



# ***tractable probabilistic modeling with probabilistic circuits***

**antonio vergari** (he/him)

 @tetraduzione

16th Oct 2024 - PIC PhD School Copenhagen

*april*

april-tools.github.io

# *april*

*autonomous &  
provably  
reliable  
intelligent  
learners*

# *april*

*about  
probabilities  
integrals &  
logic*

# *april*

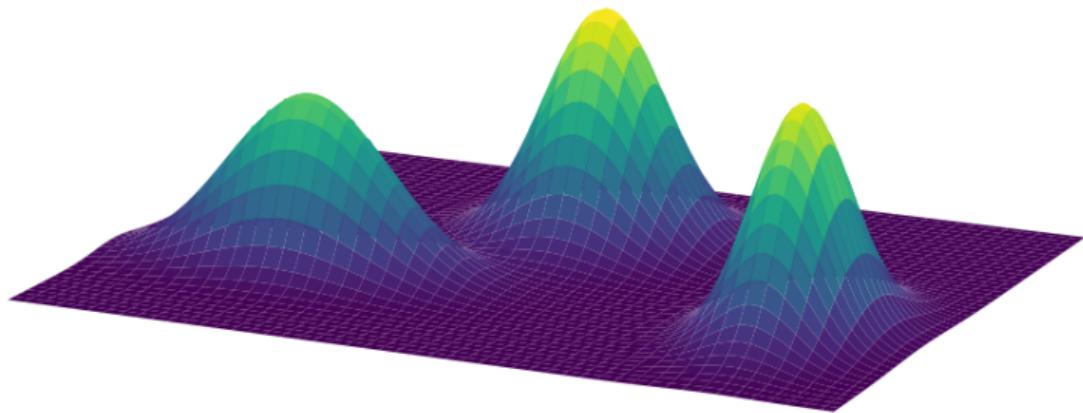
*april is  
probably a  
recursive  
identifier of a  
lab*

*deep generative models*

+

*flexible and reliable  
(logic &) probabilistic reasoning?*

- i) probabilistic circuits: syntax and semantics*
- ii) reliable and efficient neuro-symbolic AI*
- iii) beyond PCs: subtractive mixture models*



*a love letter to mixture models...*

# *Reasoning about ML models*



q<sub>1</sub>

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



q<sub>2</sub>

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



q<sub>3</sub>

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

# *Reasoning about ML models*



**q<sub>1</sub>**  $\int p(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{X}_m$   
*(missing values)*

**q<sub>2</sub>**  $\mathbb{E}_{\mathbf{x}_c \sim p(\mathbf{X}_c | X_s=0)} [f_0(\mathbf{x}_c)] - \mathbb{E}_{\mathbf{x}_c \sim p(\mathbf{X}_c | X_s=1)} [f_1(\mathbf{x}_c)]$   
*(fairness)*

**q<sub>3</sub>**  $\mathbb{E}_{\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_D)} [f(\mathbf{x} + \mathbf{e})]$   
*(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*



**q<sub>1</sub>**  $\int p(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{X}_m$   
*(missing values)*

**q<sub>2</sub>**  $\frac{\mathbb{E}_{\mathbf{x}_c \sim p(\mathbf{X}_c | X_s=0)} [f_0(\mathbf{x}_c)] - \mathbb{E}_{\mathbf{x}_c \sim p(\mathbf{X}_c | X_s=1)} [f_1(\mathbf{x}_c)]}{(fairness)}$

**q<sub>3</sub>**  $\mathbb{E}_{\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_D)} [f(\mathbf{x} + \mathbf{e})]$   
*(adversarial robust.)*

***hard to compute in general!***

# *Reasoning about ML models*



**q<sub>1</sub>**  $\int p(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{X}_m$   
*(missing values)*

**q<sub>2</sub>**  $\mathbb{E}_{\mathbf{x}_c \sim p(\mathbf{X}_c | X_s=0)} [f_0(\mathbf{x}_c)] - \mathbb{E}_{\mathbf{x}_c \sim p(\mathbf{X}_c | X_s=1)} [f_1(\mathbf{x}_c)]$   
*(fairness)*

**q<sub>3</sub>**  $\mathbb{E}_{\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_D)} [f(\mathbf{x} + \mathbf{e})]$   
*(adversarial robust.)*

*it is crucial we compute them exactly and in polytime!*

# *Reasoning about ML models*



**q<sub>1</sub>**  $\int p(\mathbf{x}_o, \mathbf{x}_m) d\mathbf{X}_m$   
*(missing values)*

**q<sub>2</sub>**  $\mathbb{E}_{\mathbf{x}_c \sim p(\mathbf{X}_c | X_s=0)} [f_0(\mathbf{x}_c)] - \mathbb{E}_{\mathbf{x}_c \sim p(\mathbf{X}_c | X_s=1)} [f_1(\mathbf{x}_c)]$   
*(fairness)*

**q<sub>3</sub>**  $\mathbb{E}_{\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_D)} [f(\mathbf{x} + \mathbf{e})]$   
*(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  $O(\text{poly}(|\mathbf{m}|))$ .

$\Rightarrow$  *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  $O(\text{poly}(|\mathbf{m}|))$ .

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

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

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

$\Rightarrow$  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]



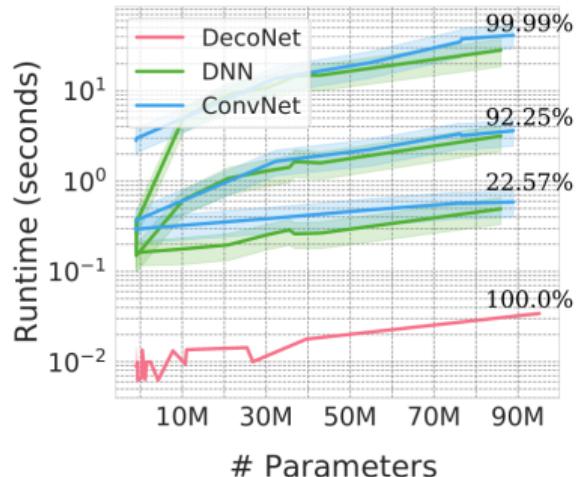
estimation error is bounded (0)  
[Choi 2022]

# **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	<b>53</b>
PC (large)	<b>1.20</b>	<b>1.24</b>	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  $\mathbf{q}$

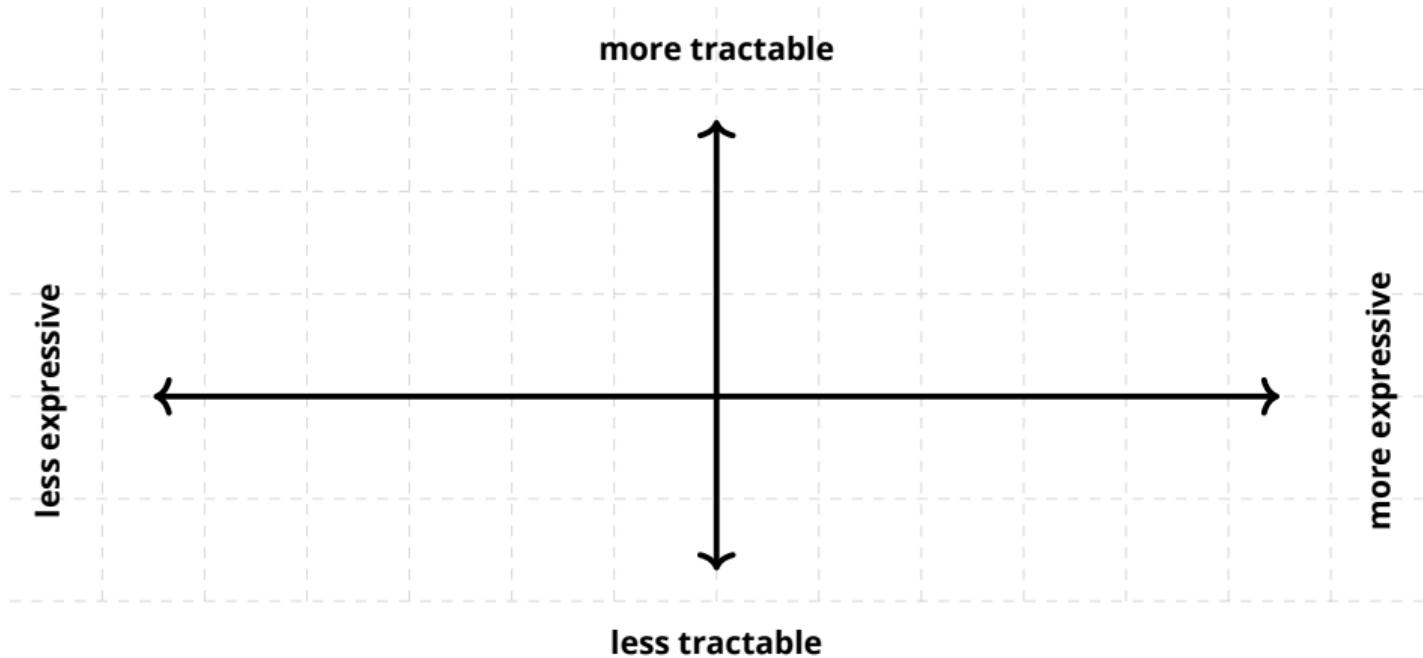
**MC**

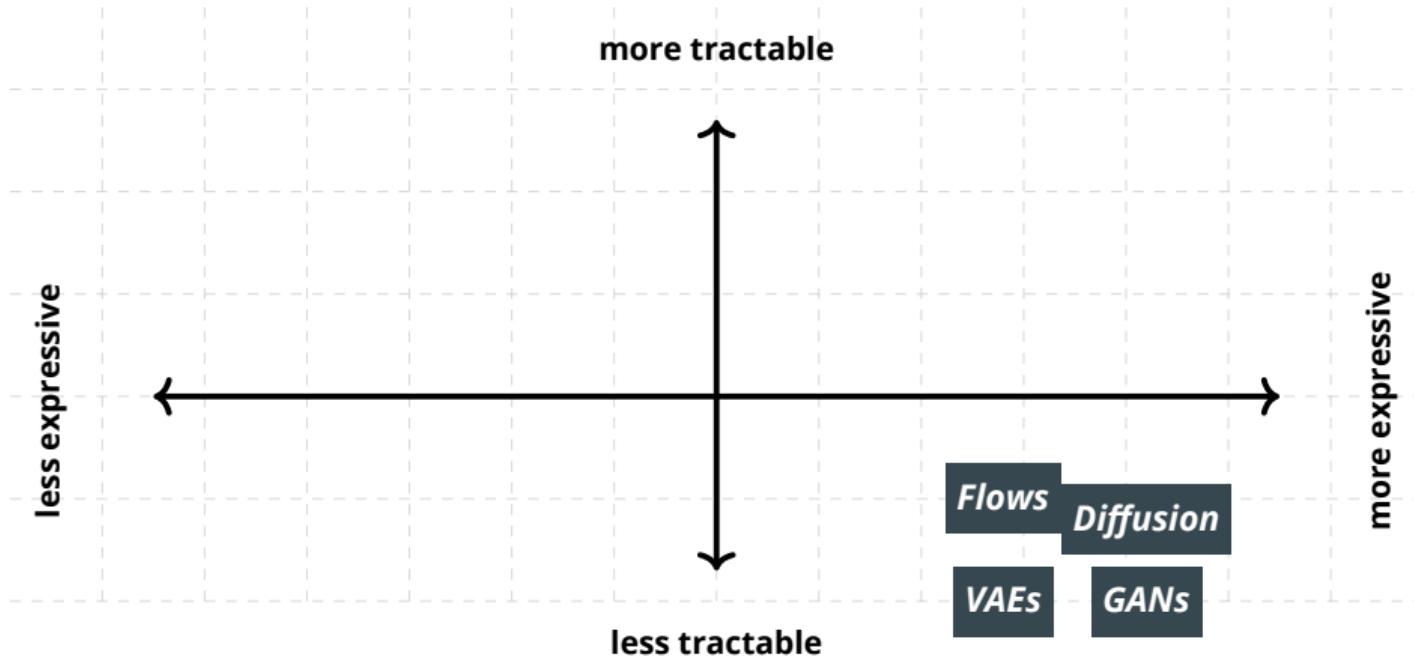
$$\mathbb{E}_{p(\mathbf{x})}[f(\mathbf{x})] \approx \frac{1}{N} \sum_{i=1}^N f(\mathbf{x}^{(i)})$$

