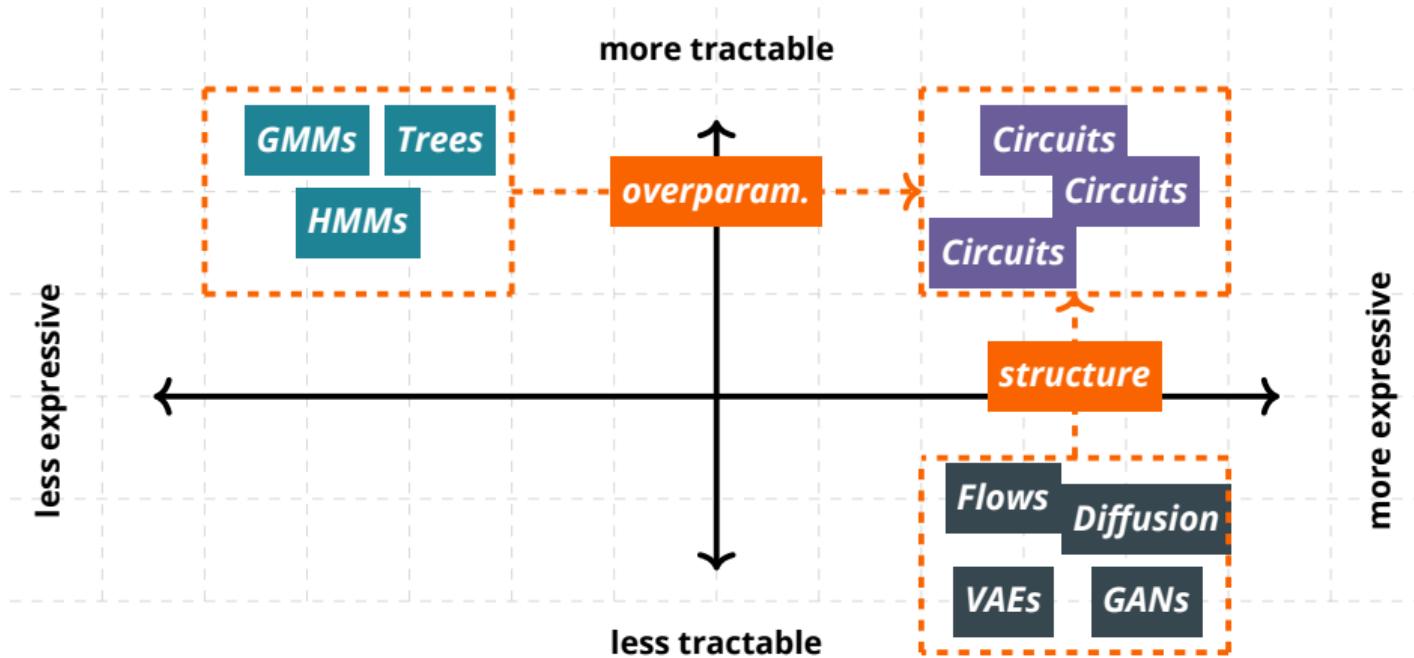
circuits can be both expressive and tractable!



start simple...



then make it more expressive!



impose structure!

Goal

*“Can we design
computational graphs
that efficiently encode inference
procedures in classical PGMs?”*

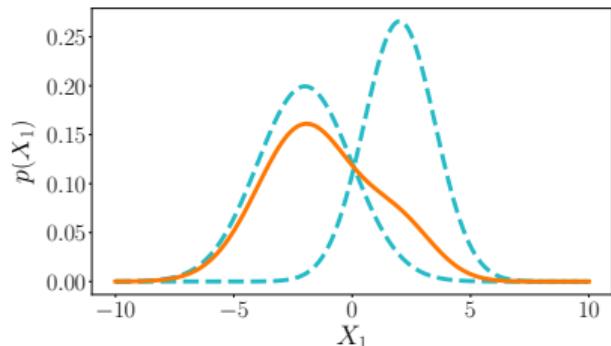
Goal

***“Can we design
computational graphs
that efficiently encode inference
procedures in classical PGMs?”***

⇒ ***yes! with circuits!***

GMMS

as computational graphs

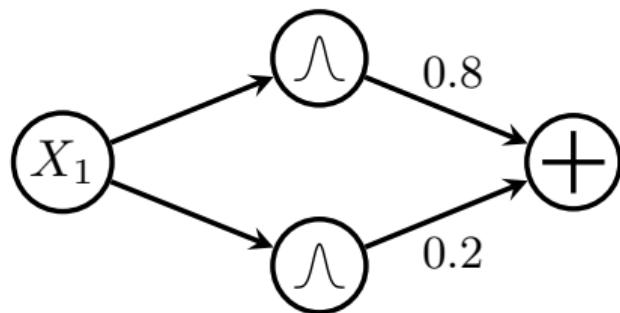


$$p(X) = w_1 \cdot p_1(X_1) + w_2 \cdot p_2(X_1)$$

⇒ translating inference to data structures...

GMMs

as computational graphs

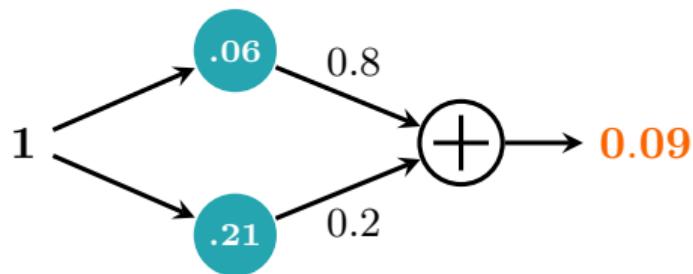


$$p(X_1) = 0.2 \cdot p_1(X_1) + 0.8 \cdot p_2(X_1)$$

⇒ ...e.g., as a weighted sum unit over Gaussian input distributions

GMMS

as computational graphs

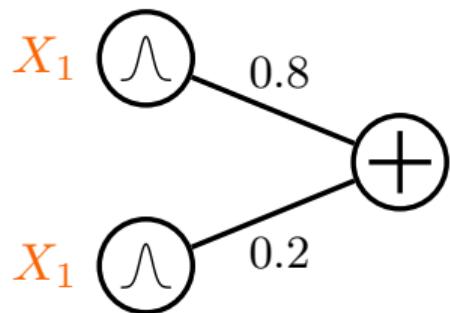


$$p(X = 1) = 0.2 \cdot p_1(X_1 = 1) + 0.8 \cdot p_2(X_1 = 1)$$

⇒ inference = feedforward evaluation

GMMs

as computational graphs

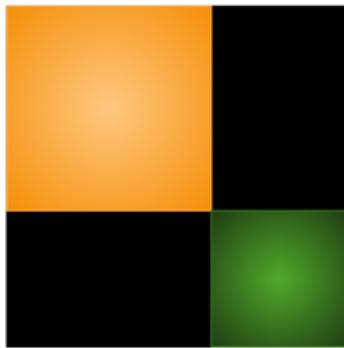
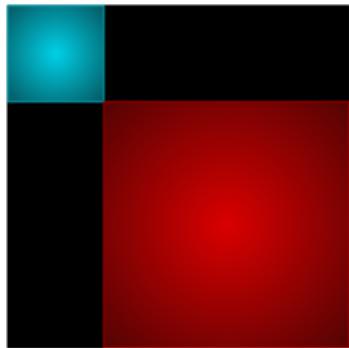


A simplified notation:

- ⇒ **scopes attached to inputs**
- ⇒ **edge directions omitted**

GMMS

as computational graphs

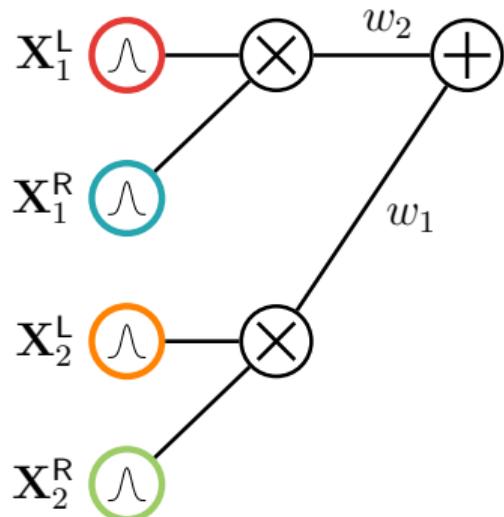


$$p(\mathbf{X}) = w_1 \cdot p_1(\mathbf{X}_1^L) \cdot p_1(\mathbf{X}_1^R) + \\ w_2 \cdot p_2(\mathbf{X}_2^L) \cdot p_2(\mathbf{X}_2^R)$$

⇒ local factorizations...

GMMs

as computational graphs



$$p(\mathbf{X}) = w_1 \cdot p_1(\mathbf{X}_1^L) \cdot p_1(\mathbf{X}_1^R) + \\ w_2 \cdot p_2(\mathbf{X}_2^L) \cdot p_2(\mathbf{X}_2^R)$$

⇒ ...are product units

Probabilistic Circuits (PCs)

A grammar for tractable computational graphs

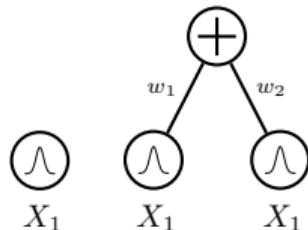
I. *A simple tractable function is a circuit*

$$\bigcirc \wedge \\ X_1$$

Probabilistic Circuits (PCs)

A grammar for tractable computational graphs

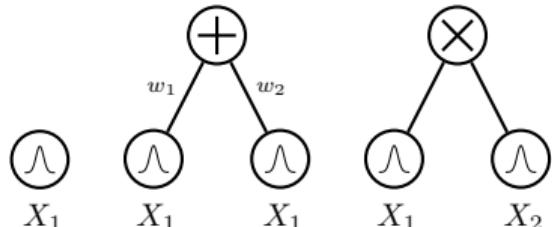
- I. A simple tractable function is a circuit
- II. A weighted combination of circuits is a circuit



Probabilistic Circuits (PCs)

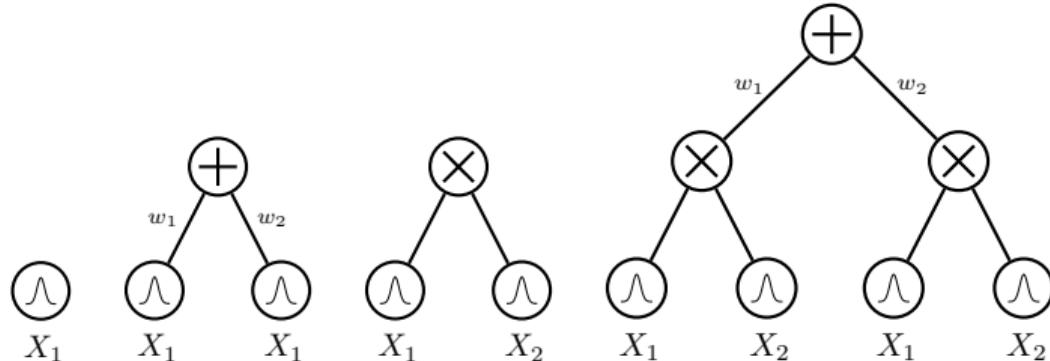
A grammar for tractable computational graphs

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



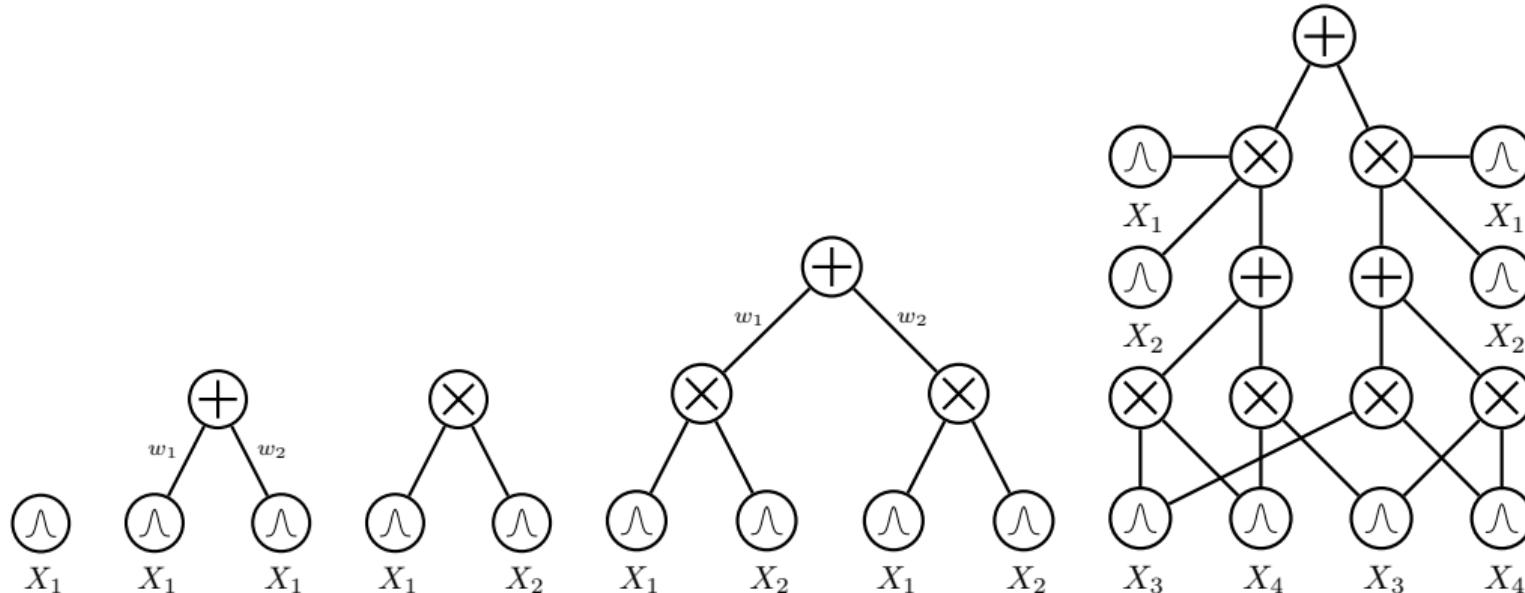
Probabilistic Circuits (PCs)

A grammar for tractable computational graphs



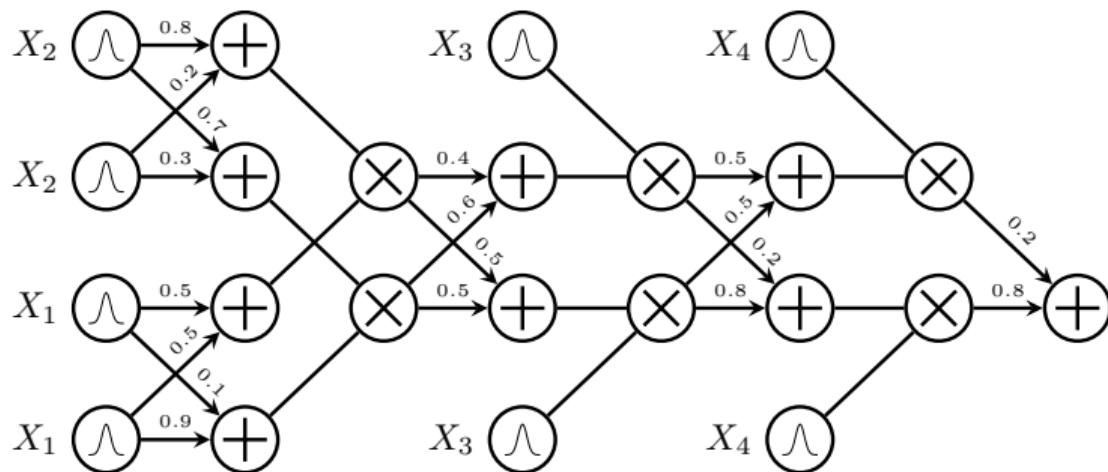
Probabilistic Circuits (PCs)

A grammar for tractable computational graphs



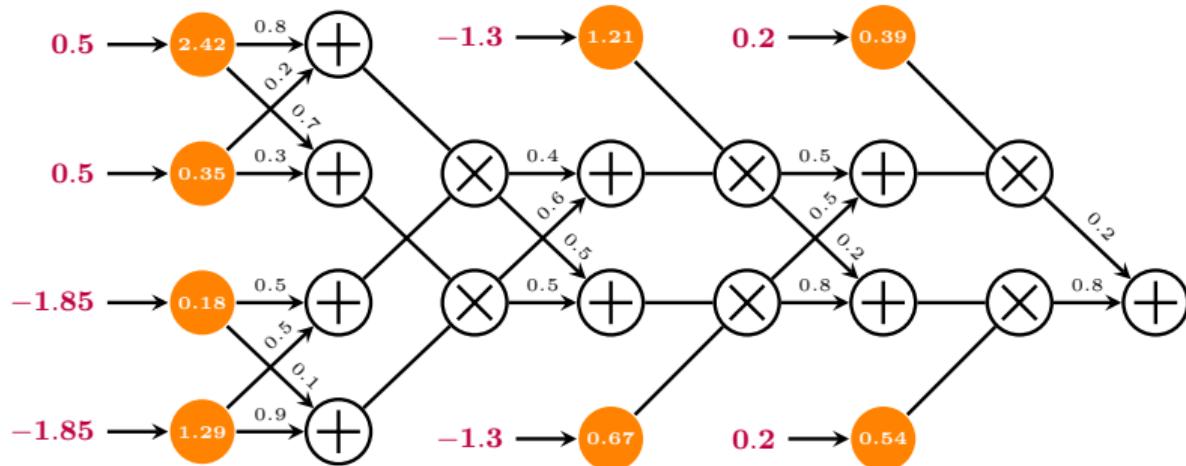
Probabilistic queries = feedforward evaluation

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



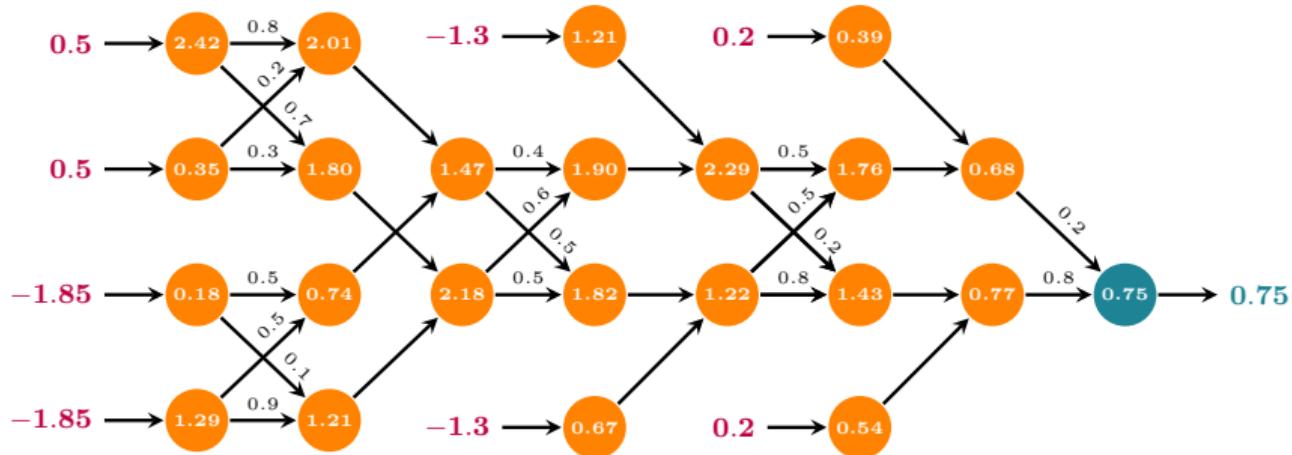
Probabilistic queries = feedforward evaluation

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



Probabilistic queries = feedforward evaluation

$$p(X_1 = -1.85, X_2 = 0.5, X_3 = -1.3, X_4 = 0.2) = 0.75$$



...why PCs?

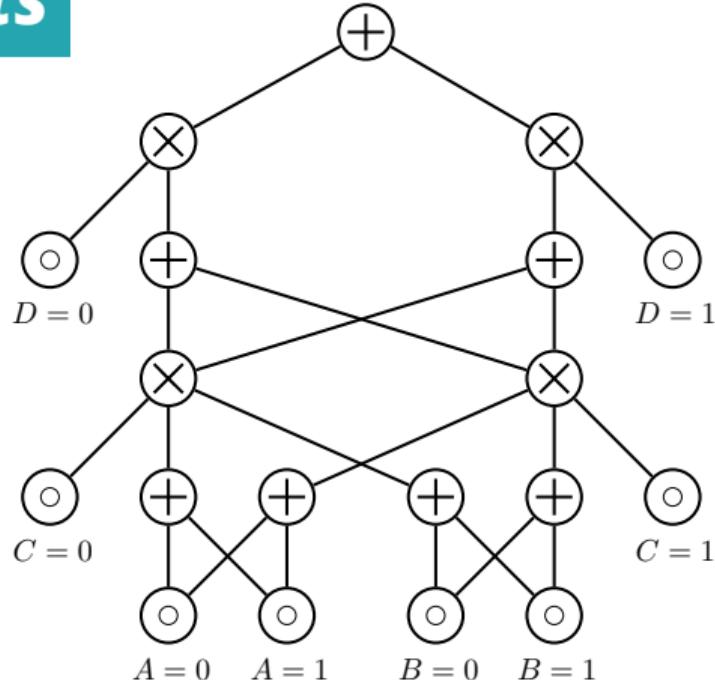
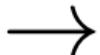
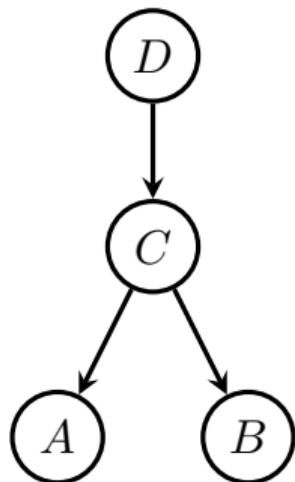
1. A grammar for tractable models

One formalism to represent many probabilistic and logical models

⇒ #HMMs #Trees #XGBoost, Tensor Networks, ...
and other PGMs...

From PGMs to circuits

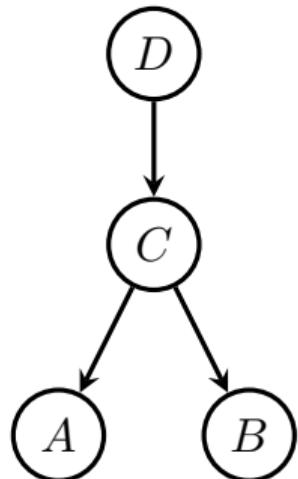
via compilation



From PGMs to circuits

via compilation

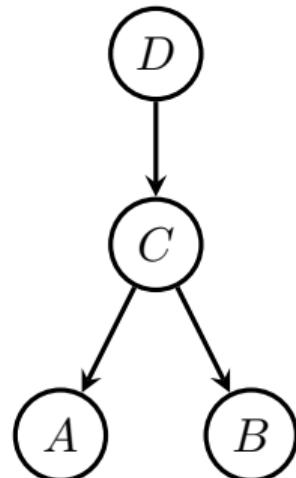
Bottom-up compilation: starting from leaves...



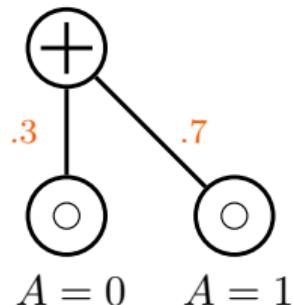
From PGMs to circuits

via compilation

...compile a leaf CPT



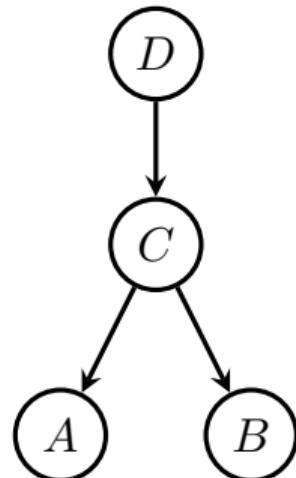
$$p(A|C = 0)$$



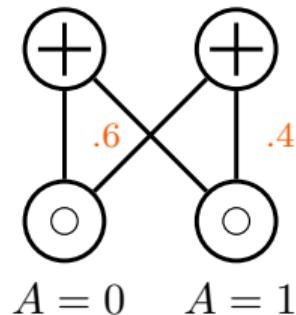
From PGMs to circuits

via compilation

...compile a leaf CPT



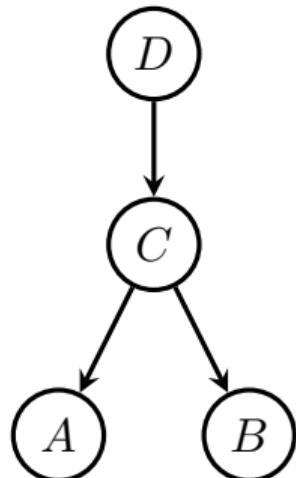
$$p(A|C = 1)$$



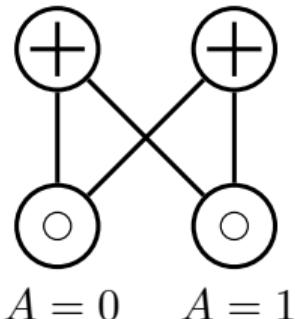
From PGMs to circuits

via compilation

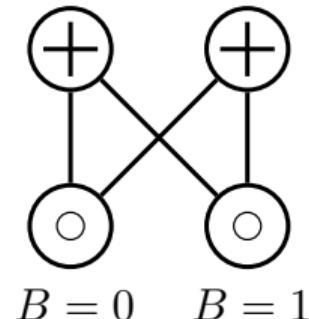
...compile a leaf CPT...for all leaves...



$$p(A|C)$$



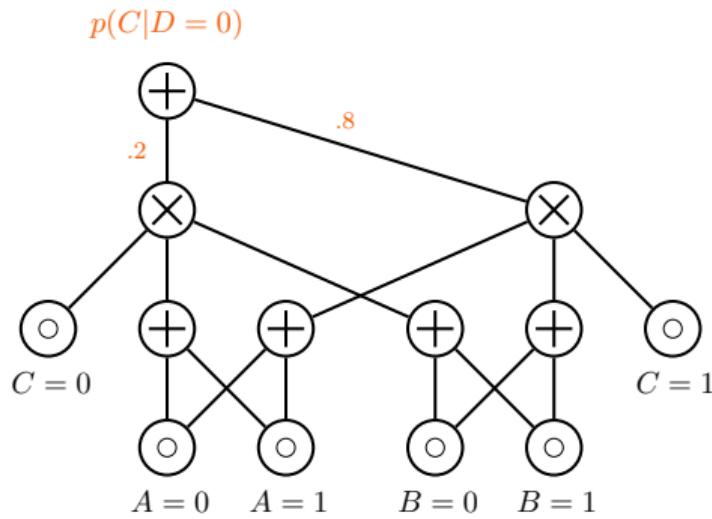
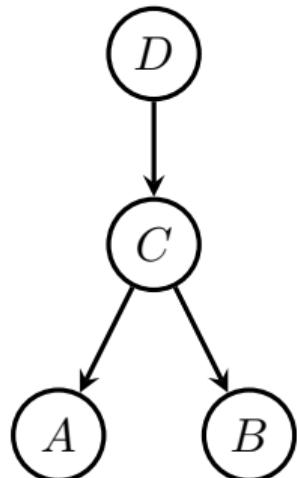
$$p(B|C)$$



From PGMs to circuits

via compilation

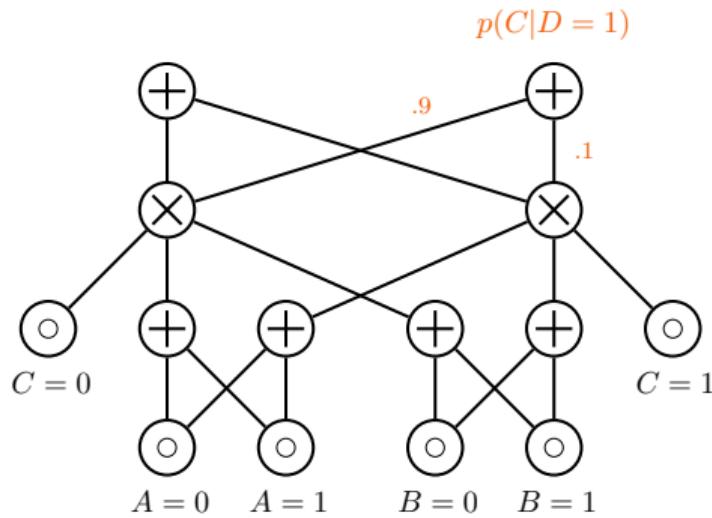
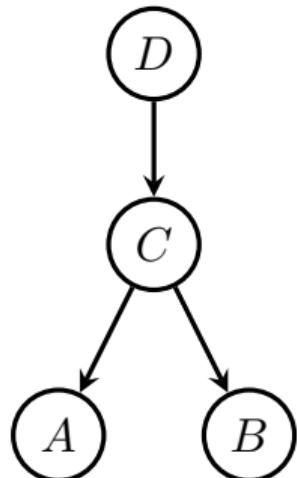
...and recurse over parents...