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

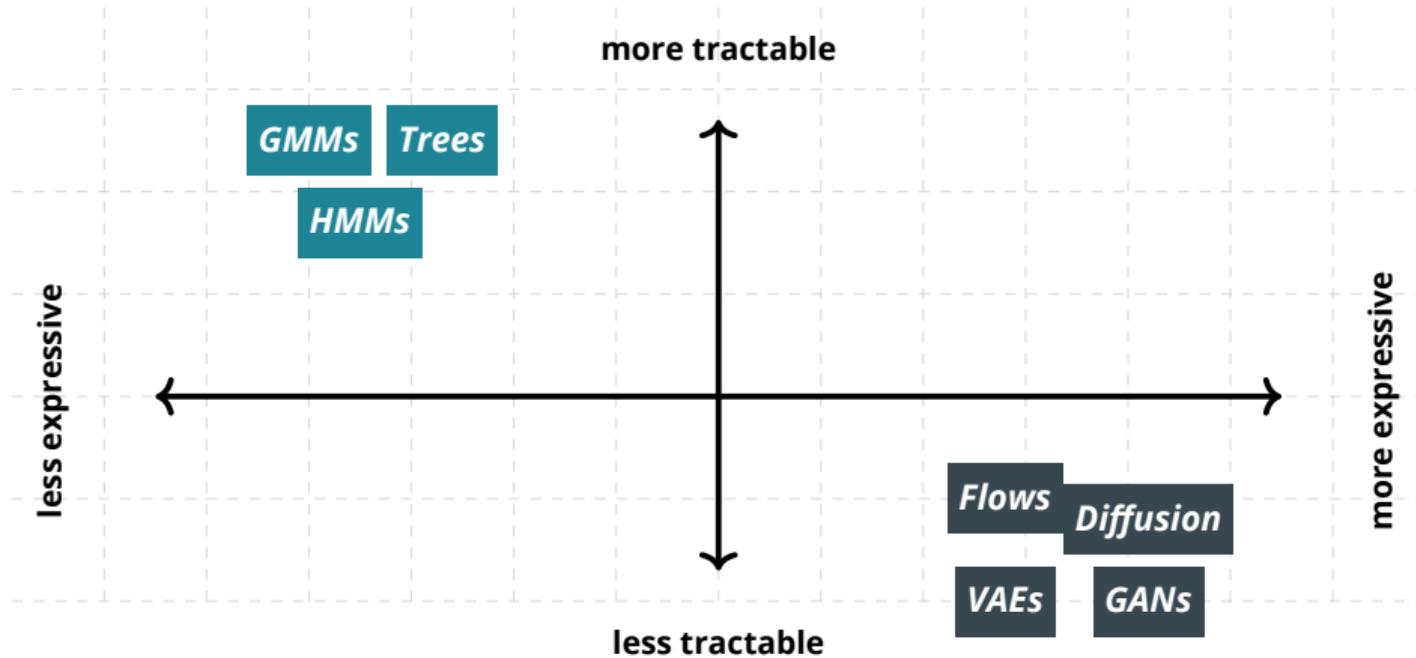
## *Goal*

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

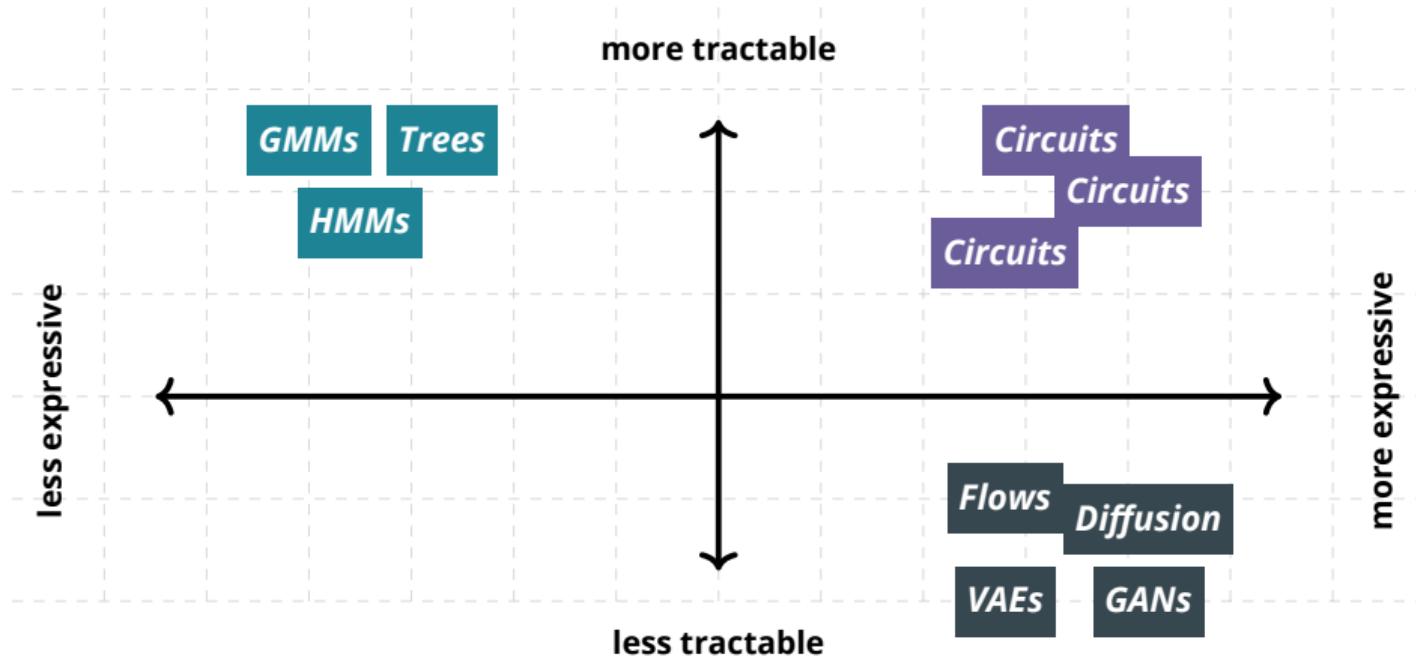




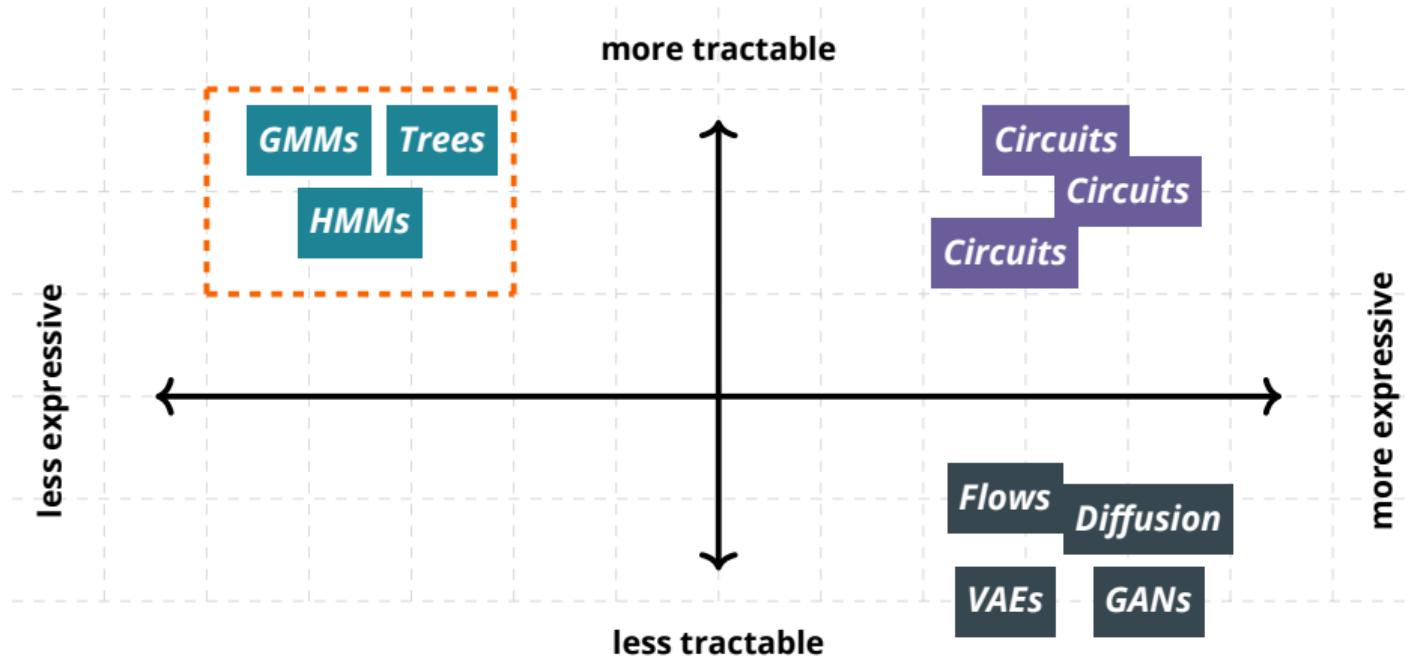
***expressive models are not very tractable...***



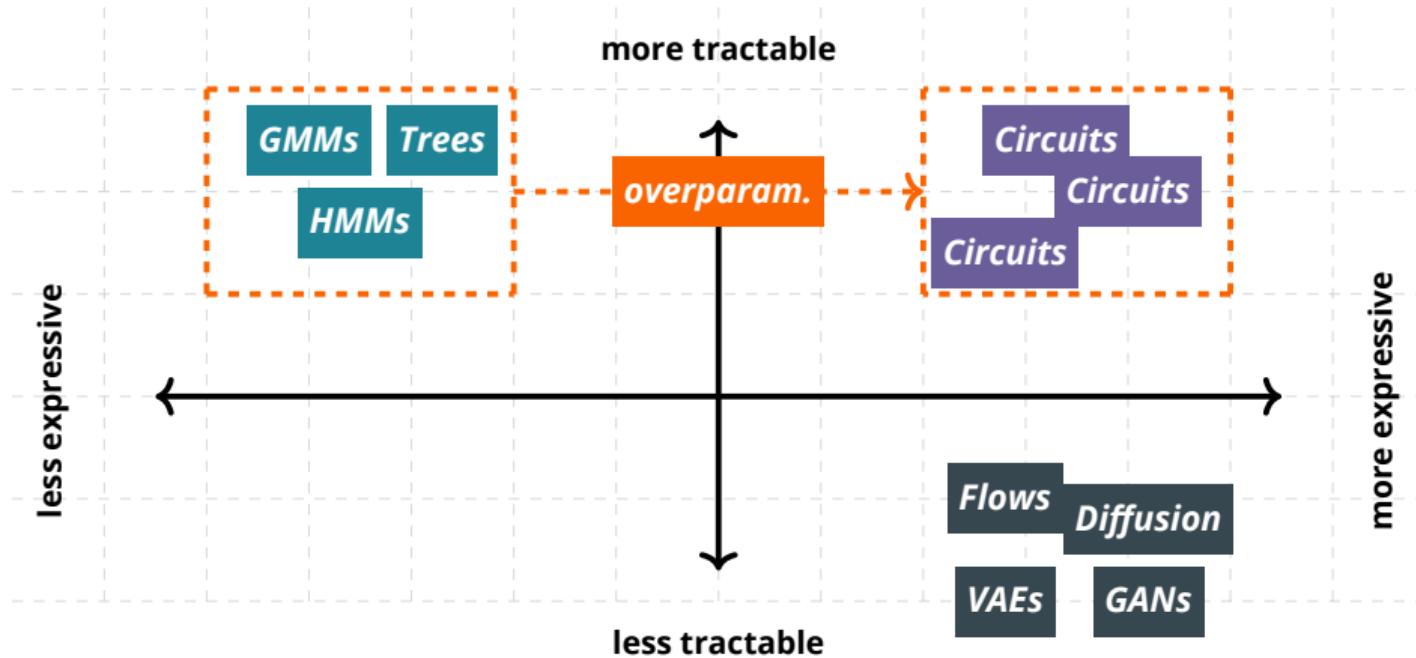
***tractable models are not that expressive...***



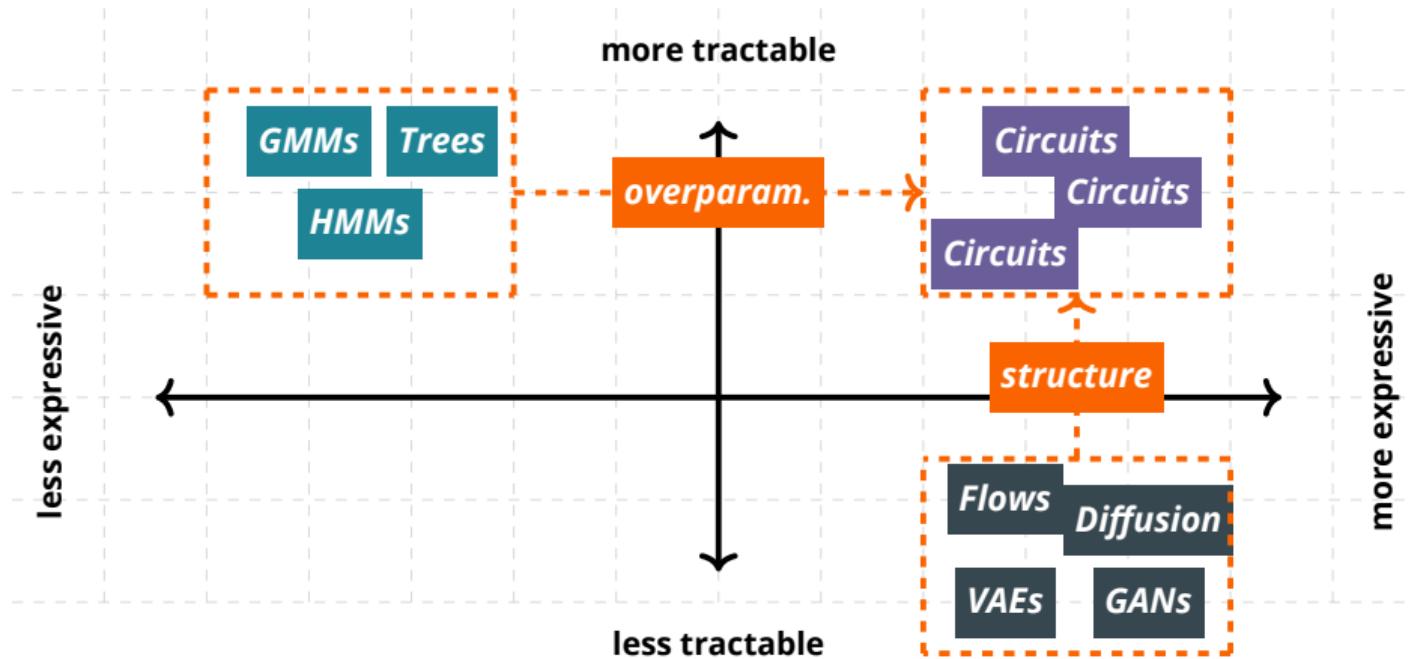
***circuits can be both expressive and tractable!***