From PGMs to circuits

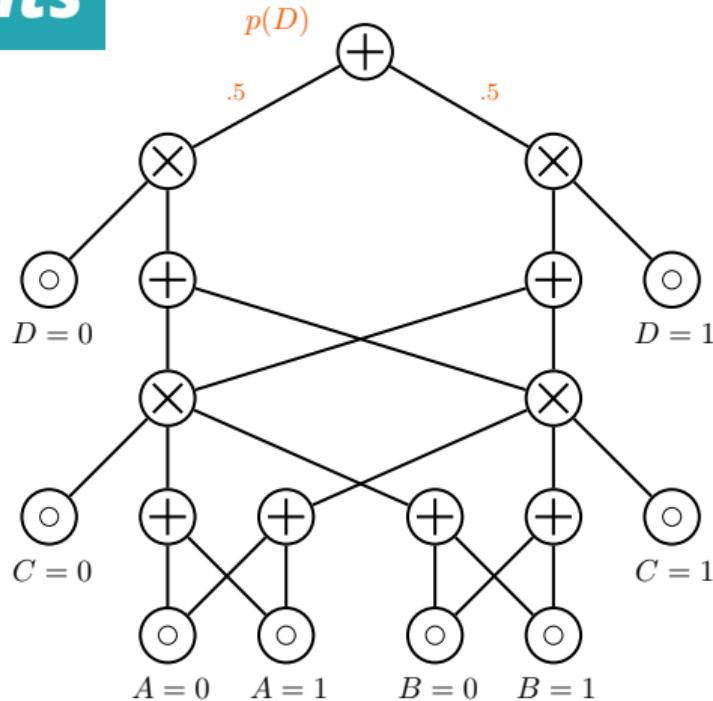
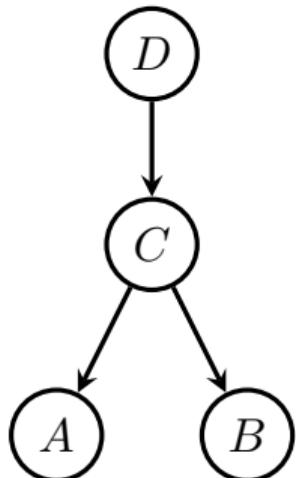
via compilation

...while reusing previously compiled nodes!...



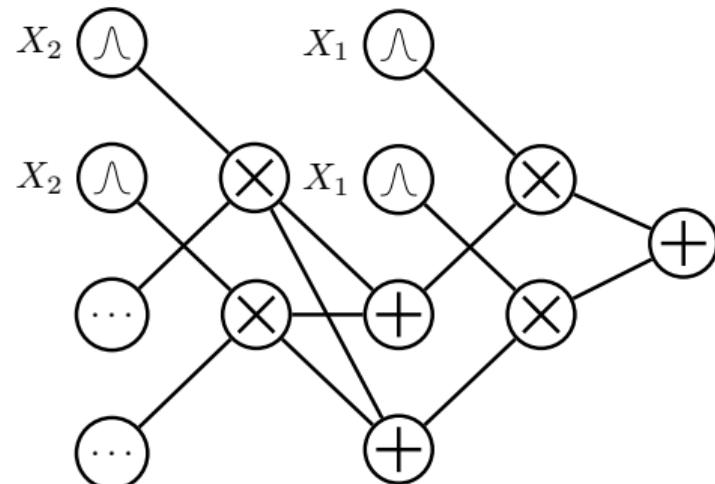
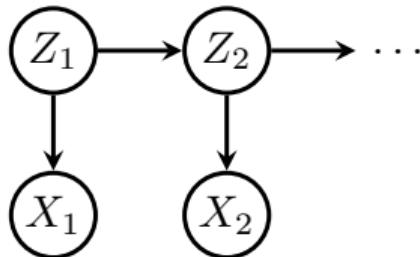
From PGMs to circuits

via compilation



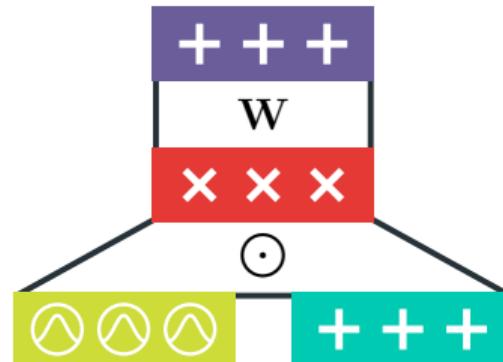
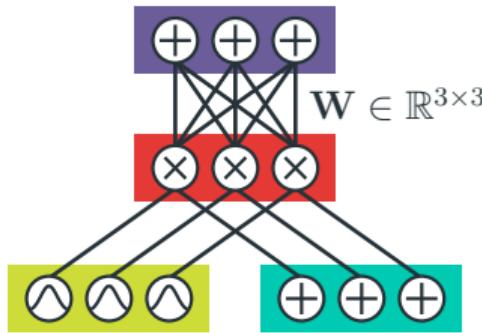
HMMs

as computational graphs



overparameterize

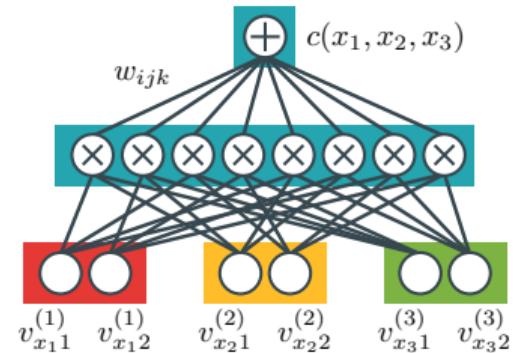
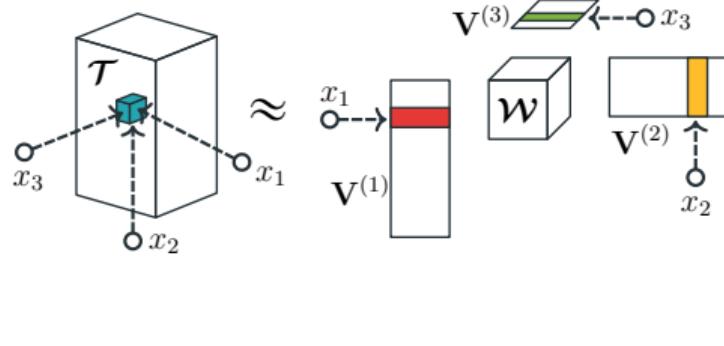
and tensorize



Mari, Vessio, and Vergari, "Unifying and Understanding Overparameterized Circuit Representations via Low-Rank Tensor Decompositions",
6th Workshop on Tractable Probabilistic Modeling, 2023

tensor factorizations

as circuits



...wait!

“Are all PGMs circuits?”

and/or

“Are all circuits PGMs?”

...not so easy!

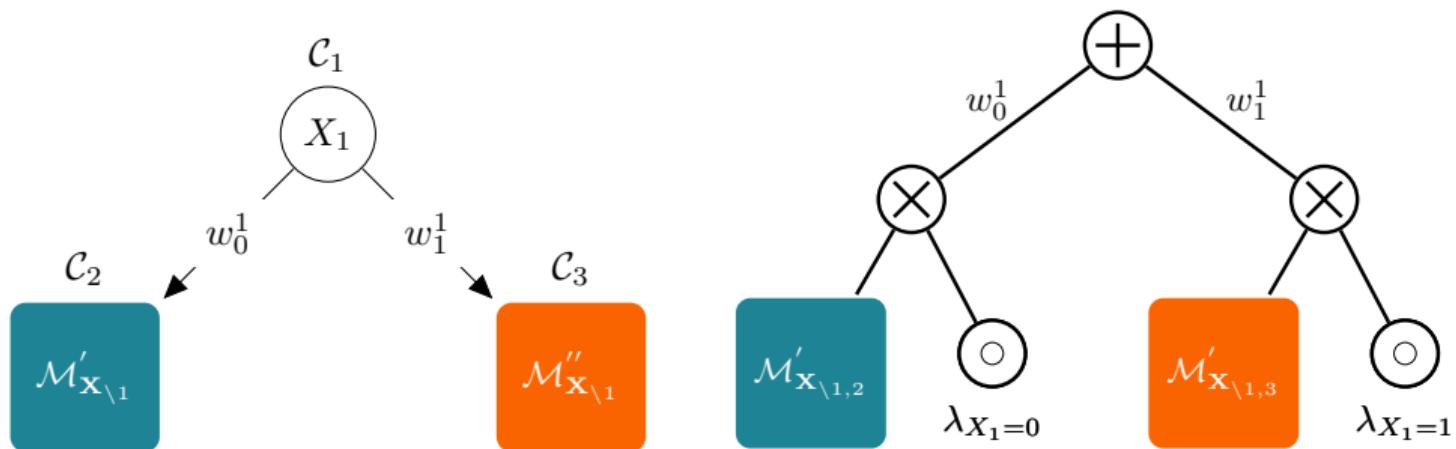
1. *Marginal inference in PGMs is exponential in the treewidth!*
but PCs can exploit context specific independence

$$\mathbf{X} \perp\!\!\!\perp \mathbf{Y} \mid Z = z_1$$

but

$$\mathbf{X} \not\perp\!\!\!\perp \mathbf{Y} \mid Z = z_2$$

Decision trees as PCs...



...not so easy!

1. Marginal inference in PGMs is exponential in the treewidth!

but PCs can exploit **context specific independence**

2. We do not know how to compile exactly an arbitrary PGM over continuous vars!

we need to extend the language of PCs to **integral units**

Probabilistic Integral Circuits

...why PCs?

1. A grammar for tractable models

One formalism to represent many probabilistic and logical models

⇒ #HMMs #Trees #XGBoost, Tensor Networks, ...
and other PGMs...

...why PCs?

1. A grammar for tractable models

One formalism to represent many probabilistic and logical models

⇒ #HMMs #Trees #XGBoost, Tensor Networks, ...
and other PGMs...

2. Expressiveness

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

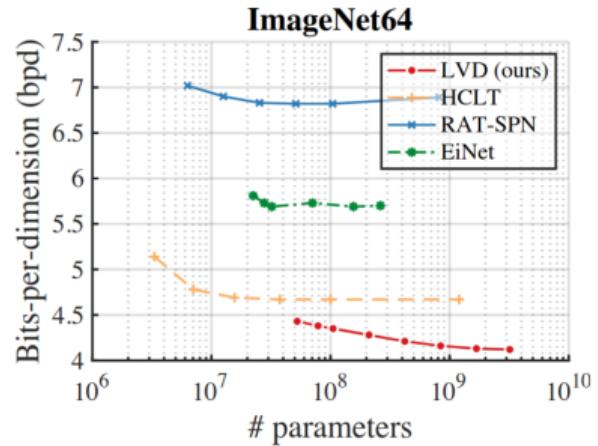
How expressive?

	QPQ	PC	Sp-PC	HCLT	RAT	IDF	BitS	BBans	McB
MNIST	1.11	1.17	1.14	1.21	1.67	1.90	1.27	1.39	1.98
F-MNIST	3.16	3.32	3.27	3.34	4.29	3.47	3.28	3.66	3.72
EMN-MN	1.55	1.64	1.52	1.70	2.56	2.07	1.88	2.04	2.19
EMN-LE	1.54	1.62	1.58	1.75	2.73	1.95	1.84	2.26	3.12
EMN-BA	1.59	1.66	1.60	1.78	2.78	2.15	1.96	2.23	2.88
EMN-BY	1.53	1.47	1.54	1.73	2.72	1.98	1.87	2.23	3.14

competitive with Flows and VAEs!

How scalable?

Dataset	TPMs				DGMs		
	LVD (ours)	HCLT	EiNet	RAT-SPN	Glow	RealNVP	BIVA
ImageNet32	4.39\pm0.01	4.82	5.63	6.90	4.09	4.28	3.96
ImageNet64	4.12\pm0.00	4.67	5.69	6.82	3.81	3.98	-
CIFAR	4.38\pm0.02	4.61	5.81	6.95	3.35	3.49	3.08



up to billions of parameters

...why PCs?

1. A grammar for tractable models

One formalism to represent many probabilistic and logical models

⇒ #HMMs #Trees #XGBoost, Tensor Networks, ...
and other PGMs...