*start simple...*



***then make it more expressive!***



***impose structure!***

## *Goal*

***“Can we design  
computational graphs  
that efficiently encode inference?”***

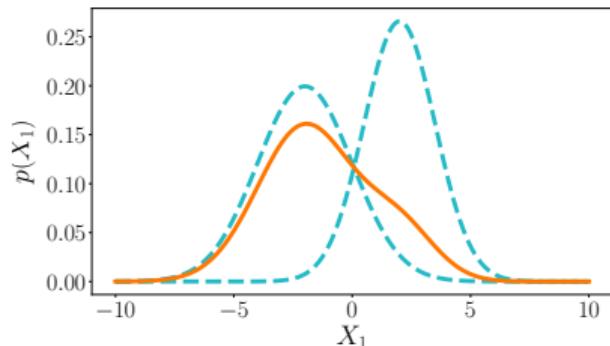
## *Goal*

***“Can we design  
computational graphs  
that efficiently encode inference?”***

⇒ *yes! with circuits!*

# GMMS

as computational graphs

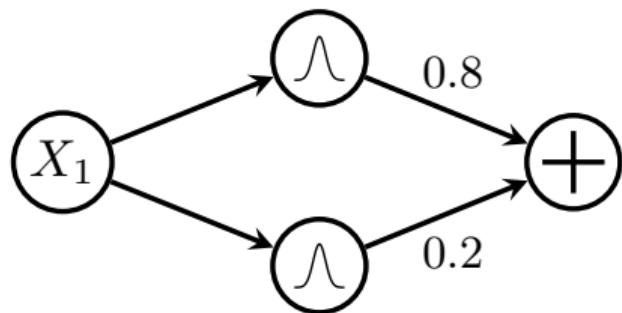


$$p(X) = w_1 \cdot p_1(X_1) + w_2 \cdot p_2(X_1)$$

⇒ translating inference to data structures...

# GMMs

as computational graphs

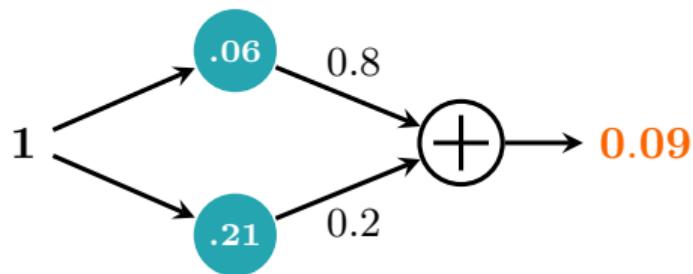


$$p(X_1) = 0.2 \cdot p_1(X_1) + 0.8 \cdot p_2(X_1)$$

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

# GMMS

as computational graphs

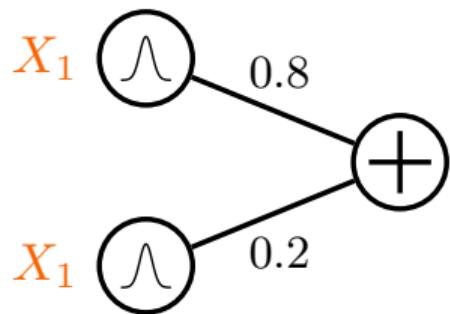


$$p(X = 1) = 0.2 \cdot p_1(X_1 = 1) + 0.8 \cdot p_2(X_1 = 1)$$

⇒ inference = feedforward evaluation

# GMMs

as computational graphs

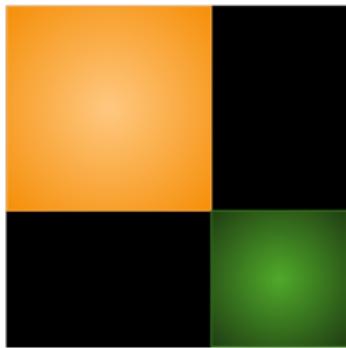
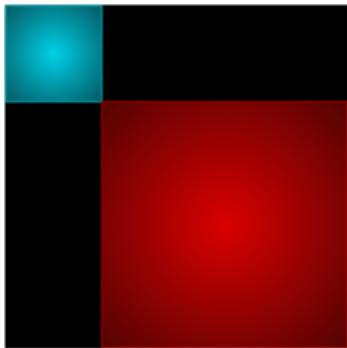


A simplified notation:

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

# GMMS

as computational graphs

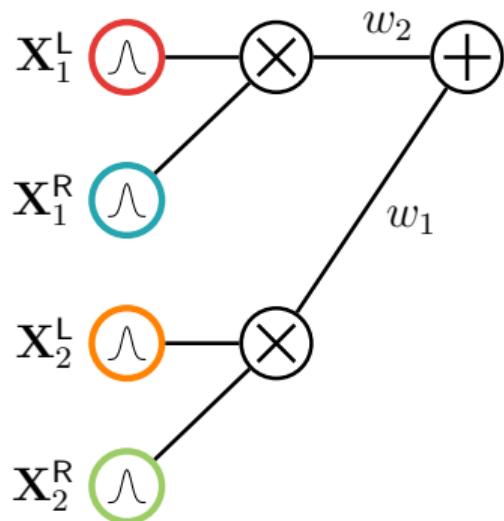


$$p(\mathbf{X}) = w_1 \cdot p_1(\mathbf{X}_1^L) \cdot p_1(\mathbf{X}_1^R) + \\ w_2 \cdot p_2(\mathbf{X}_2^L) \cdot p_2(\mathbf{X}_2^R)$$

⇒ 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*

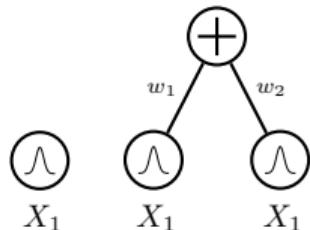
I. A simple tractable function is a circuit

$$\bigcirc \wedge \\ X_1$$

# Probabilistic Circuits (PCs)

A grammar for tractable computational graphs

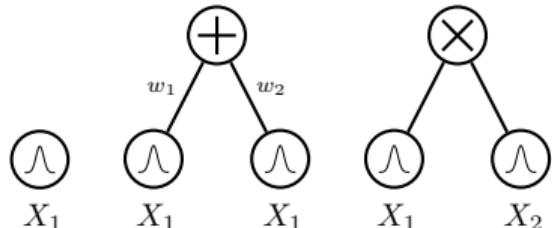
- I. A simple tractable function is a circuit
- II. A weighted combination of circuits is a circuit



# Probabilistic Circuits (PCs)

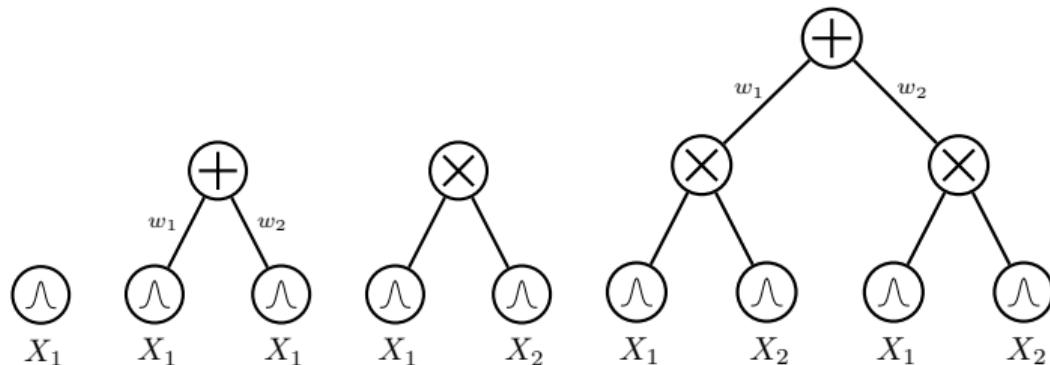
A grammar for tractable computational graphs

- 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



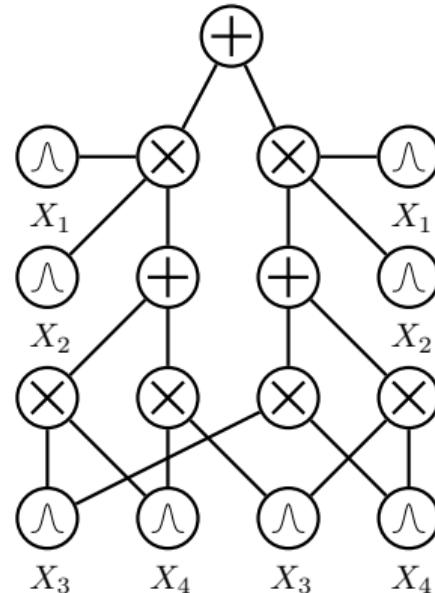
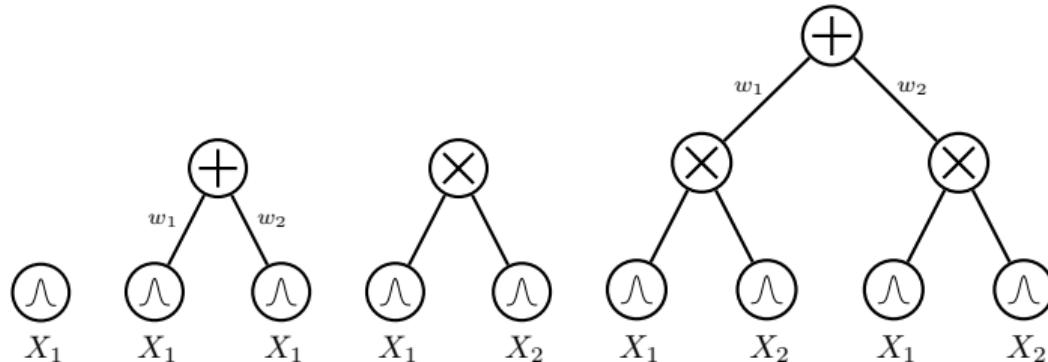
# Probabilistic Circuits (PCs)

A grammar for tractable computational graphs



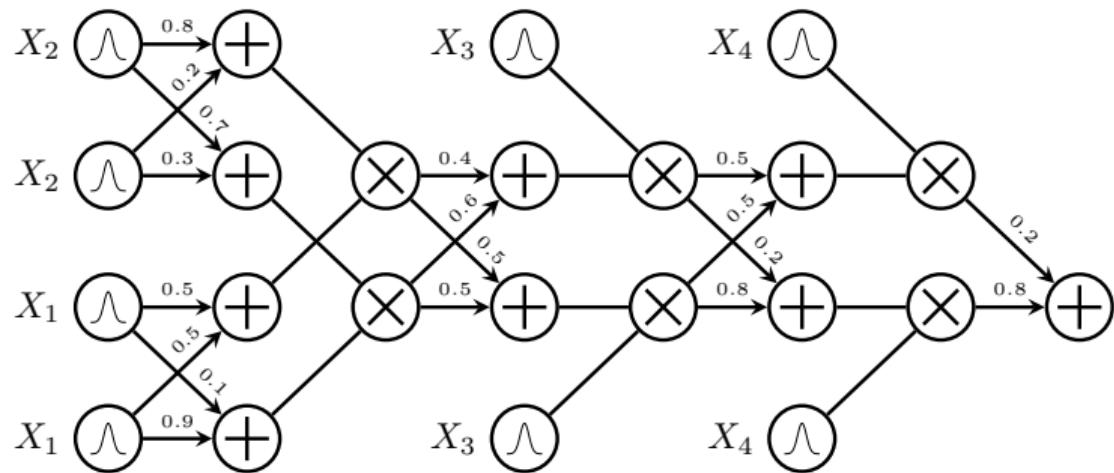
# Probabilistic Circuits (PCs)

A grammar for tractable computational graphs



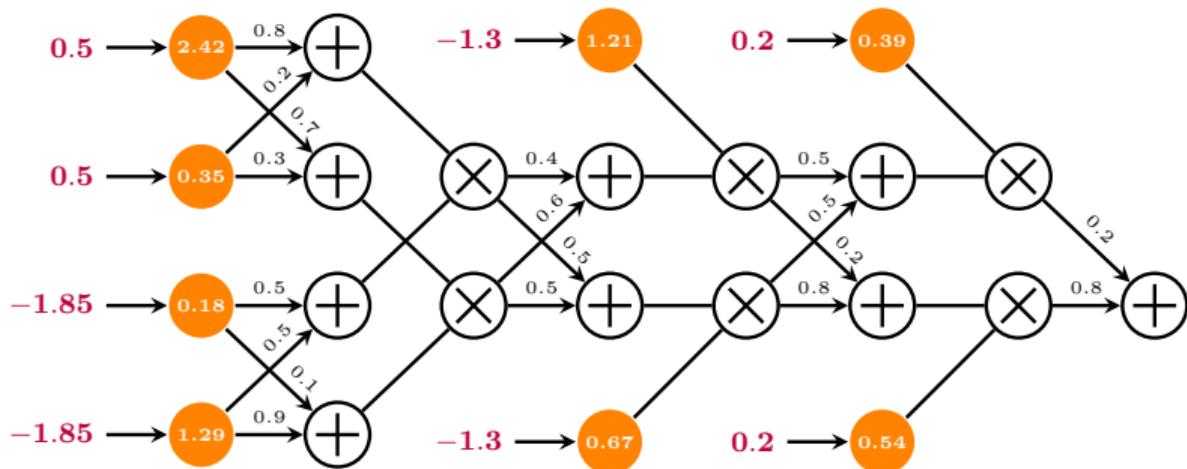
## Probabilistic queries = feedforward evaluation

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



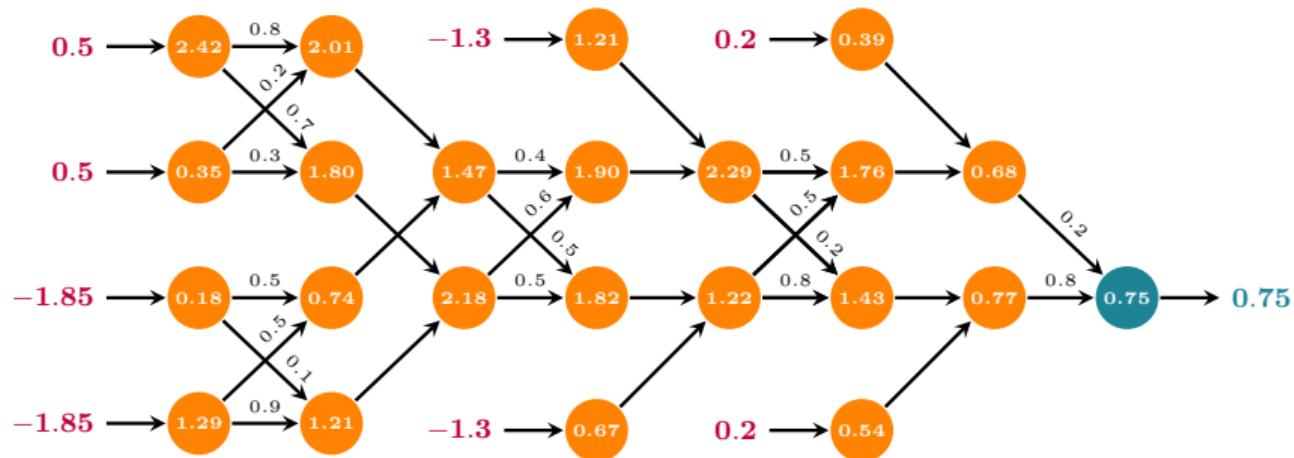
## Probabilistic queries = feedforward evaluation

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



## 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?**

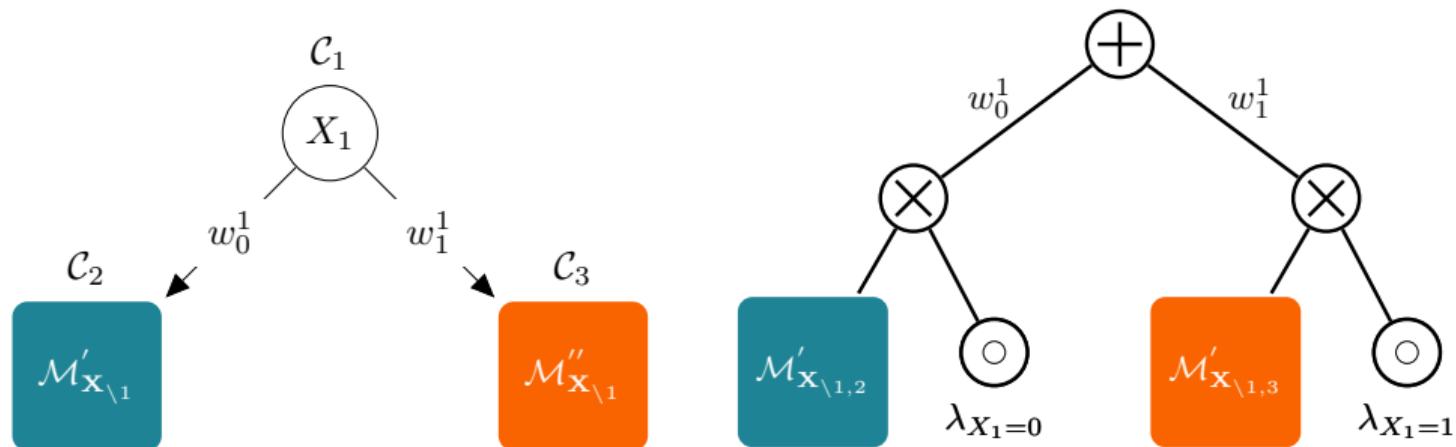
## **1. A grammar for tractable models**

One formalism to represent many probabilistic and logical models

⇒ #HMMs #Trees #XGBoost, Tensor Networks, ...  
O/BDDs, SDDs and other PGMs...

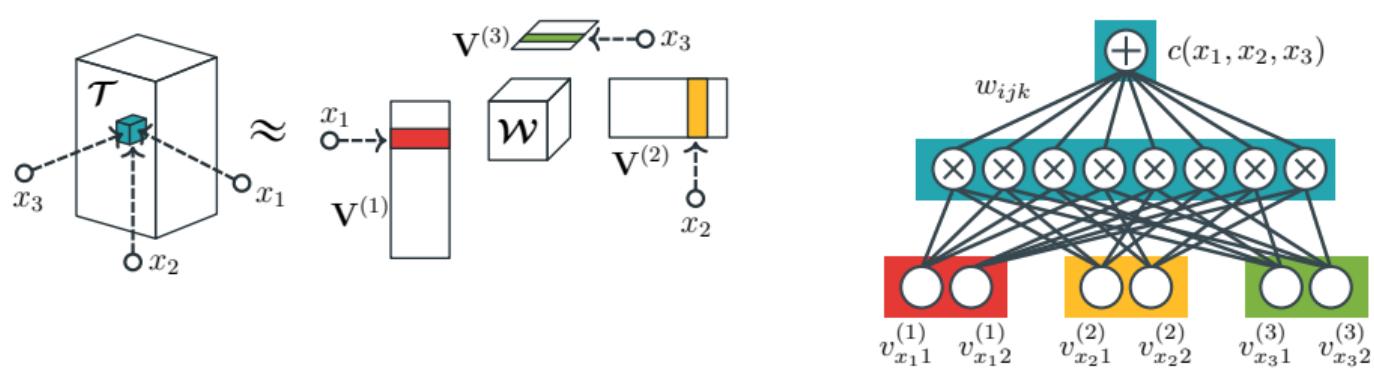
# **decision trees**

*but also OB/S/DDs, as circuits*



# *tensor factorizations*

*as circuits*



---

Loconte et al., "What is the Relationship between Tensor Factorizations and Circuits (and How Can We Exploit it)?", arXiv, 2024

# **...why PCs?**

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

One formalism to represent many probabilistic and logical models

⇒ #HMMs #Trees #XGBoost, Tensor Networks, ...  
O/BDDs, SDDs 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, ...  
O/BDDs, SDDs and other PGMs...

## **2. Expressiveness**

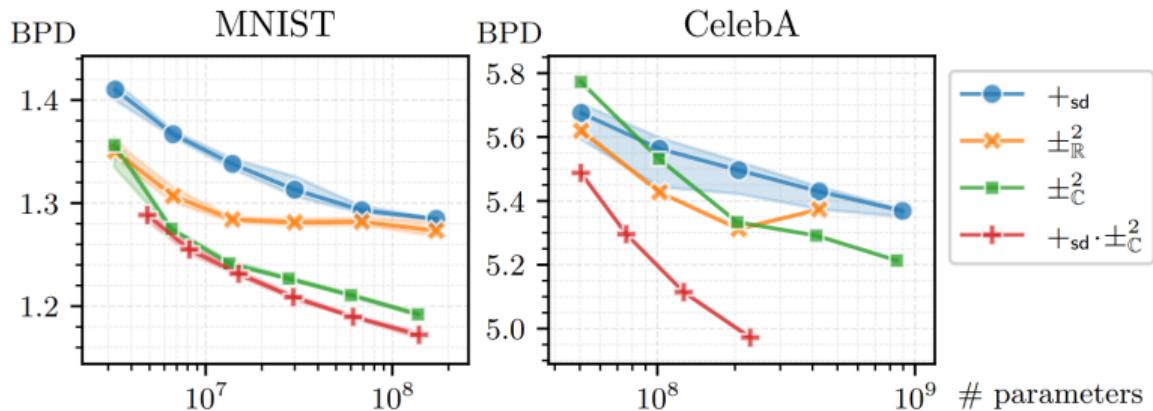
Competitive with intractable models, VAEs, Flow...#hierachical #mixtures #polynomials

# How expressive?

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

**competitive with Flows and VAEs!**

# How scalable?



*up to billions of parameters*

# ***How to build & learn probabilistic circuits?***



***learning & reasoning with circuits in pytorch***

<https://github.com/april-tools/cirkit>

```
1 from cirkit.templates import circuit_templates  
2  
3 symbolic_circuit = circuit_templates.image_data(  
4     (1, 28, 28),                      # The shape of MNIST  
5     region_graph='quad-graph',  
6     input_layer='categorical',          # input distributions  
7     sum_product_layer='cp',            # CP, Tucker, CP-T  
8     num_input_units=64,                # overparameterizing  
9     num_sum_units=64,  
10    sum_weight_param=circuit_templates.Parameterization(  
11        activation='softmax',  
12        initialization='normal'  
13    )  
14 )
```

```
1 from cirkit.pipeline import compile
2 circuit = compile(symbolic_circuit)
3
4 with torch.no_grad():
5     test_lls = 0.0
6     for batch, _ in test_dataloader:
7         batch = batch.to(device).unsqueeze(dim=1)
8         log_likelihoods = circuit(batch)
9         test_lls += log_likelihoods.sum().item()
10 average_ll = test_lls / len(data_test)
11 bpd = -average_ll / (28 * 28 * np.log(2.0))
12 print(f"Average LL: {average_ll:.3f}") # Average LL:
13     → -682.916
14 print(f"Bits per dim: {bpds:.3f}") # Bits per dim: 1.257
```

# **...why PCs?**

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

One formalism to represent many probabilistic and logical models

⇒ #HMMs #Trees #XGBoost, Tensor Networks, ...  
O/BDDs, SDDs and other PGMs...

## **2. Expressiveness**

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

# **...why PCs?**

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

One formalism to represent many probabilistic and logical models

⇒ #HMMs #Trees #XGBoost, Tensor Networks, ...  
O/BDDs, SDDs 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**

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

**property C**

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

**property D**

$\Rightarrow$  *tractable partition function*  
 $\Rightarrow$  also any *conditional* is tractable