2. Expressiveness

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

3. Tractability == Structural Properties!!!

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

Structural properties

smoothness

decomposability

determinism

compatibility

Structural properties

property A

property B

property C

property D

Structural properties

property A

tractable computation of *arbitrary integrals*

property B

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

property C

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

property D

\Rightarrow tractable partition function

Structural properties

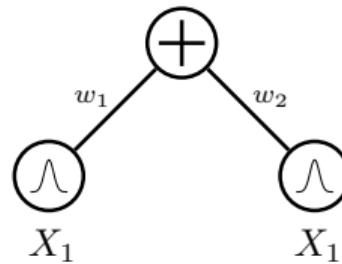
smoothness

the inputs of sum units are defined over the same variables

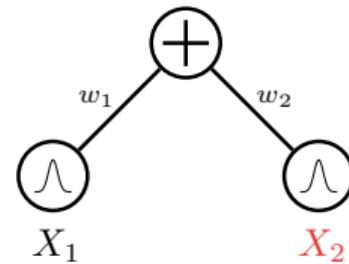
decomposability

compatibility

determinism



smooth circuit



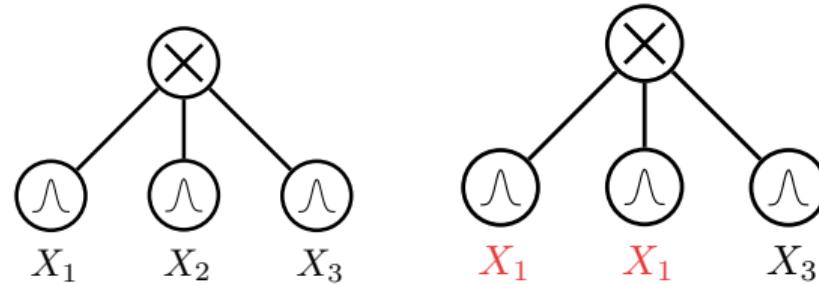
non-smooth circuit

Structural properties

smoothness

the inputs of prod units are defined over disjoint variable sets

decomposability



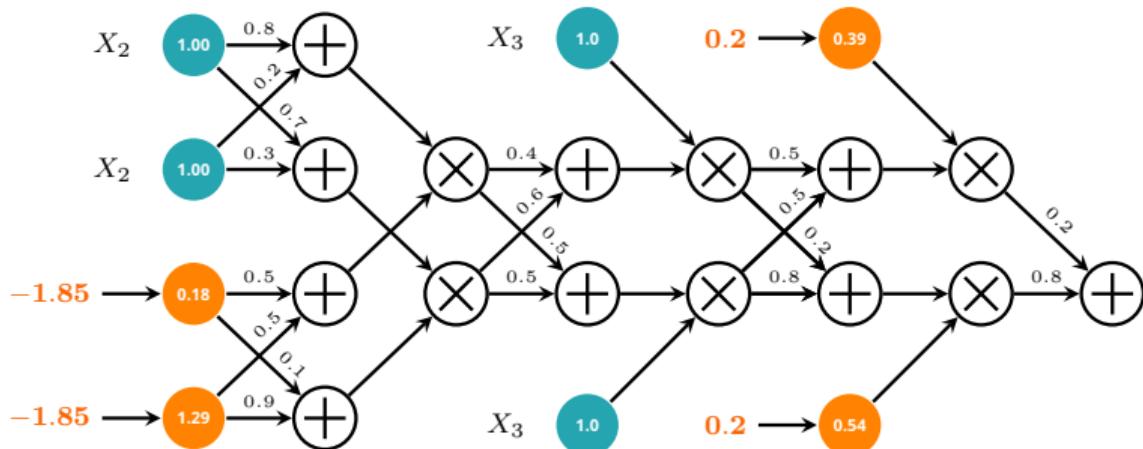
compatibility

determinism

decomposable circuit non-decomposable circuit

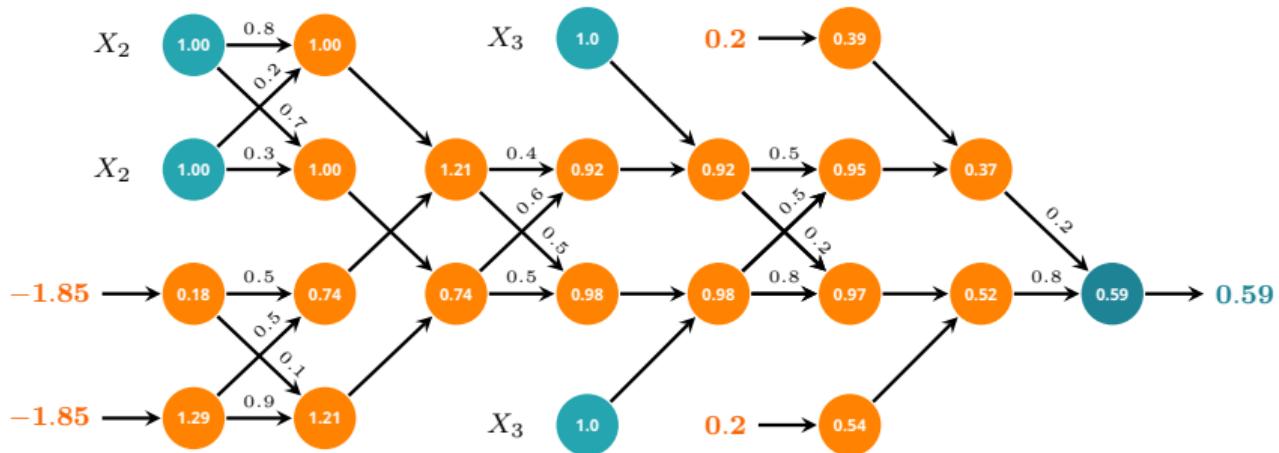
Probabilistic queries = feedforward evaluation

$$p(X_1 = -1.85, X_4 = 0.2)$$



Probabilistic queries = feedforward evaluation

$$p(X_1 = -1.85, X_4 = 0.2)$$



***smooth* + *decomposable* circuits = ...**

Computing arbitrary integrations (or summations)

⇒ *linear in circuit size!*

E.g., suppose we want to compute Z:

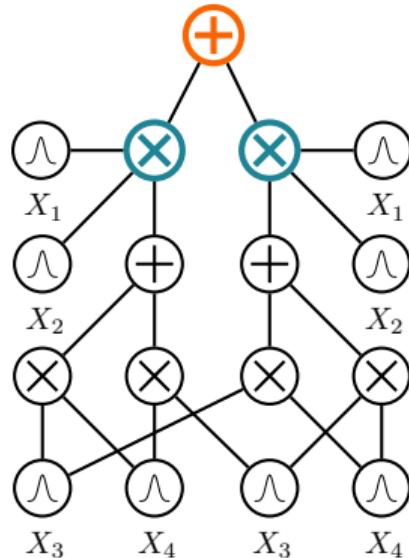
$$\int p(\mathbf{x}) d\mathbf{x}$$

***smooth* + *decomposable* circuits = ...**

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

$$\begin{aligned}\int \mathbf{p}(\mathbf{x}) d\mathbf{x} &= \int \sum_i w_i \mathbf{p}_i(\mathbf{x}) d\mathbf{x} = \\ &= \sum_i w_i \int \mathbf{p}_i(\mathbf{x}) d\mathbf{x}\end{aligned}$$

⇒ integrals are “pushed down” to inputs

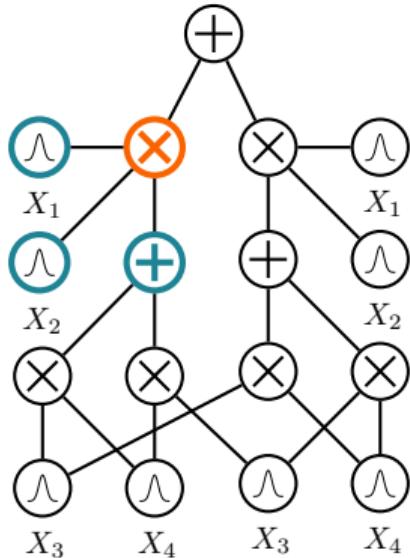


***smooth + decomposable* circuits = ...**

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

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

⇒ integrals decompose into easier ones



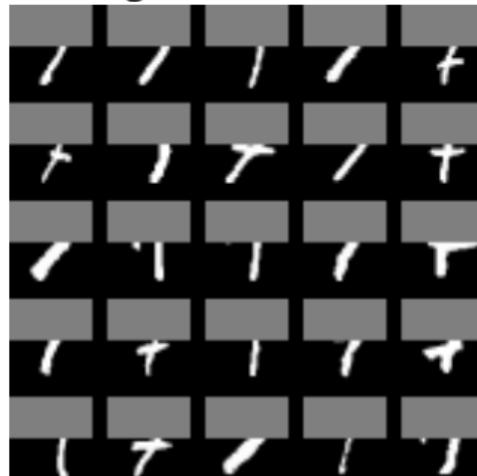
Tractable inference on PCs

Einsum networks

Original

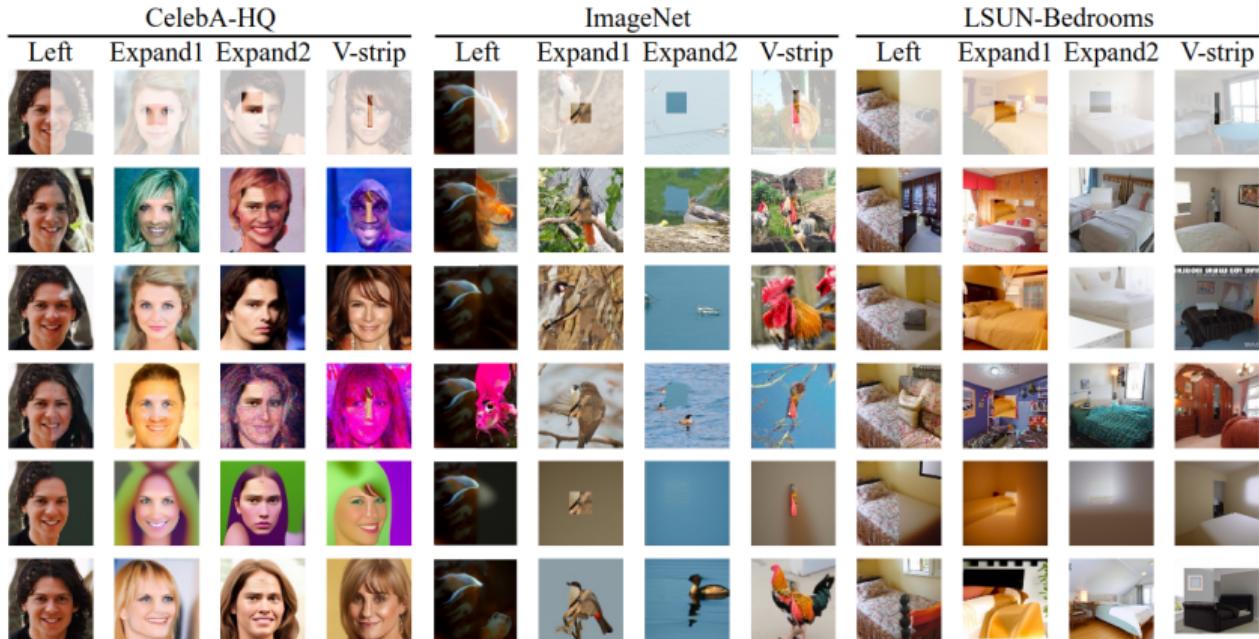


Missing



Conditional sample





Which structural properties

for complex reasoning



smooth + decomposable

Which structural properties

for complex reasoning



smooth + decomposable

??????

Adversarial smoothing

Certify robustness for inputs \mathbf{x} by
smoothing it by computing

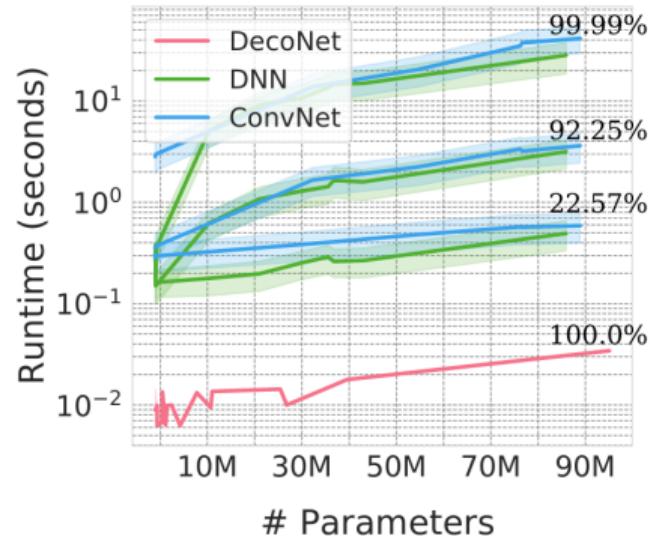
$$g_\sigma(\mathbf{x}) = \mathbb{E}_{\mathbf{e} \sim \mathcal{N}(0, \sigma \mathbf{I})} [f(\mathbf{x} + \mathbf{e})]$$



Adversarial smoothing

Certify robustness for inputs \mathbf{x} by
smoothing it by computing

$$g_\sigma(\mathbf{x}) = \mathbb{E}_{\mathbf{e} \sim \mathcal{N}(0, \sigma \mathbf{I})} [f(\mathbf{x} + \mathbf{e})]$$

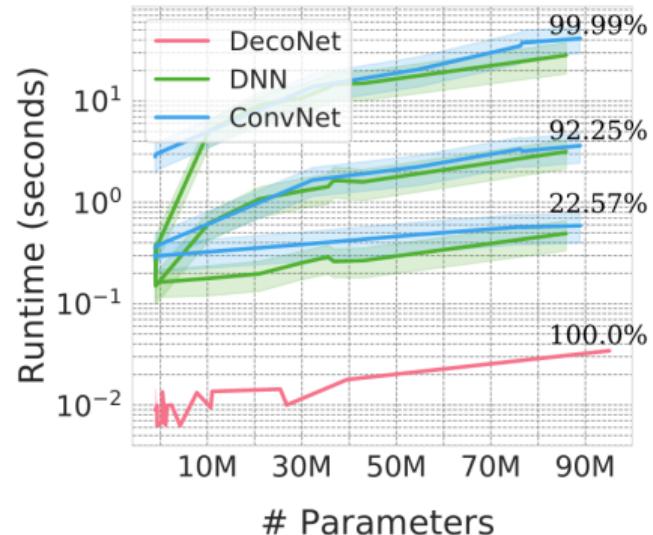


Adversarial smoothing

Certify robustness for inputs \mathbf{x} by **smoothing** it by computing

$$g_\sigma(\mathbf{x}) = \mathbb{E}_{\mathbf{e} \sim \mathcal{N}(0, \sigma \mathbf{I})} [f(\mathbf{x} + \mathbf{e})]$$

in a single feed-forward evaluation, if we **impose some structure** over a computational graph

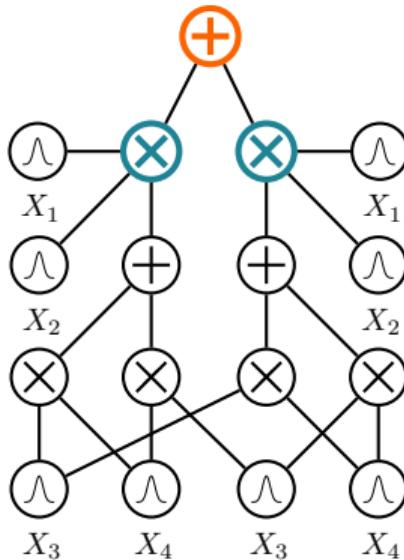


decomposable circuits = tractable *adv smoothing*

If $\mathbf{f}(\mathbf{x}) = \sum_i w_i \mathbf{f}_i(\mathbf{x})$:

$$\int \mathcal{N}(\mathbf{e}) \mathbf{f}(\mathbf{x} + \mathbf{e}) d\mathbf{e} = \sum_i w_i \mathbb{E}_{\mathcal{N}(\mathbf{e})} [\mathbf{f}_i(\mathbf{x} + \mathbf{e})]$$

⇒ expectations are “pushed down” to inputs



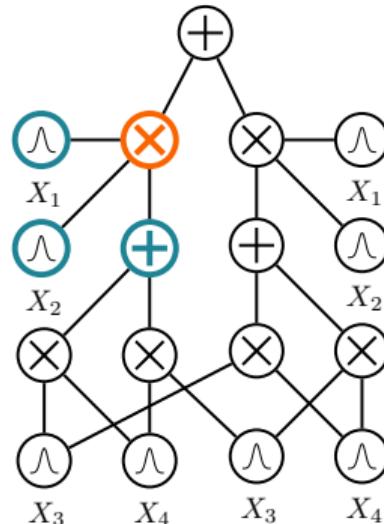
decomposable circuits = tractable *adv smoothing*

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

$$\int \mathcal{N}(\mathbf{e}_x)\mathcal{N}(\mathbf{e}_y)\mathcal{N}(\mathbf{e}_z) \mathbf{f}(\mathbf{x} + \mathbf{e}_x, \mathbf{y} + \mathbf{e}_y, \mathbf{z} + \mathbf{e}_z) d\mathbf{e}_x d\mathbf{e}_y d\mathbf{e}_z$$

$$\mathbb{E}_{\mathbf{e}_x}[\mathbf{f}(\mathbf{x} + \mathbf{e}_x)] \cdot \mathbb{E}_{\mathbf{e}_y}[\mathbf{f}(\mathbf{y} + \mathbf{e}_y)] \cdot \mathbb{E}_{\mathbf{e}_z}[\mathbf{f}(\mathbf{z} + \mathbf{e}_z)]$$

\Rightarrow expectations decompose into easier ones



Which structural properties

for complex reasoning



smooth + decomposable

decomposable

Which structural properties

for complex reasoning



smooth + decomposable



???????



decomposable

General expectations

Integrals involving two or more functions:

$$\int \textcolor{orange}{p}(\mathbf{x}) \textcolor{teal}{f}(\mathbf{x}) d \mathbf{X}$$



General expectations

Integrals involving two or more functions:

$$\int \textcolor{orange}{p}(\mathbf{x}) \textcolor{teal}{f}(\mathbf{x}) d\mathbf{X}$$

represent both $\textcolor{orange}{p}$ and $\textcolor{teal}{f}$ as circuits...but with which structural properties? E.g.,



General expectations

Integrals involving two or more functions:

$$\int \textcolor{orange}{p}(\mathbf{x}) \textcolor{teal}{f}(\mathbf{x}) d\mathbf{X}$$

represent both $\textcolor{orange}{p}$ and $\textcolor{teal}{f}$ as circuits...but with which structural properties? E.g.,

$$\mathbb{E}_{\mathbf{x}_c \sim p(\mathbf{x}_c | X_s=0)} [f_0(\mathbf{x}_c)] - \mathbb{E}_{\mathbf{x}_c \sim p(\mathbf{x}_c | X_s=1)} [f_1(\mathbf{x}_c)]$$



Structural properties

smoothness

decomposability

compatibility

determinism

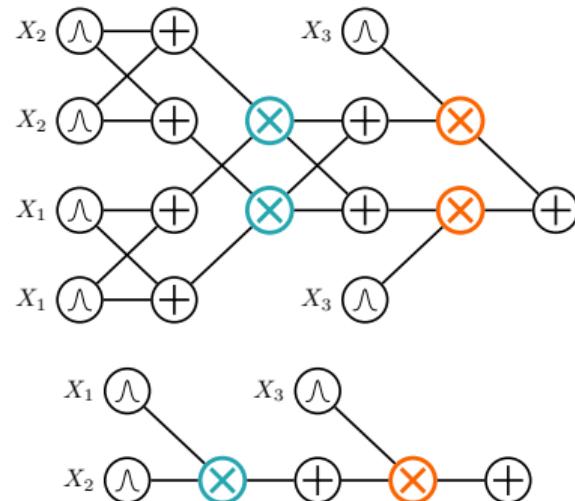
Structural properties

smoothness

decomposability

compatibility

determinism



compatible circuits

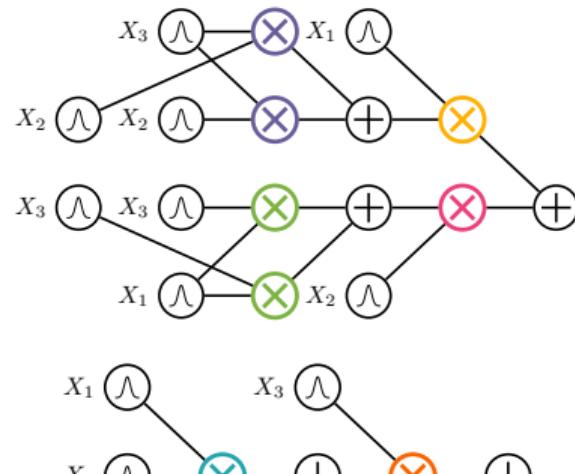
Structural properties

smoothness

decomposability

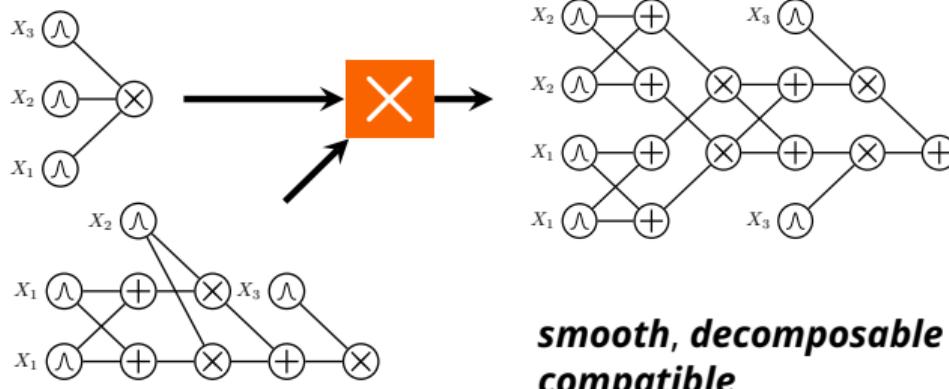
compatibility

determinism



non-compatible circuits

Tractable products

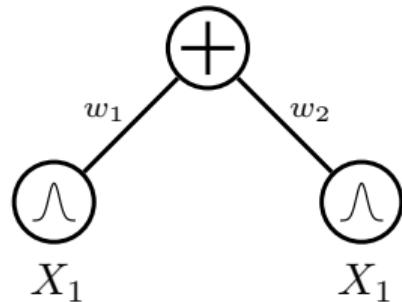


***smooth, decomposable
compatible***

exactly compute $\int \mathbf{p}(\mathbf{x}) \mathbf{f}(\mathbf{x}) d\mathbf{X}$ in time $O(|\mathbf{p}| |\mathbf{f}|)$

Tractable products

$$p(\mathbf{X}) = \sum_{i=1}^K w_i p_i(\mathbf{X}), \quad \sum_{i=1}^K w_i = 1, \quad w_i \geq 0,$$



Tractable products

$$p(\mathbf{X}) = \sum_{i=1}^K w_i p_i(\mathbf{X}), \quad \sum_{i=1}^K w_i = 1, \quad w_i \geq 0,$$

$$q(\mathbf{X}) = \sum_{j=1}^L w'_j q_j(\mathbf{X}), \quad \sum_{j=1}^L w'_j = 1, \quad w'_j \geq 0,$$

Tractable products

$$p(\mathbf{X}) = \sum_{i=1}^K w_i p_i(\mathbf{X}), \quad \sum_{i=1}^K w_i = 1, \quad w_i \geq 0,$$

$$q(\mathbf{X}) = \sum_{j=1}^L w'_j q_j(\mathbf{X}), \quad \sum_{j=1}^L w'_j = 1, \quad w'_j \geq 0,$$

$$p(\mathbf{X}) \times q(\mathbf{X}) = \sum_{i=1}^K \sum_{j=1}^L w_i w'_j p_i(\mathbf{X}) q_j(\mathbf{X}),$$

cartesian product of “local” mixtures

Which structural properties

for complex reasoning



smooth + decomposable



compatibility



decomposable

Structural properties

smoothness

decomposability

compatibility

determinism

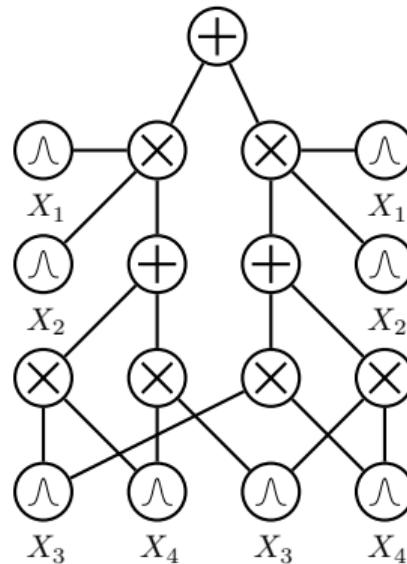
Determinism + **decomposability** = **tractable MAP**

Computing maximization with arbitrary evidence e

⇒ *linear in circuit size!*

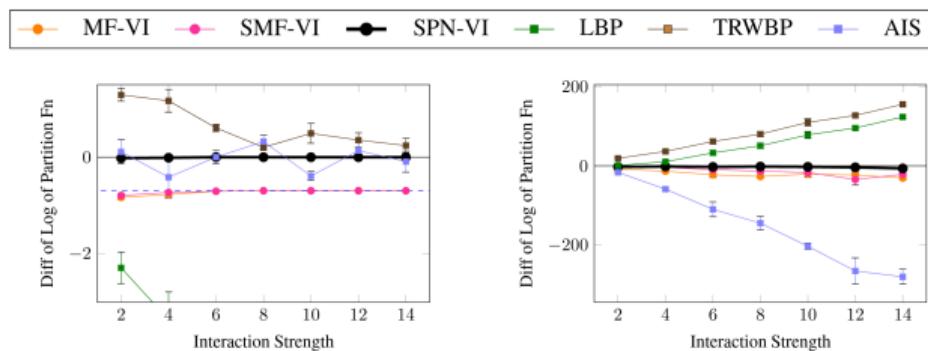
E.g., suppose we want to compute:

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



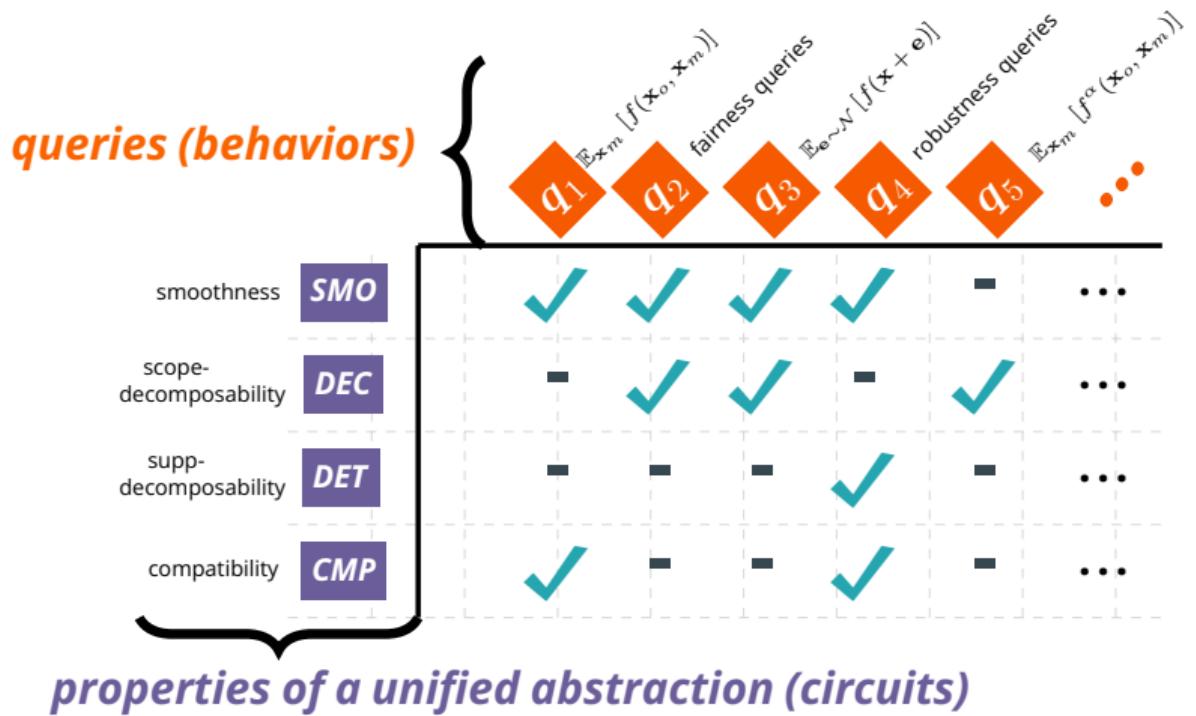
Determinism + decomposability = tractable ELBO

Using deterministic and decomposable PCs as expressive variational family \mathcal{Q} for discrete polynomial log-densities, i.e. $\operatorname{argmax}_{q \in \mathcal{Q}} \mathbb{E}_{\mathbf{x} \sim q} [\log w(\mathbf{x})] + \mathbb{H}(q)$



Closed-form computation for the entropy \mathbb{H} [Vergari et al. 2021]

Shih and Ermon, "Probabilistic Circuits for Variational Inference in Discrete Graphical Models", NeurIPS, 2020



<i>atomic queries</i>	\times	$p \times q$	$p + q$	pow_N	pow_R	p^n	p^q	$/$	p/q	\log	$\log p$	\exp	$\exp p$
smoothness	SMO	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
scope-decomposability	DEC	-	-	-	✓	-	-	-	-	-	-	-	-
supp-decomposability	DET	-	-	-	✓	✓	✓	✓	✓	✓	✓	-	-
compatibility	CMP	-	✓	✓	-	-	✓	-	-	-	-	-	-

✓ necess. conditions for reliability and efficiency

properties of a unified abstraction (circuits)

(De)composing queries

"What is the **expected prediction** for a patient with unavailable records?"

$$\int [p(\mathbf{x}_m \mid \mathbf{x}_o) \times f(\mathbf{x}_o, \mathbf{x}_m)]$$

(De)composing queries

"What is the **expected prediction** for a patient with unavailable records?"

$$\int [p(\mathbf{x}_m \mid \mathbf{x}_o) \times f(\mathbf{x}_o, \mathbf{x}_m)]$$

"What's the **variance** of my prediction?"

$$\int [p(\mathbf{x}_m \mid \mathbf{x}_o) \times \text{pow}(f(\mathbf{x}_o, \mathbf{x}_m), 2)] - \text{pow} \left(\int [p(\mathbf{x}_m \mid \mathbf{x}_o) \times f(\mathbf{x}_o, \mathbf{x}_m)], 2 \right)$$

(De)composing queries

"What is the **expected prediction** for a patient with unavailable records?"

$$\int [p(\mathbf{x}_m \mid \mathbf{x}_o) \times f(\mathbf{x}_o, \mathbf{x}_m)]$$

"What's the **variance** of my prediction?"