# *Structural properties*

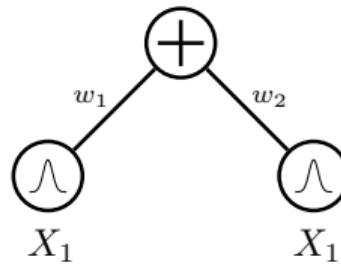
**smoothness**

the inputs of sum units are defined over the same variables

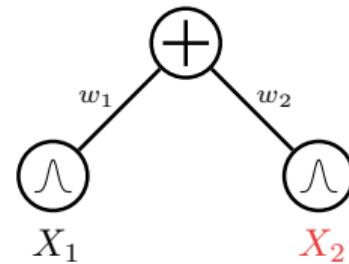
**decomposability**

**compatibility**

**determinism**



*smooth circuit*



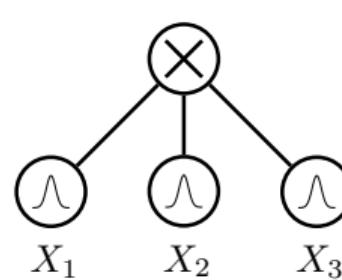
*non-smooth circuit*

# *Structural properties*

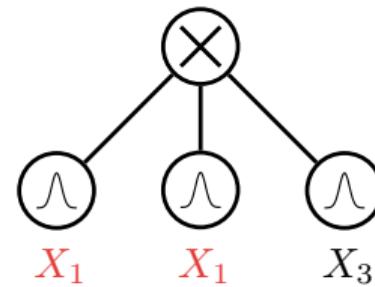
**smoothness**

the inputs of prod units are defined over disjoint variable sets

**decomposability**



**compatibility**

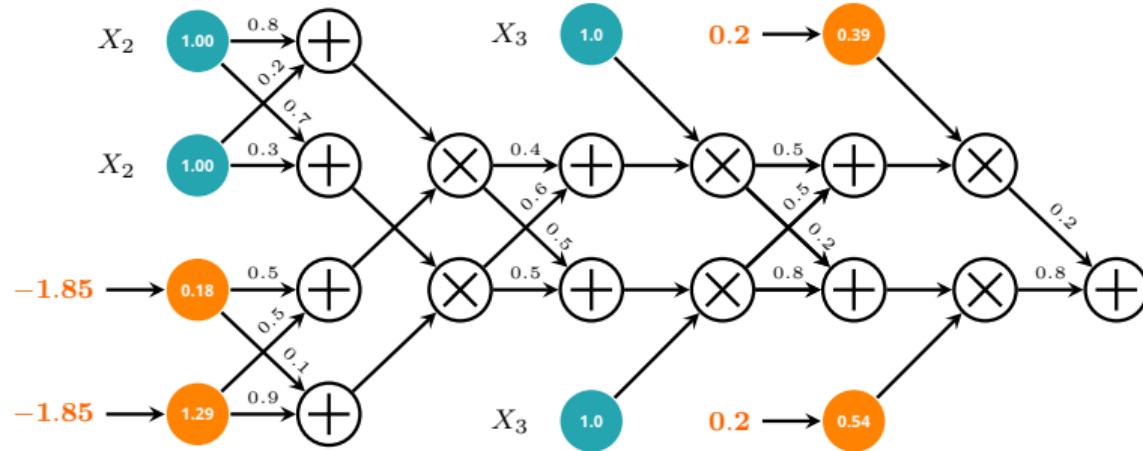


**determinism**

**decomposable circuit    non-decomposable circuit**

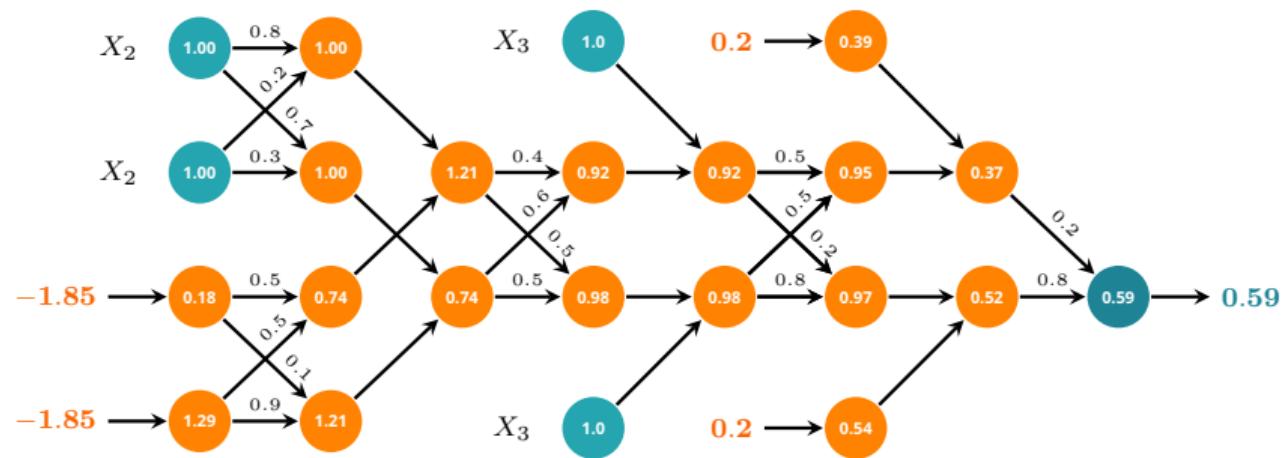
## Probabilistic queries = feedforward evaluation

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



**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:

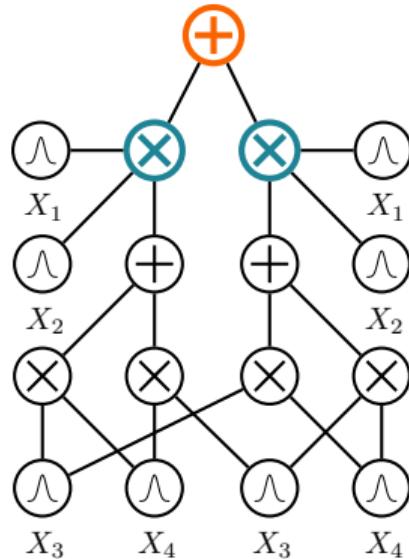
$$\int p(\mathbf{x}) d\mathbf{x}$$

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

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

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

$\Rightarrow$  integrals are “pushed down” to inputs

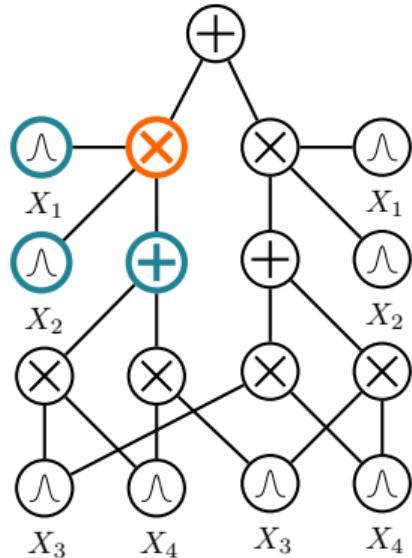


## ***smooth + decomposable*** circuits = ...

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



```
1 from cirkit.backend.torch.queries import IntegrateQuery
2 marginal_query = IntegrateQuery(circuit)
3
4 with torch.no_grad():
5     test_marginal_lls = 0.0
6
7     for batch, _ in test_dataloader:
8         batch = batch.to(device).unsqueeze(dim=1)
9         marginal_log_likelihoods = marginal_query(batch,
10             ↪ integrate_vars=vars_to_marginalize)
11         test_marginal_lls +=
12             ↪ marginal_log_likelihoods.sum().item()
13
14     marg_ll = test_marginal_lls / len(data_test)
15     print(f"marg LL: {marg_ll:.3f}") # marg LL:: -378.417
```

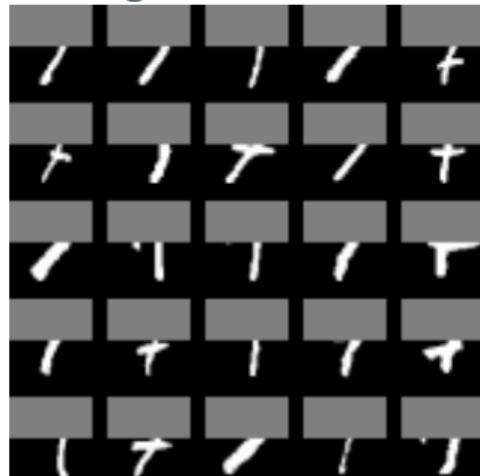
# **Tractable inference on PCs**

*Einsum networks*

Original

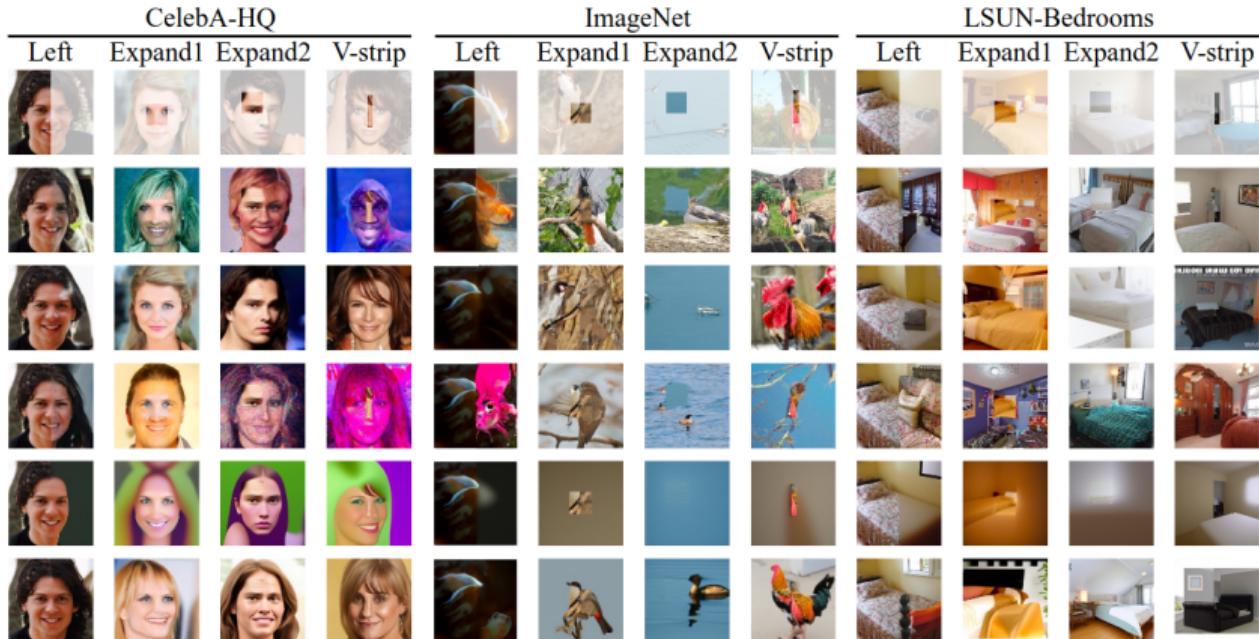


Missing



Conditional sample





# **Which structural properties**

*for complex reasoning*



*smooth + decomposable*



???????



???????

# *General expectations*

Integrals involving two or more functions:

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

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



## *General expectations*

Integrals involving two or more functions:

$$\int \textcolor{orange}{p}(x) \textcolor{teal}{f}(x) dX$$

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



# *Structural properties*

*smoothness*

*decomposability*

*compatibility*

*determinism*

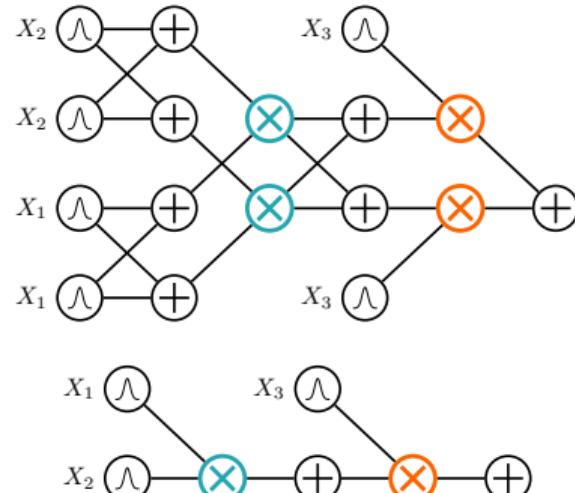
# *Structural properties*

**smoothness**

**decomposability**

**compatibility**

**determinism**



*compatible circuits*

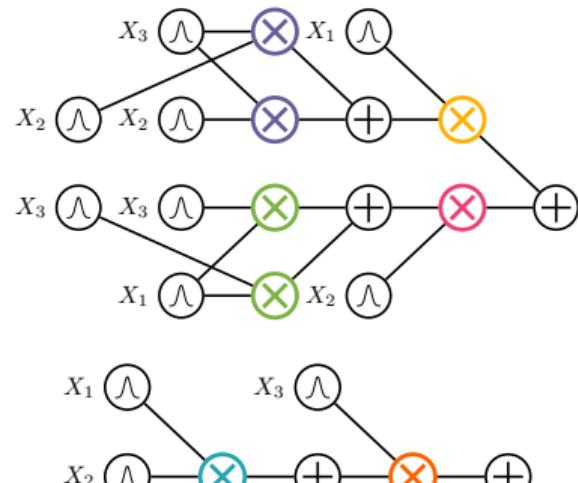
# *Structural properties*

**smoothness**

**decomposability**

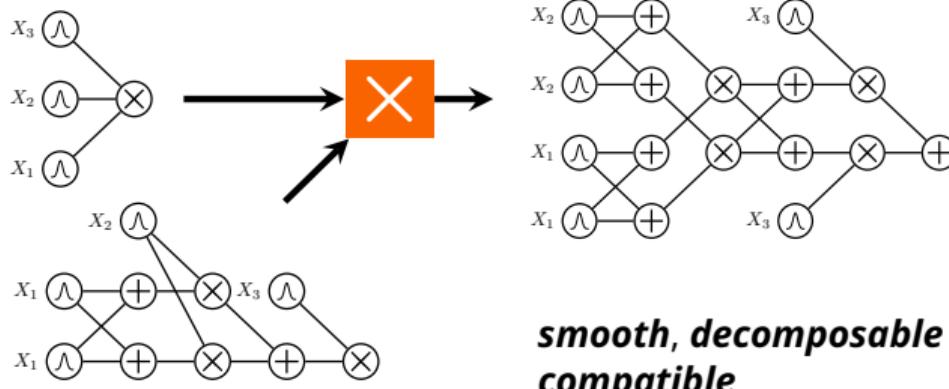
**compatibility**

**determinism**



*non-compatible circuits*

# Tractable products



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

---

# Semantic Probabilistic Layers for Neuro-Symbolic Learning

---

**Kareem Ahmed**

CS Department  
UCLA

[ahmedk@cs.ucla.edu](mailto:ahmedk@cs.ucla.edu)

**Stefano Teso**

CIMeC and DISI  
University of Trento

[stefano.teso@unitn.it](mailto:stefano.teso@unitn.it)

**Kai-Wei Chang**

CS Department  
UCLA

[kwchang@cs.ucla.edu](mailto:kwchang@cs.ucla.edu)

**Guy Van den Broeck**

CS Department  
UCLA

[guyvdb@cs.ucla.edu](mailto:guyvdb@cs.ucla.edu)

**Antonio Vergari**

School of Informatics  
University of Edinburgh

[avergari@ed.ac.uk](mailto:avergari@ed.ac.uk)

*circuit products for reliable NeSy: tomorrow*

```
1 from cirkit.templates import circuit_templates  
2  
3 symbolic_circuit = circuit_templates.image_data(  
4     (1, 28, 28),                      # The shape of MNIST  
5     region_graph='quad-graph',  
6     input_layer='categorical',          # input distributions  
7     sum_product_layer='cp',            # CP, Tucker, CP-T  
8     num_input_units=64,                # overparameterizing  
9     num_sum_units=64,  
10    sum_weight_param=circuit_templates.Parameterization(  
11        activation='softmax',  
12        initialization='normal'  
13    )  
14 )
```

# *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...***

*which loss?*

*how to learn normalized weights?*

*how to exploit structural properties?*

# **maximum likelihood**

*the go-to objective in ProbML*

Given a dataset  $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$  and your parametric model  $p_\theta(\mathbf{X})$  solve

$$\hat{\theta}_{\text{ML}} = \max_{\theta} \prod_{i=1}^N p_\theta(\mathbf{x}^{(i)}) = \min_{\theta} - \sum_{i=1}^N \log p_\theta(\mathbf{x}^{(i)})$$

$\Rightarrow$  minimize the negative log-likelihood (NLL)

```
1 from cirkit.templates import circuit_templates  
2  
3 symbolic_circuit = circuit_templates.image_data(  
4     (1, 28, 28),                      # The shape of MNIST  
5     region_graph='quad-graph',  
6     input_layer='categorical',          # input distributions  
7     sum_product_layer='cp',            # CP, Tucker, CP-T  
8     num_input_units=64,                # overparameterizing  
9     num_sum_units=64,  
10    sum_weight_param=circuit_templates.Parameterization(  
11        activation='softmax',  
12        initialization='normal'  
13    )  
14 )
```

# **which parameters?**

*how to reparameterize circuits*

Input distributions.

Sum unit parameters.

# **which parameters?**

*how to reparameterize circuits*

**Input distributions.** Each input can be a different parametric distribution

⇒ *Bernoullis, Categoricals, Gaussians, exponential families, small NNs, ...*

**Sum unit parameters.**

# which parameters?

how to reparameterize circuits

**Input distributions.** Each input can be a different parametric distribution

**Sum unit parameters.** Enforce them to be non-negative, i.e.,  $w_i \geq 0$  but unnormalized

$$w_i = \exp(\alpha_i), \quad \alpha_i \in \mathbb{R}, \quad i = 1, \dots, K$$

and renormalize the loss

$$\min_{\theta} - \left( \sum_{i=1}^N \log \tilde{p}_{\theta}(\mathbf{x}^{(i)}) - \log \int \tilde{p}_{\theta}(\mathbf{x}^{(i)}) d\mathbf{X} \right)$$

or just renormalize the weights, i.e.,  $\sum_i w_i = 1$

$$\mathbf{w} = \text{softmax}(\boldsymbol{\alpha}), \quad \boldsymbol{\alpha} \in \mathbb{R}^K$$

```
1 from cirkit.templates import circuit_templates  
2  
3 symbolic_circuit = circuit_templates.image_data(  
4     (1, 28, 28),                      # The shape of MNIST  
5     region_graph='quad-graph',  
6     input_layer='categorical',          # input distributions  
7     sum_product_layer='cp',            # CP, Tucker, CP-T  
8     num_input_units=64,                # overparameterizing  
9     num_sum_units=64,  
10    sum_weight_param=circuit_templates.Parameterization(  
11        activation='softmax',  
12        initialization='normal'  
13    )  
14 )
```

# **Probabilistic Circuits (PCs)**

*the unit-wise definition*

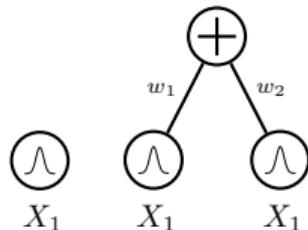
I. A simple tractable function is a circuit

$$\bigcirc \wedge \\ X_1$$

# Probabilistic Circuits (PCs)

*the unit-wise definition*

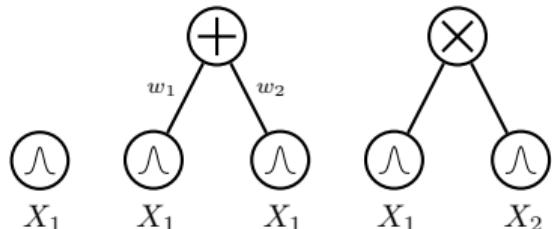
- I. A simple tractable function is a circuit
- II. A weighted combination of circuits is a circuit



# Probabilistic Circuits (PCs)

*the unit-wise definition*

- 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)**

*the layer-wise definition*

- I. A set of tractable functions is a circuit layer



# Probabilistic Circuits (PCs)

*the layer-wise definition*

- I. A set of tractable functions is a circuit layer
- II. A linear projection of a layer is a circuit layer

$$c(\mathbf{x}) = \mathbf{W}l(\mathbf{x})$$



# Probabilistic Circuits (PCs)

*the layer-wise definition*

- I. A set of tractable functions is a circuit layer
- II. A linear projection of a layer is a circuit layer

$$c(\mathbf{x}) = \mathbf{W}\mathbf{l}(\mathbf{x})$$



# Probabilistic Circuits (PCs)

the layer-wise definition

- I. A set of tractable functions is a circuit layer
- II. A linear projection of a layer is a circuit layer
- III. The product of two layers is a circuit layer

$$c(\mathbf{x}) = \mathbf{l}(\mathbf{x}) \odot \mathbf{r}(\mathbf{x}) \quad // \text{Hadamard}$$



# Probabilistic Circuits (PCs)

*the layer-wise definition*

- I. A set of tractable functions is a circuit layer
- II. A linear projection of a layer is a circuit layer
- III. The product of two layers is a circuit layer

$$c(\mathbf{x}) = \mathbf{l}(\mathbf{x}) \odot \mathbf{r}(\mathbf{x}) \quad // \text{Hadamard}$$

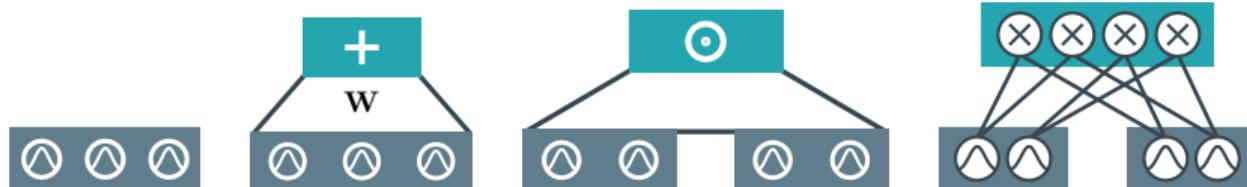


# Probabilistic Circuits (PCs)

the layer-wise definition

- I. A set of tractable functions is a circuit layer
- II. A linear projection of a layer is a circuit layer
- III. The product of two layers is a circuit layer

$$c(\mathbf{x}) = \text{vec}(\mathbf{l}(\mathbf{x})\mathbf{r}(\mathbf{x})^\top) \quad // \text{ Kronecker}$$

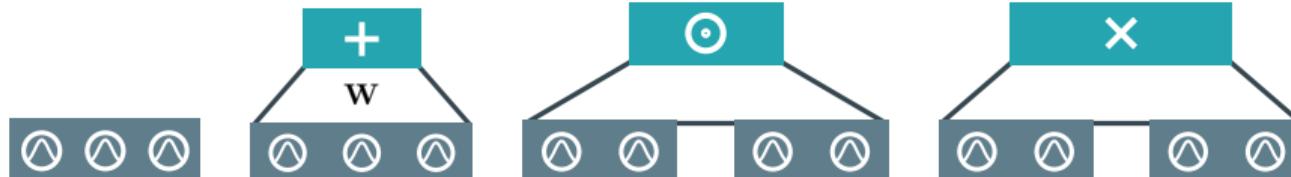


# Probabilistic Circuits (PCs)

*the layer-wise definition*

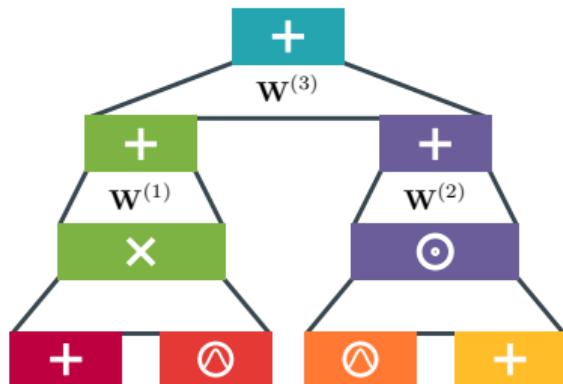
- I. A set of tractable functions is a circuit layer
- II. A linear projection of a layer is a circuit layer
- III. The product of two layers is a circuit layer

$$c(\mathbf{x}) = \text{vec}(\mathbf{l}(\mathbf{x})\mathbf{r}(\mathbf{x})^\top) \quad // \text{ Kronecker}$$



# Probabilistic Circuits (PCs)

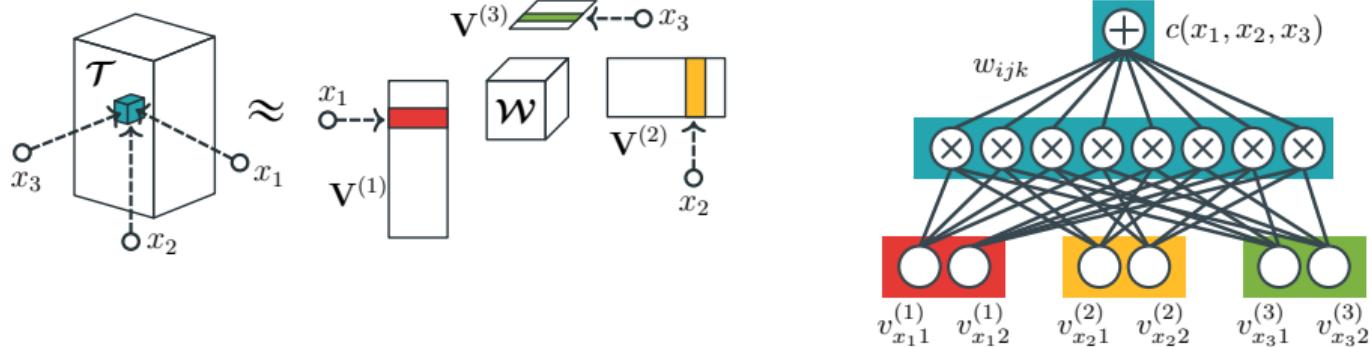
the layer-wise definition



- I. A set of tractable functions is a circuit layer
  - II. A linear projection of a layer is a circuit layer
  - III. The product of two layers is a circuit layer
- stack layers to build a deep circuit!**

# *circuits layers*

*as tensor factorizations*

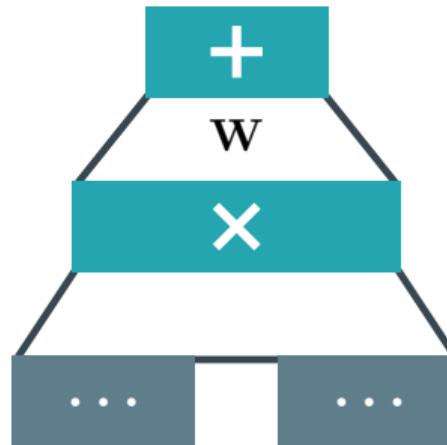
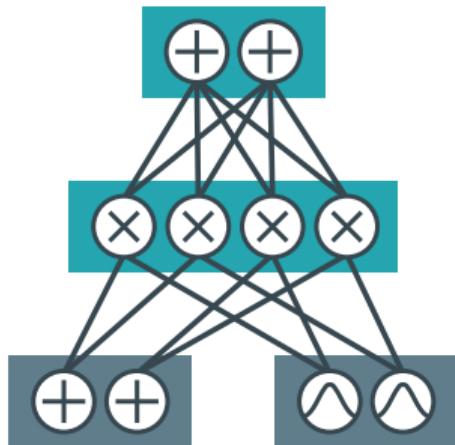


---

Loconte et al., "What is the Relationship between Tensor Factorizations and Circuits (and How Can We Exploit it)?", arXiv, 2024

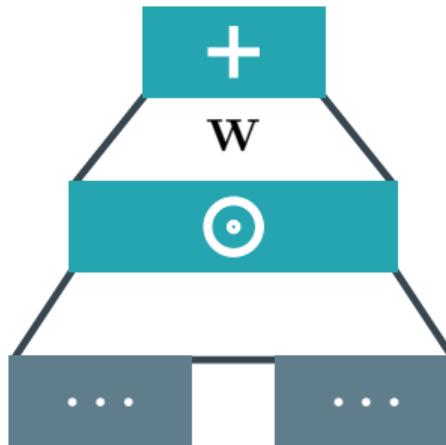
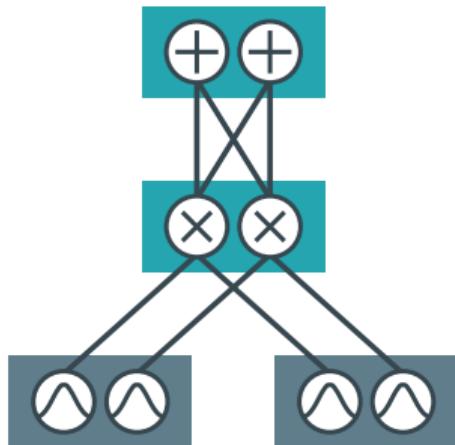
# **more layers**

*Tucker decomposition*

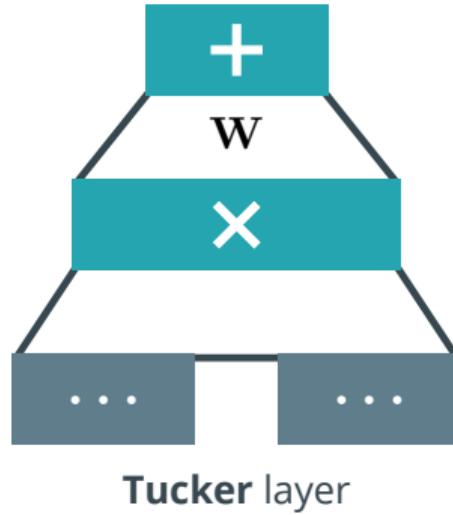
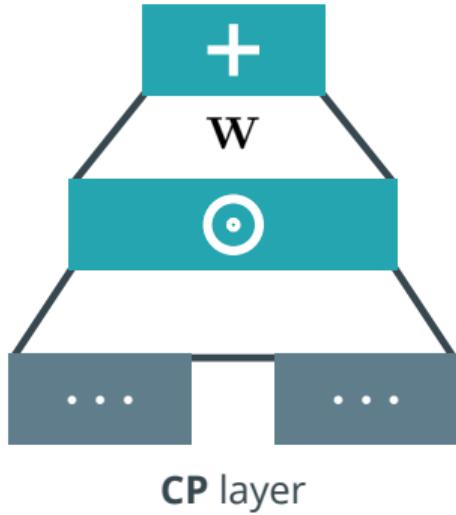