$$\int [p(\mathbf{x}_m \mid \mathbf{x}_o) \times \text{pow}(f(\mathbf{x}_o, \mathbf{x}_m), 2)] - \text{pow} \left(\int [p(\mathbf{x}_m \mid \mathbf{x}_o) \times f(\mathbf{x}_o, \mathbf{x}_m)], 2 \right)$$

"How **fair** is the prediction is a certain protected attribute changes?"

$$\int [p(\mathbf{x}_c \mid X_s = 0) \times f(\mathbf{x}_c, 0)] - \int [p(\mathbf{x}_c \mid X_s = 1) \times f(\mathbf{x}_c, 1)]$$

A language for queries

Integral expressions that can be formed by composing these operators

`+ , × , pow , log , exp` and `/`

⇒ *many divergences and information-theoretic queries*

A language for queries

Integral expressions that can be formed by composing these operators

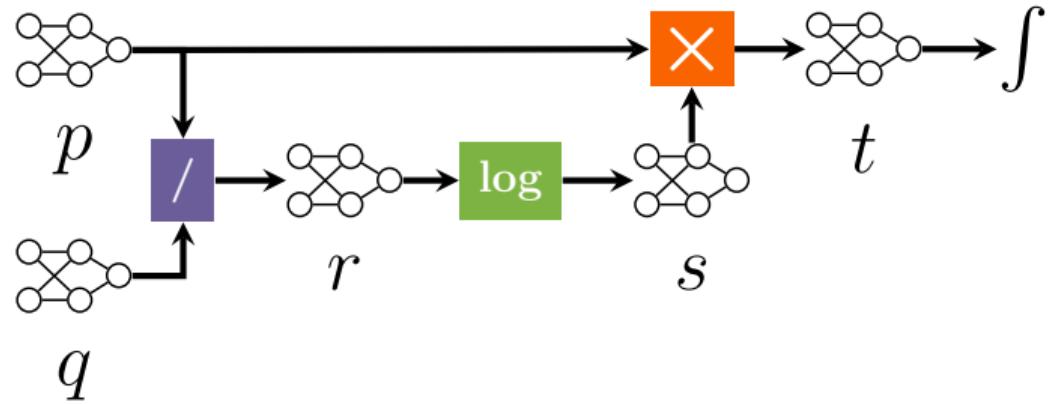
`+ , × , pow , log , exp` and `/`

⇒ *many divergences and information-theoretic queries*

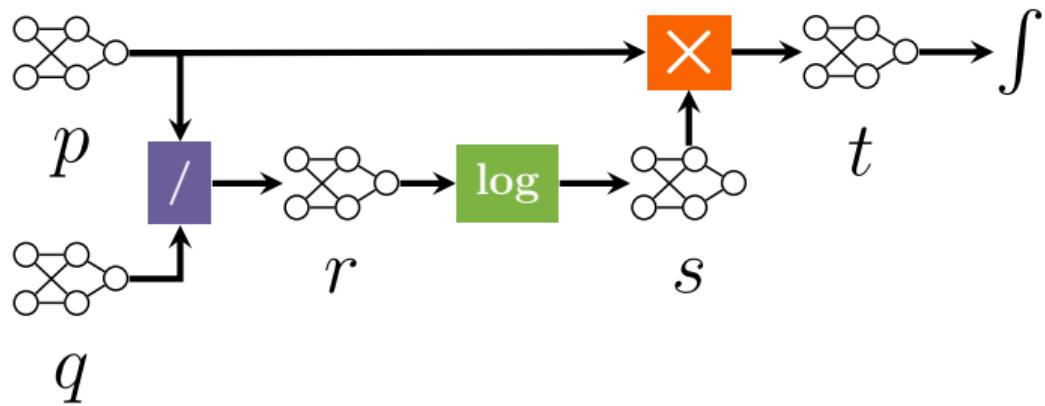
Represented as ***higher-order computational graphs***—pipelines—operating over circuits!

⇒ *re-using intermediate transformations across queries*

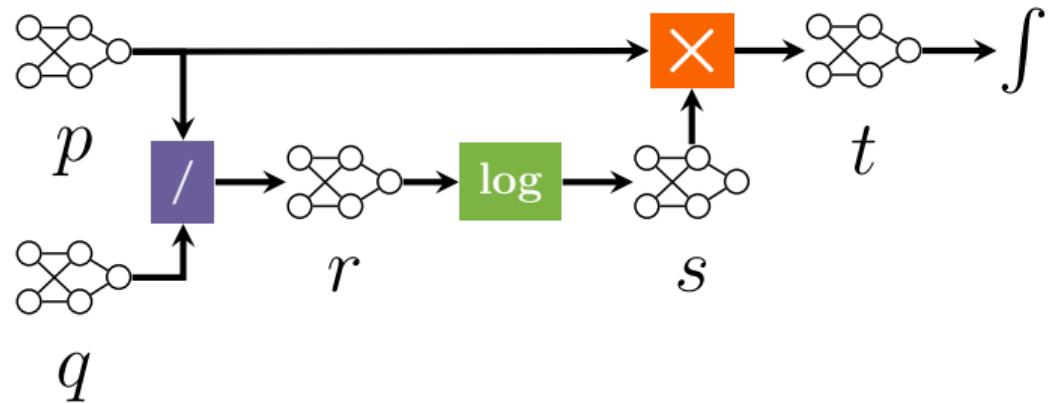
$$\mathbb{KLD}(p \parallel q) = \int_{\text{val}(\mathbf{X})} p(\mathbf{x}) \times \log(p(\mathbf{x})/q(\mathbf{x})) \, d\mathbf{X}$$



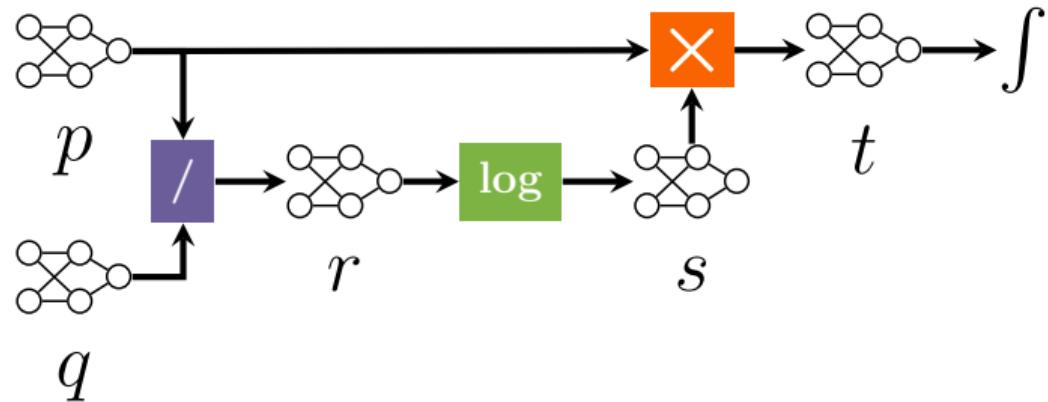
$$\text{KLD}(p \parallel q) = \int_{\text{val}(\mathbf{X})} p(\mathbf{x}) \times \log \left(p(\mathbf{x}) / q(\mathbf{x}) \right) d\mathbf{X}$$



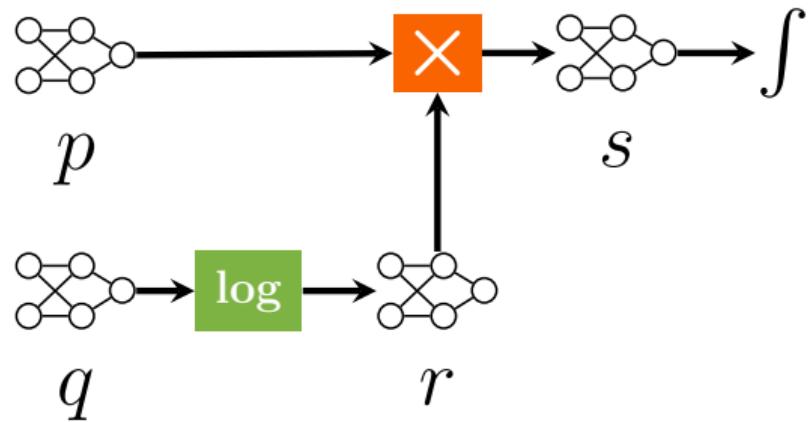
$$\mathbb{KLD}(p \parallel q) = \int_{\text{val}(\mathbf{X})} p(\mathbf{x}) \times \log(p(\mathbf{x})/q(\mathbf{x})) d\mathbf{X}$$



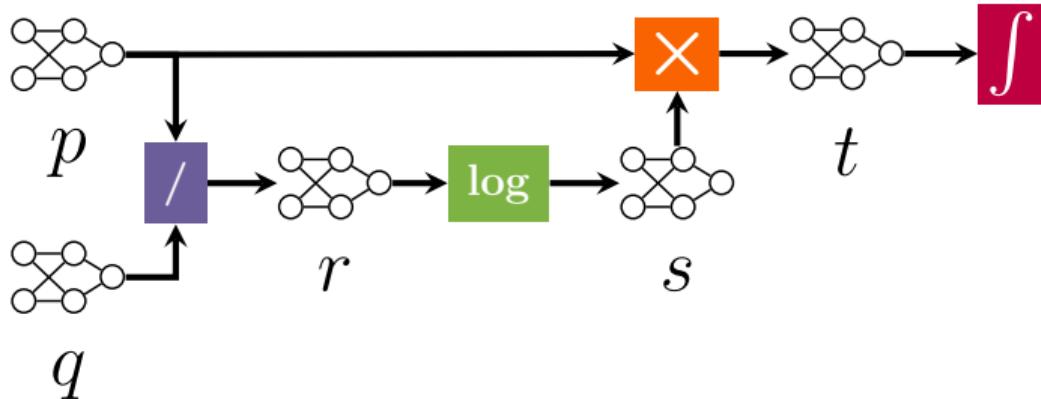
$$\mathbb{KLD}(p \parallel q) = \int_{\text{val}(\mathbf{X})} p(\mathbf{x}) \times \log(p(\mathbf{x})/q(\mathbf{x})) d\mathbf{X}$$



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

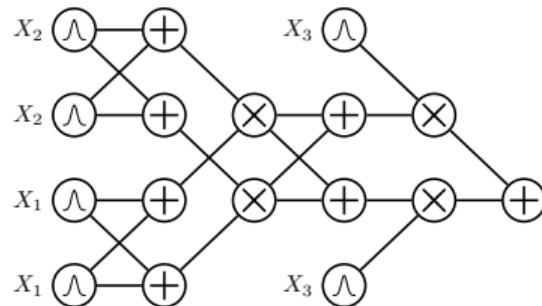
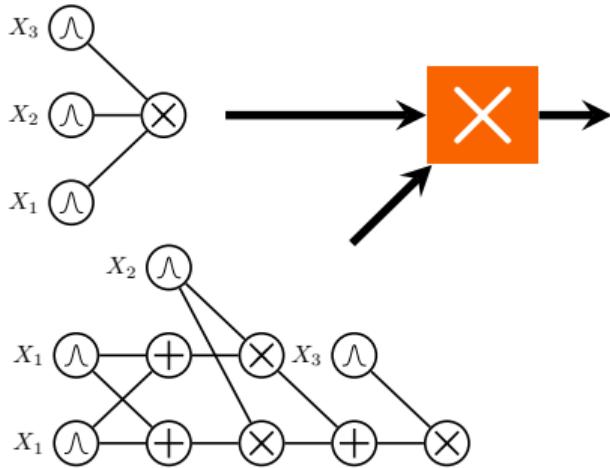


$$\int p(\mathbf{x}) \times \log \left(p(\mathbf{x}) / q(\mathbf{x}) \right) d\mathbf{X}$$



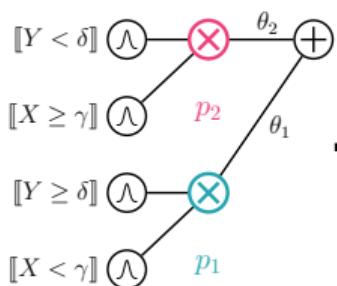
build a LEGO-like query calculus...