# **more layers**

*Candecomp Parafac (CP) decomposition*



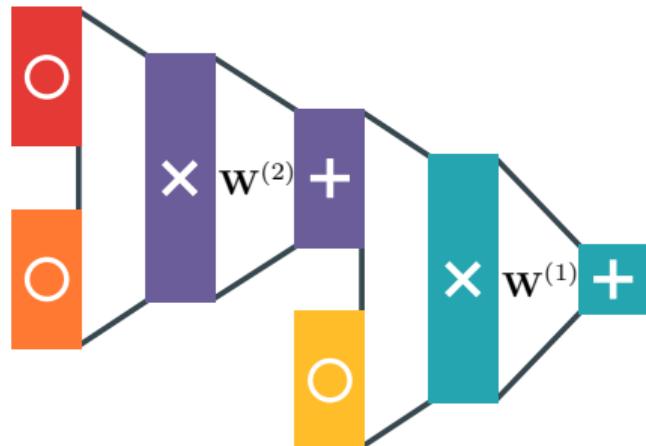
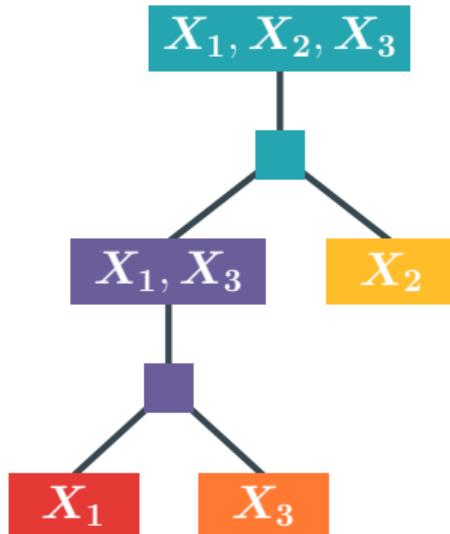
## *more layers*



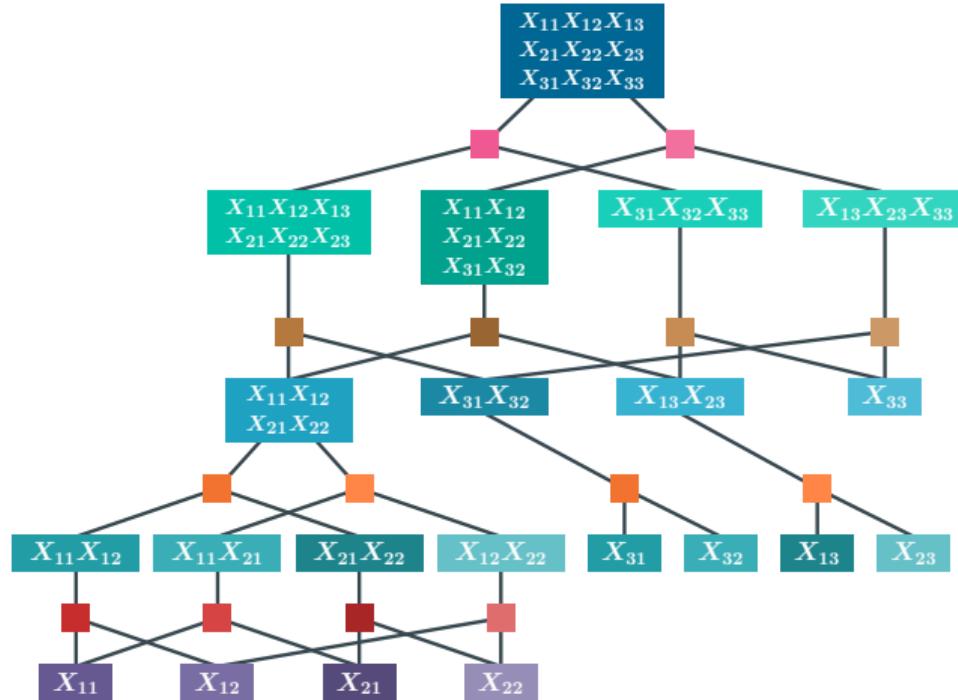
```
1 from cirkit.templates import circuit_templates  
2  
3 symbolic_circuit = circuit_templates.image_data(  
4     (1, 28, 28),  
5     region_graph='quad-graph',  
6     input_layer='categorical',    # input distributions  
7     sum_product_layer='cp',      # CP, Tucker, CP-T  
8     num_input_units=64,          # overparameterizing  
9     num_sum_units=64,  
10    sum_weight_param=circuit_templates.Parameterization(  
11        activation='softmax',  
12        initialization='normal'  
13    )  
14 )
```

# *region graphs*

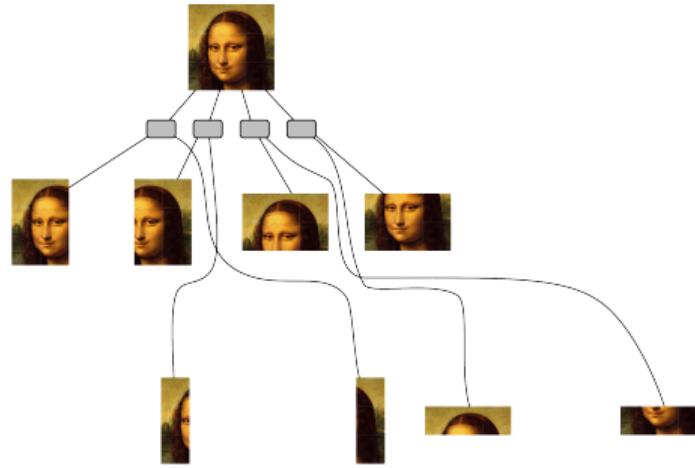
*a template for smooth&decomposable PCs*

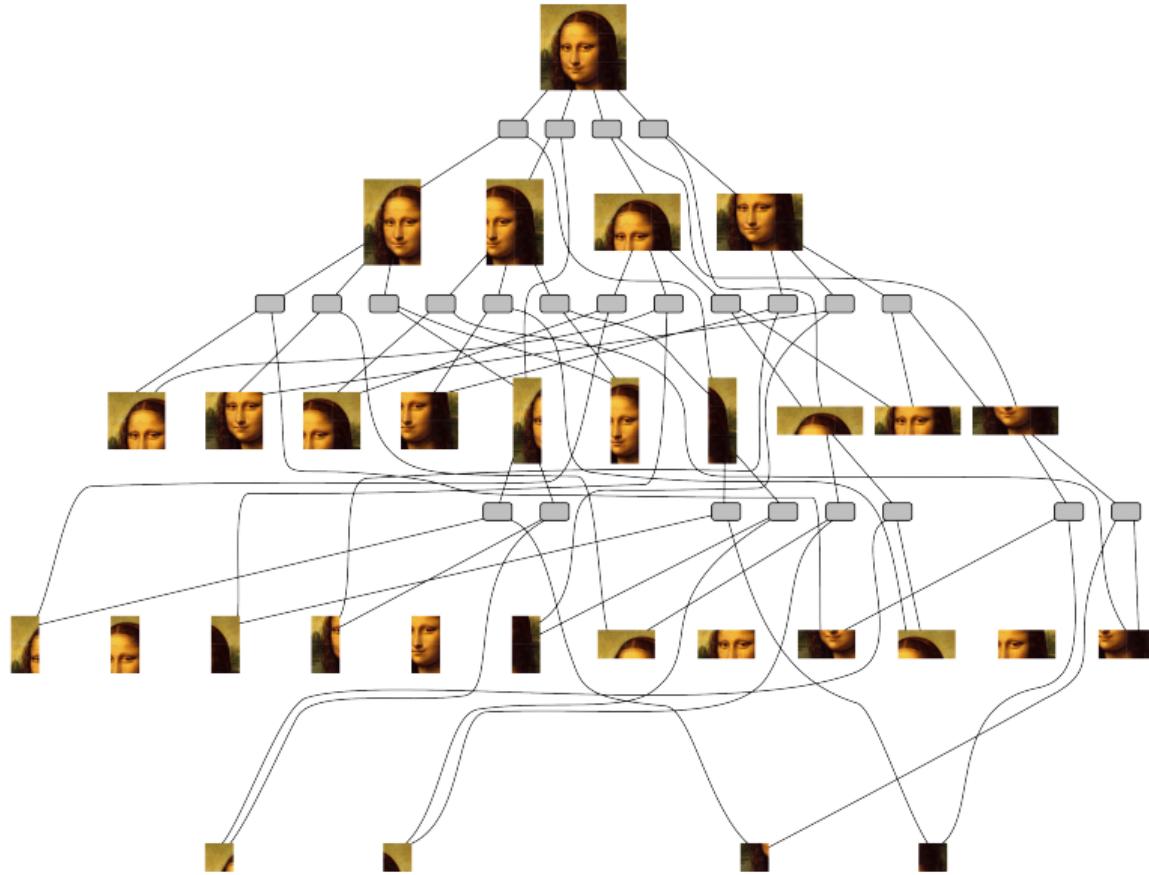


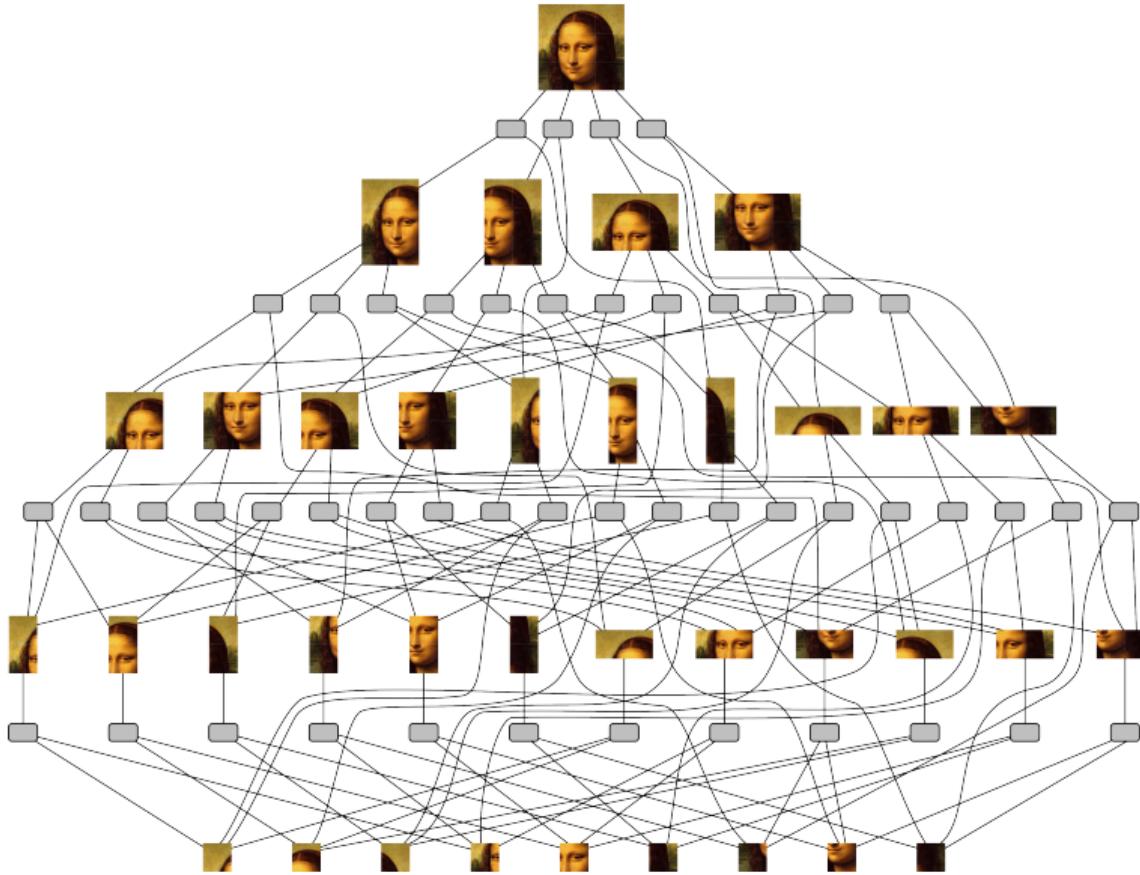
# *which region graph?*







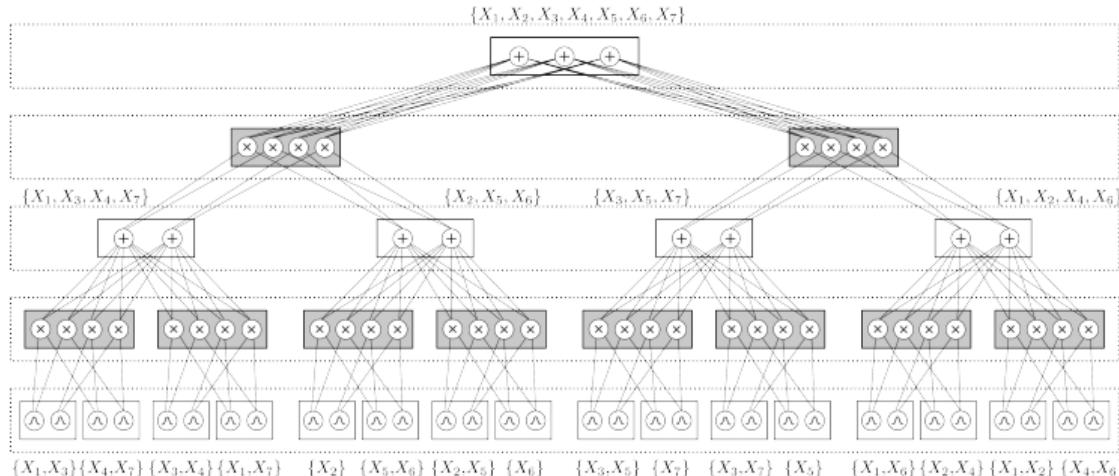




# *random regions graphs*

The “no-learning” option

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



The screenshot shows a Jupyter Notebook interface. At the top, there's a header bar with a file icon, a dropdown menu labeled "main", and the path "cirkit / notebooks / region-graphs-and-parametrisation.ipynb". To the right of the path is a search bar with the placeholder "Go to file" and a refresh icon. Below the header, a commit message from "loreloc" is displayed: "updated notebooks with respect to API changes" with a timestamp "e3e7e80 · 2 days ago" and a clock icon. Underneath the commit message, there are buttons for "Preview", "Code", and "Blame", followed by the statistics "1082 lines (1082 loc) · 793 KB". On the far right of the header, there are buttons for "Raw", "Copy", and "Download".

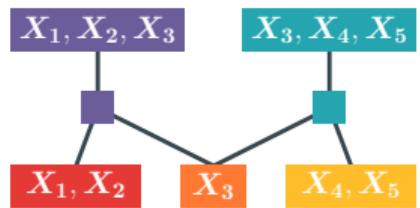
# Notebook on Region Graphs and Sum Product Layers

## Goals

By the end of this tutorial you will:

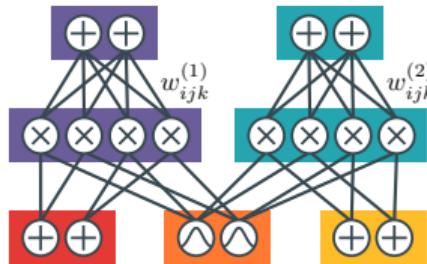
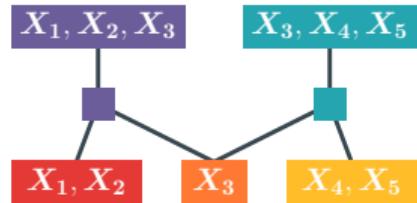
- know what a region graph is
- know how to choose between region graphs for your circuit
- understand how to parametrize a circuit by choosing a sum product layer
- build circuits to tractably estimate a probability distribution over images<sup>1</sup>

# *learning recipe*



1) Build a *region graph*

# *learning recipe*

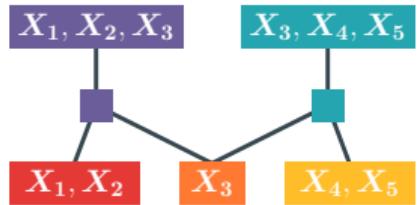


1) Build a *region graph*

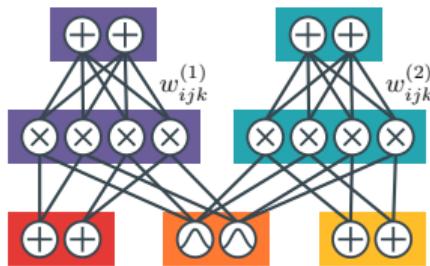
2) Overparameterize

- 2.1) pick a (composite) layer type**
- 2.2) choose how many units per layer**

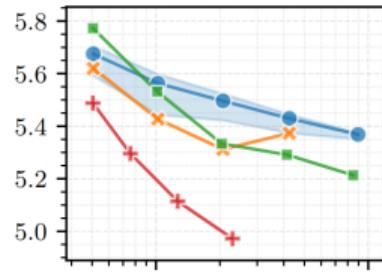
# *learning recipe*



1) Build a *region graph*



2) Overparameterize



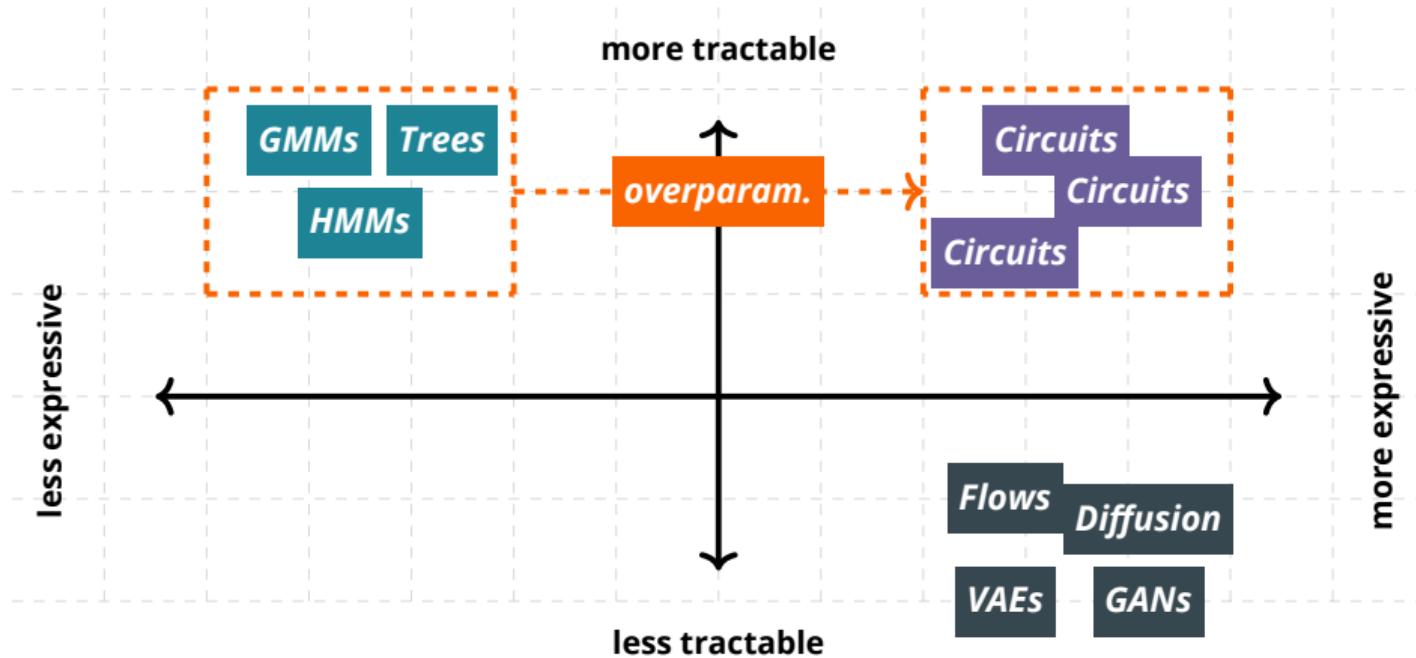
3) Learn parameters

*use any optimizer in pytorch*

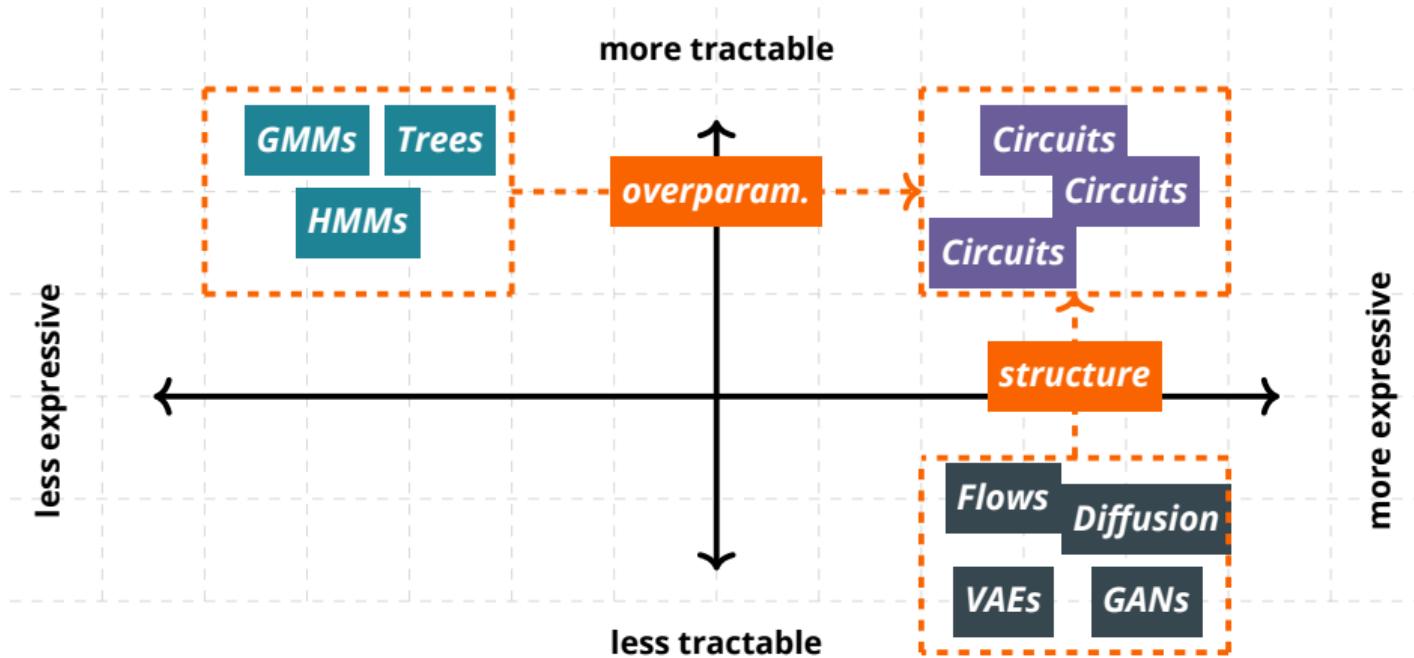


***learning & reasoning with circuits in pytorch***

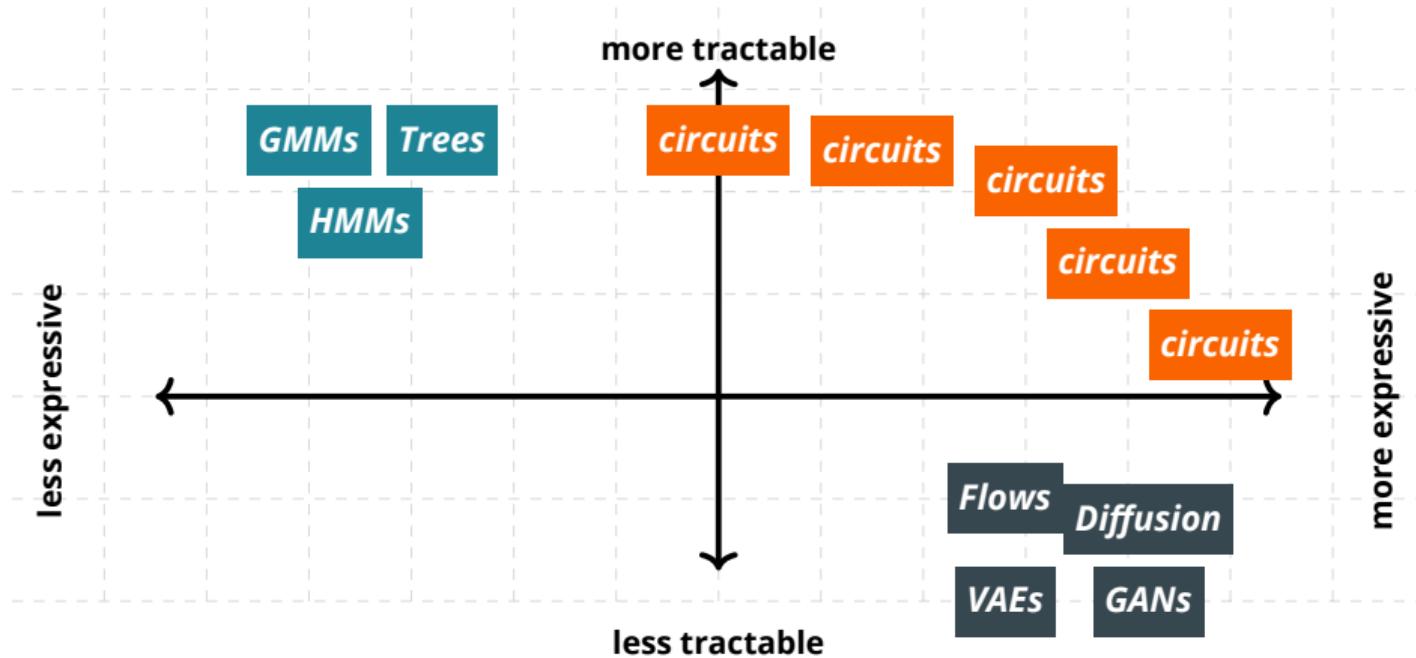
<https://github.com/april-tools/cirkit>



***make it more expressive!***



***impose structure!***



***navigate the spectrum!***