Tractable operators

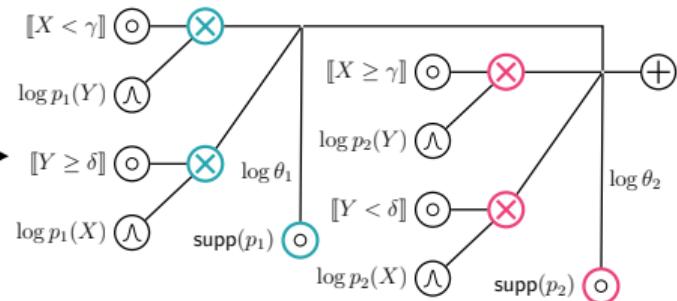


***smooth, decomposable
compatible***

Tractable operators

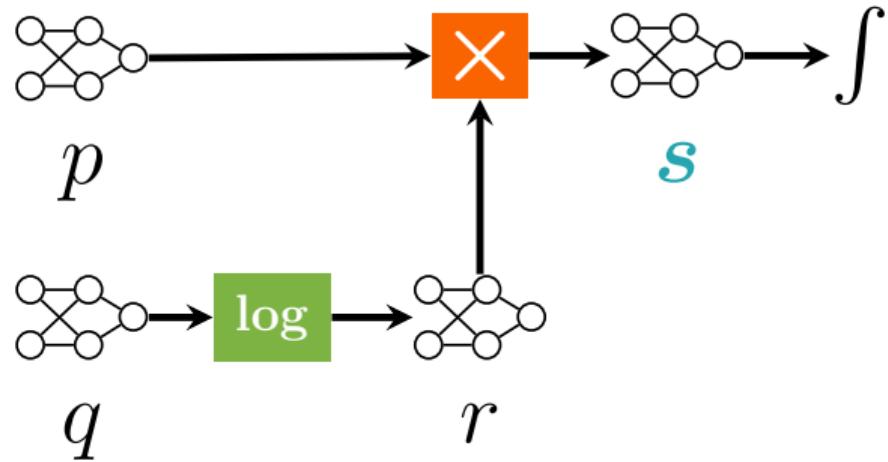


log

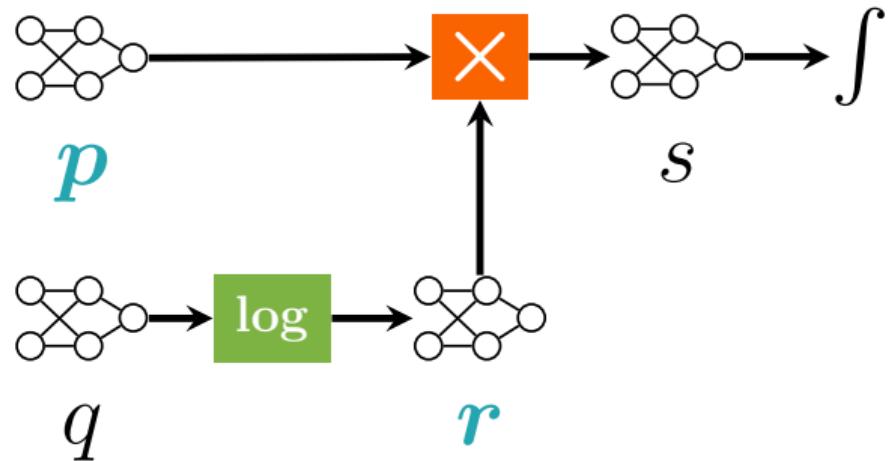


*smooth, decomposable
deterministic*

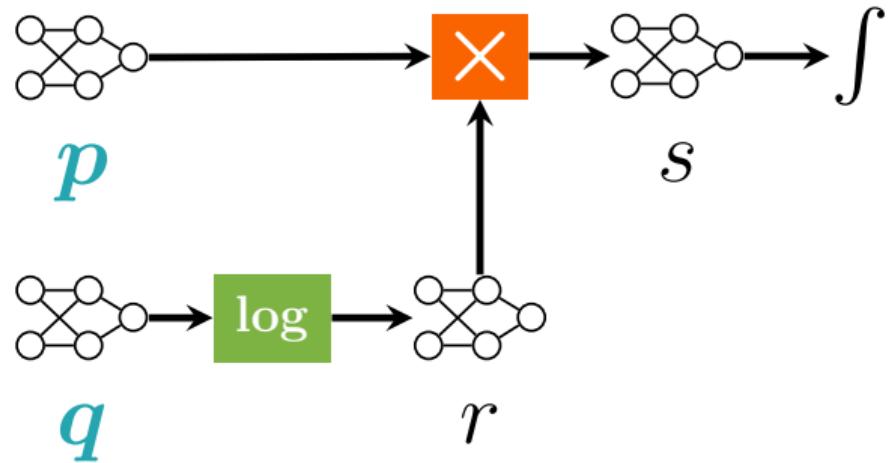
smooth, decomposable



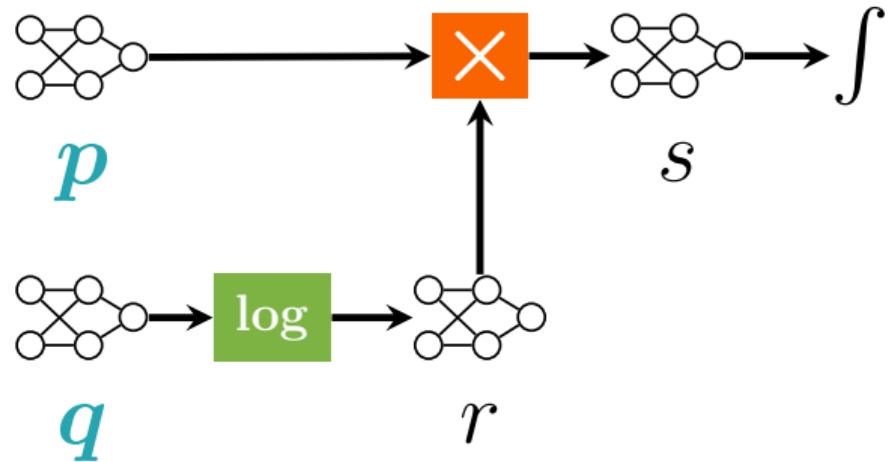
To perform tractable integration we need s to be *smooth and decomposable*...



hence we need p and r to be smooth, decomposable and **compatible**...

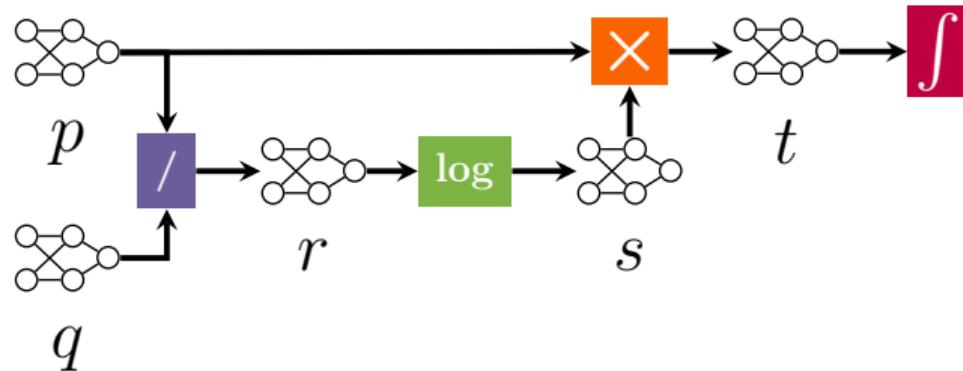


therefore q must be smooth, decomposable and **deterministic**...



we can compute XENT tractably if p and q are smooth, decomposable, compatible and q is deterministic...

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



build a LEGO-like query calculus!

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

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

and prove hardness when some input properties are not satisfied

Composable tractable sub-routines

```
function kld(p, q)          function xent(p, q)
    r = quotient(p, q)      r = log(q)
    s = log(r)              s = product(p, r)
    t = product(p, s)       return -integrate(s)
    return integrate(t)     end

function ent(p)             function alphadiv(p, q, alpha=1.5)
    q = log(p)             r = real_pow(p, alpha)
    r = product(p, q)       s = real_pow(q, 1.0-alpha)
    return -integrate(s)   t = product(r, s)
    end                     return log(integrate(t)) / (1.0-alpha)
                           end
```

*Efficient inference algorithms in a couple lines of Julia code!*¹

¹<https://github.com/UCLA-StarAI/circuit-ops-atlas>

Learning probabilistic circuits

Learning probabilistic circuits

Probabilistic circuits are (peculiar) neural networks... ***just backprop with SGD!***

Learning probabilistic circuits

Probabilistic circuits are (peculiar) neural networks... ***just backprop with SGD!***

...end of Learning section!

Learning probabilistic circuits

Probabilistic circuits are (peculiar) neural networks... ***just backprop with SGD!***

wait but...

SGD is slow to converge... can we do better?

How to learn normalized weights?

Can we exploit structural properties somehow?

Expectation-Maximization

for smooth and decomposable circuits

$$w_{i,j}^{new} \leftarrow \frac{\sum_{\mathbf{x} \in \mathcal{D}} p[ctx_i = 1, Z_i = j \mid \mathbf{x}; \mathbf{w}^{old}]}{\sum_{\mathbf{x} \in \mathcal{D}} p[ctx_i = 1 \mid \mathbf{x}; \mathbf{w}^{old}]}$$

Darwiche, "A Differential Approach to Inference in Bayesian Networks",, 2003
Peharz et al., "On the Latent Variable Interpretation in Sum-Product Networks",, 2016

Expectation-Maximization

for smooth and decomposable circuits

$$w_{i,j}^{new} \leftarrow \frac{\sum_{\mathbf{x} \in \mathcal{D}} p[ctx_i = 1, Z_i = j \mid \mathbf{x}; \mathbf{w}^{old}]}{\sum_{\mathbf{x} \in \mathcal{D}} p[ctx_i = 1 \mid \mathbf{x}; \mathbf{w}^{old}]}$$

We get **all** the required statistics with a single backprop pass:

$$p[ctx_i = 1, Z_i = j \mid \mathbf{x}; \mathbf{w}^{old}] = \frac{1}{p(\mathbf{x})} \frac{\partial p(\mathbf{x})}{\partial S_i(\mathbf{x})} N_j(\mathbf{x}) w_{i,j}^{old}$$

Darwiche, "A Differential Approach to Inference in Bayesian Networks",, 2003
Peharz et al., "On the Latent Variable Interpretation in Sum-Product Networks",, 2016

Bayesian parameter learning

Formulate a prior $p(\mathbf{w}, \boldsymbol{\theta})$ over sum-weights and parameters of input units. Then perform posterior inference:

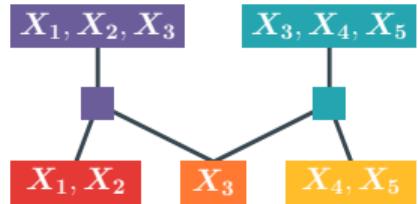
$$p(\mathbf{w}, \boldsymbol{\theta} | \mathcal{D}) \propto p(\mathbf{w}, \boldsymbol{\theta}) p(\mathcal{D} | \mathbf{w}, \boldsymbol{\theta})$$

Moment matching (oBMM) [Jaini et al. 2016; Rashwan, Zhao, and Poupart 2016]

Collapsed variational inference algorithm [Zhao et al. 2016]

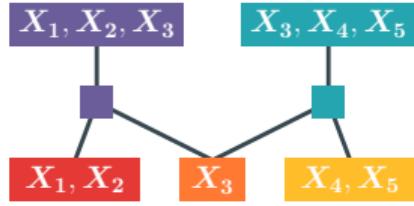
Gibbs sampling [Trapp et al. 2019; Vergari et al. 2019]

Learning recipe

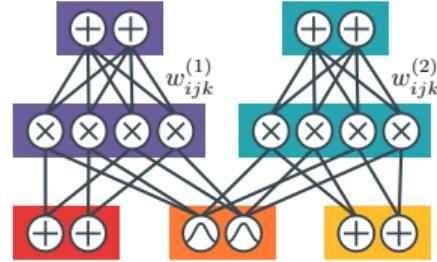


1) Build a *region graph*

Learning recipe

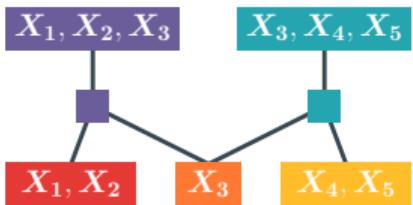


1) Build a *region graph*

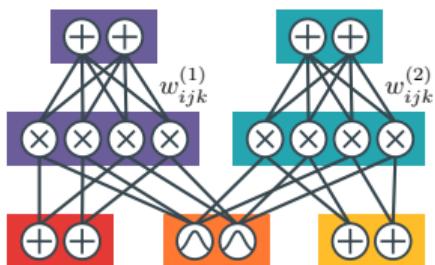


2) Overparameterize

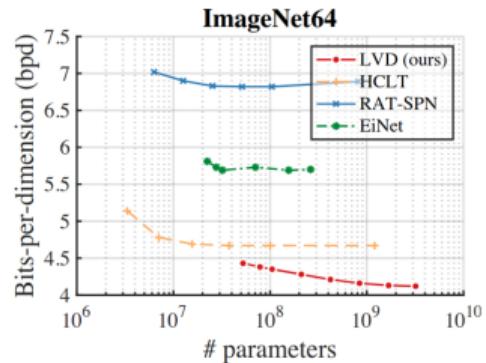
Learning recipe



1) Build a *region graph*

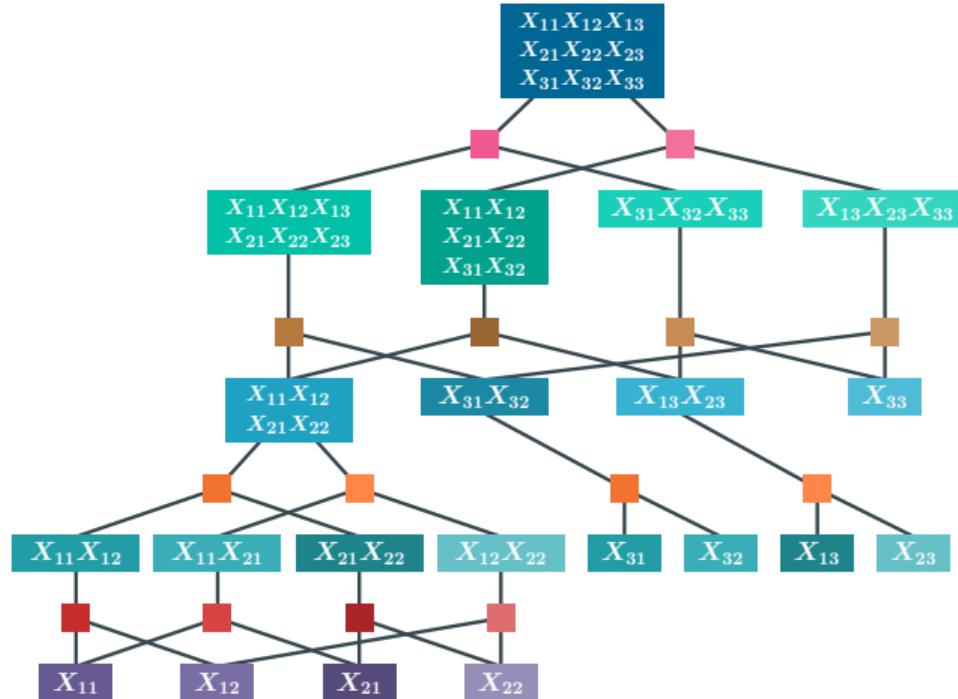


2) Overparameterize

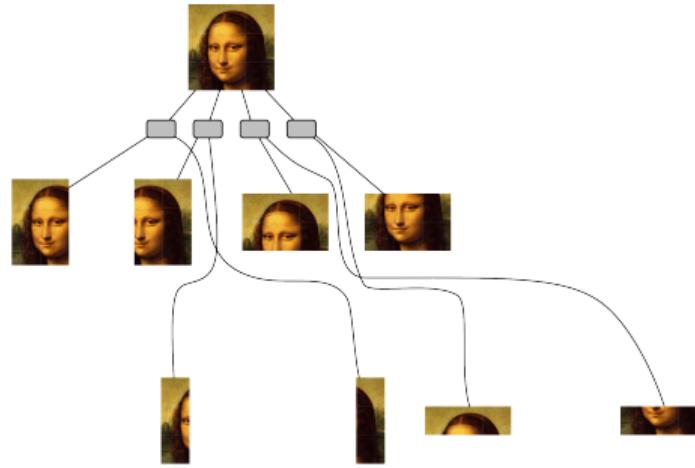


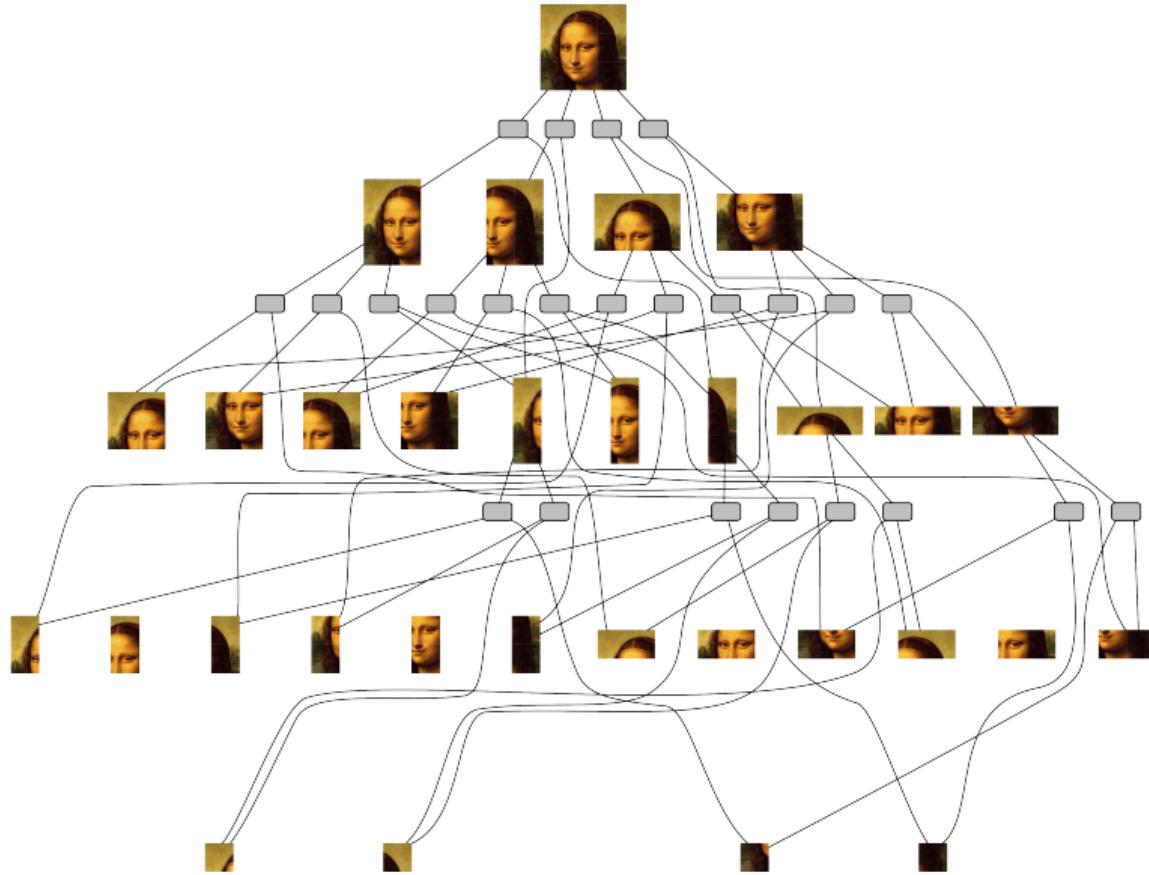
3) Learn parameters

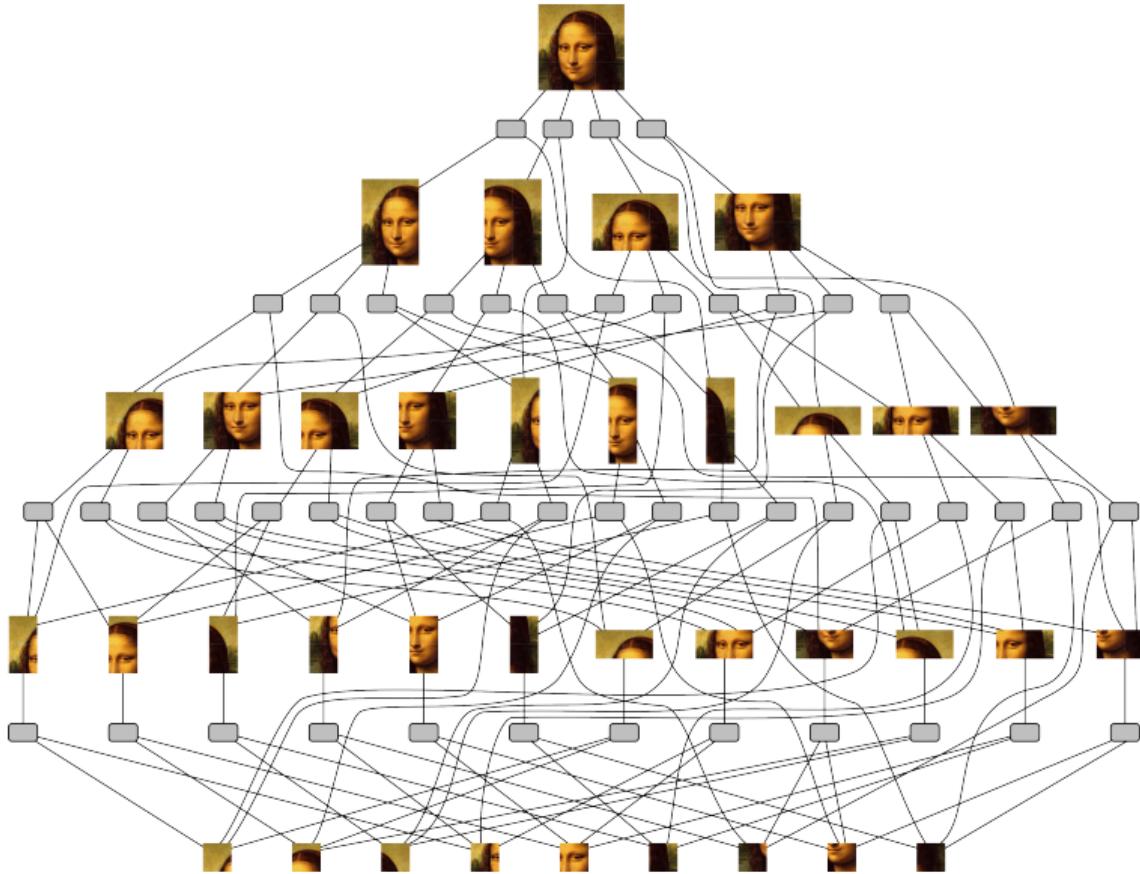
Which region graph?









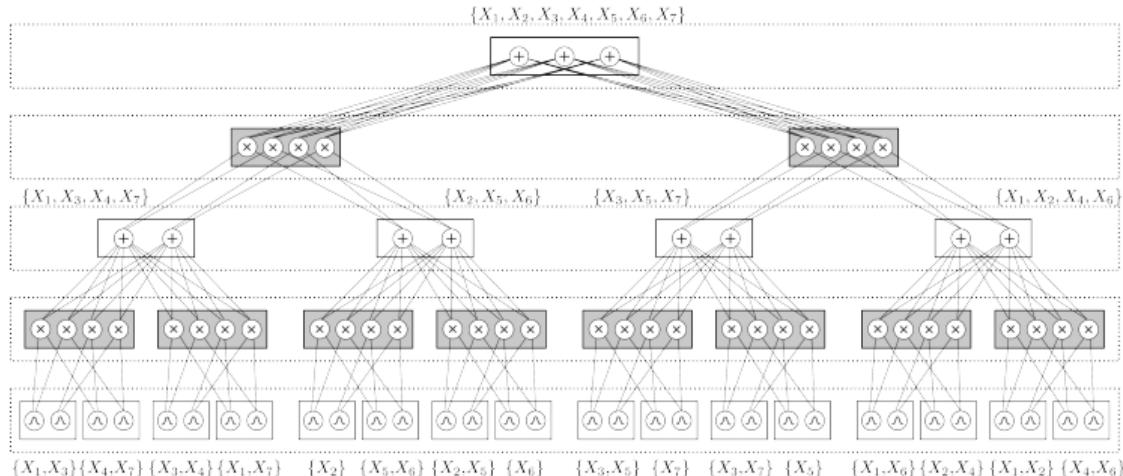


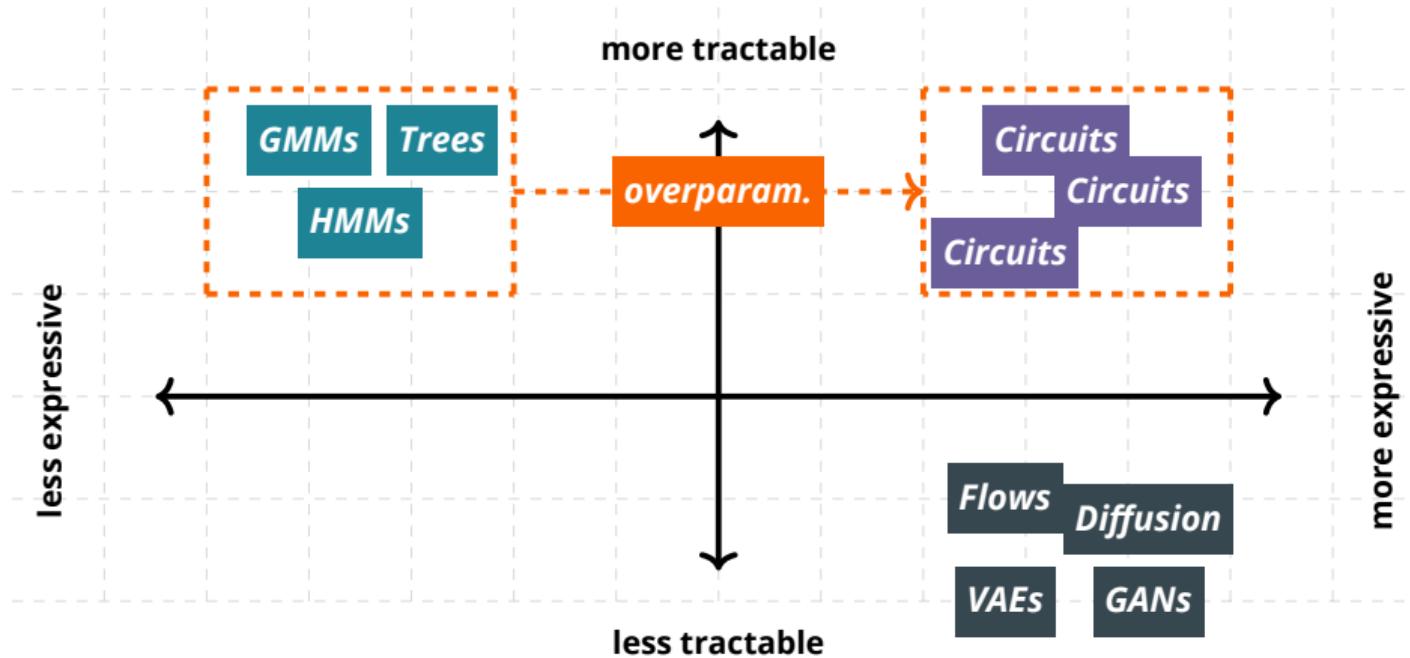
Random regions graphs

The “no-learning” option

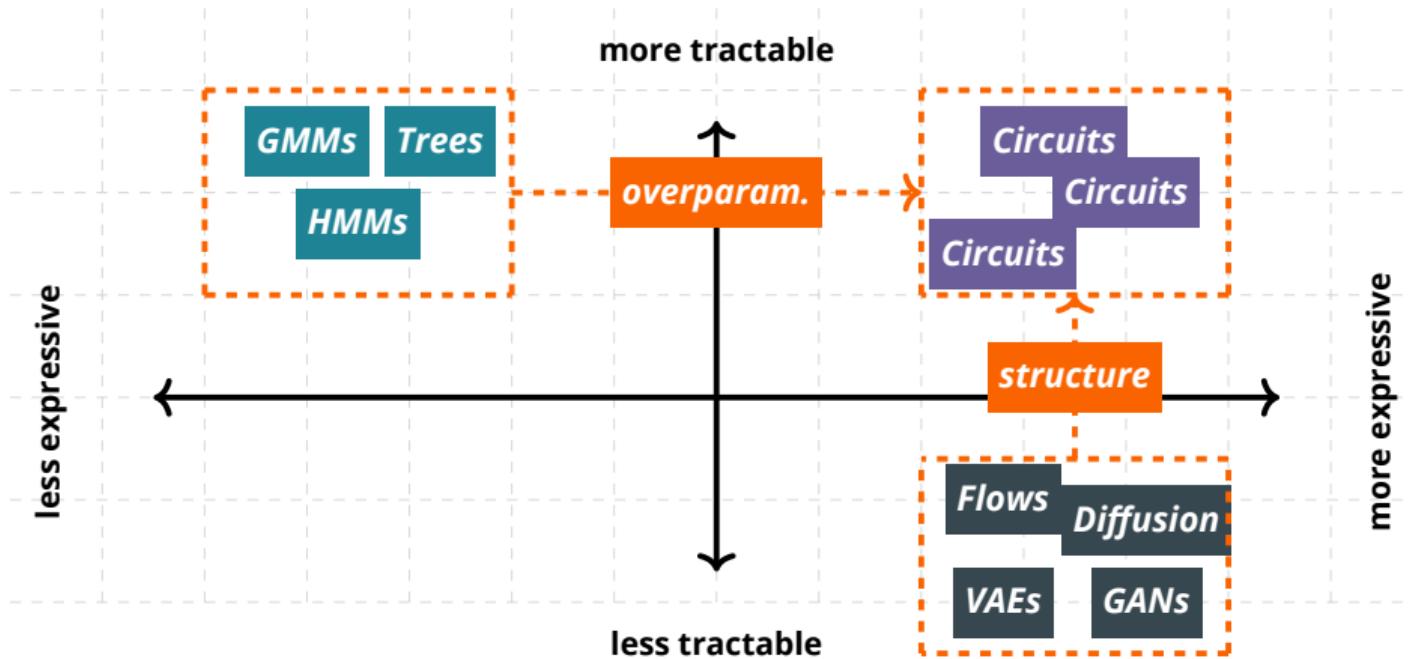
[Peharz et al. 2019]

Generating a random region graph, by recursively splitting \mathbf{X} into two random parts:





then make it more expressive!



impose structure!

<i>atomic queries</i>	\times	$+$	$p \times q$	$p + q$	pow_N	pow_R	p^n	p^q	$/$	p/q	\log	$\log p$	\exp	$\exp p$
smoothness	SMO	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
scope-decomposability	DEC	-	-	-	-	✓	-	-	-	-	-	-	-	-
supp-decomposability	DET	-	-	-	-	✓	-	✓	✓	✓	✓	-	-	-
compatibility	CMP	-	✓	✓	-	-	-	✓	-	-	-	-	-	-

✓ necess. conditions for reliability and efficiency

properties of a unified abstraction (circuits)

what really matters is structure!