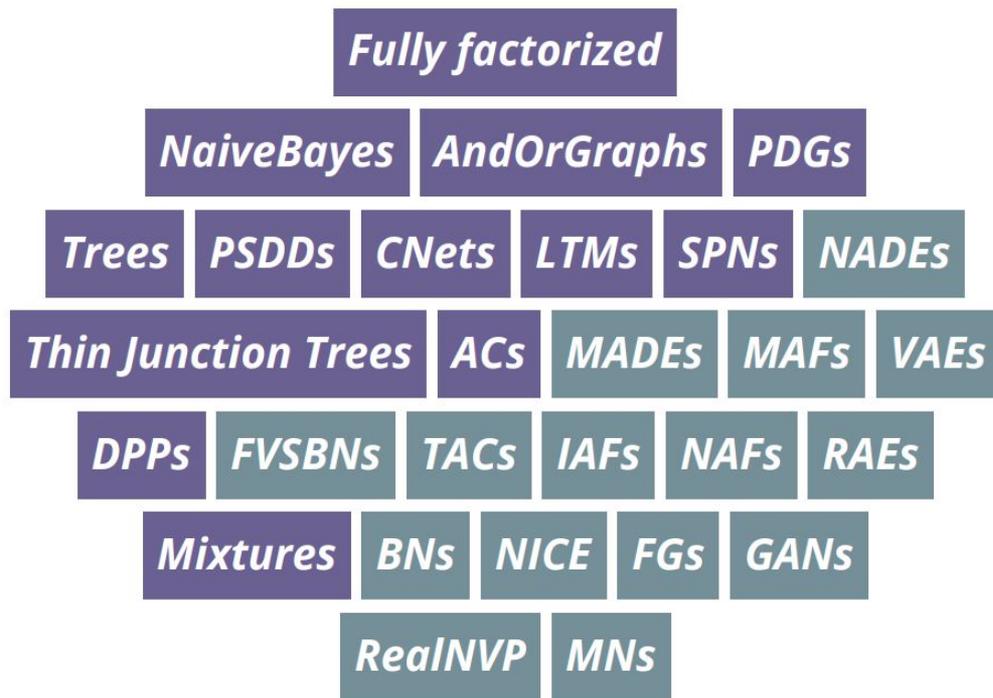# Tractable Probabilistic Circuits

Guy Van den Broeck
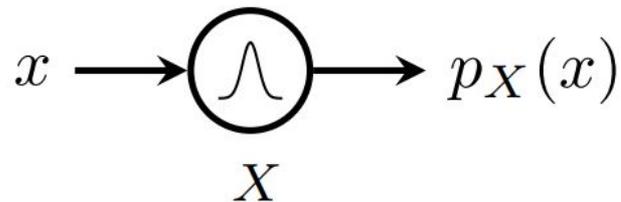
Beyond Bayes: Paths Towards Universal Reasoning Systems  - Jul 21, 2022

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
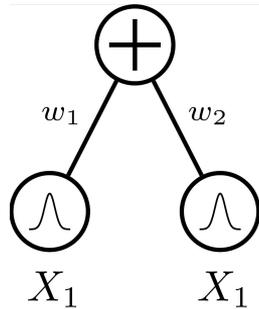
# Probabilistic circuits

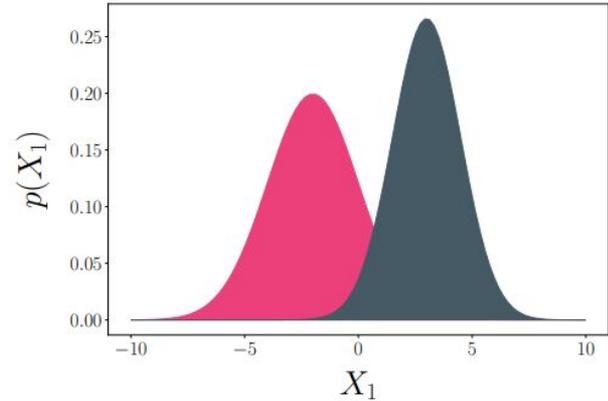*computational graphs* that recursively define distributions



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

$$\Rightarrow$$

*mixtures*

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

# Probabilistic circuits
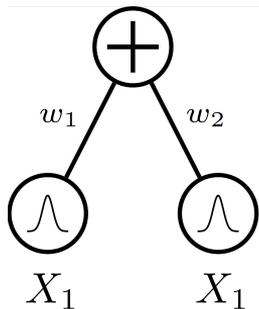
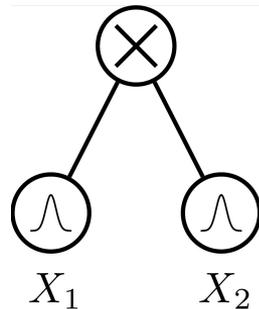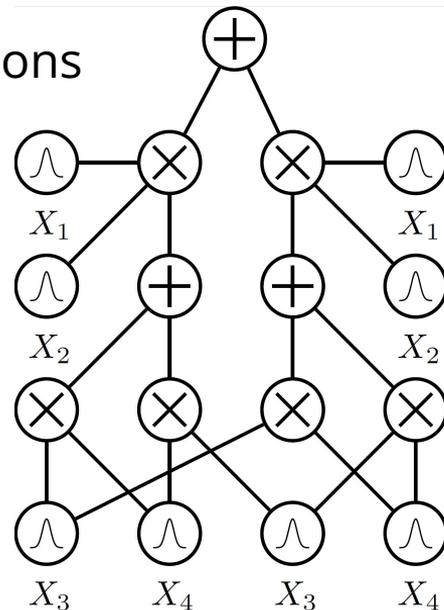*computational graphs* that recursively define distributions



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

$$\Rightarrow$$

*mixtures*

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

$$\Rightarrow$$

*factorizations*

# 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

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

# 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**

Darwiche and Marquis, "A Knowledge Compilation Map", 2002

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

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

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

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

$\implies$ *integrals are "pushed down" to children*

[Darwiche & Marquis JAIR 2001, Poon & Domingos UAI11]

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

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

$$\int \int \int p(\mathbf{x}, \mathbf{y}, \mathbf{z})d\mathbf{x}d\mathbf{y}d\mathbf{z} =$$

$$= \int \int \int p(\mathbf{x})p(\mathbf{y})p(\mathbf{z})d\mathbf{x}d\mathbf{y}d\mathbf{z} =$$

$$= \int p(\mathbf{x})d\mathbf{x} \int p(\mathbf{y})d\mathbf{y} \int p(\mathbf{z})d\mathbf{z}$$

$\Longrightarrow$  *integrals decompose into easier ones*

Forward pass evaluation for MAR

$\Longrightarrow$ *linear in circuit size!*

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

- ■ leafs over $X_1$ and $X_3$ output $Z_i = \int p(x_i) dx_i$

  $\Longrightarrow$ *for normalized leaf distributions:* 1.0

- ■ leafs over $X_2$ and $X_4$ output *EVI*

  feedforward evaluation (bottom-up)

Forward pass evaluation for MAR

$\Rightarrow$ *linear in circuit size!*

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

■ leafs over $X_1$ and $X_3$ output $\mathbf{Z}_i = \int p(x_i) dx_i$

$\Rightarrow$ *for normalized leaf distributions:* $1.0$

■ leafs over $X_2$ and $X_4$ output *EVI*

■ feedforward evaluation (bottom-up)

# Learning Expressive Probabilistic Circuits

## Hidden Chow-Liu Trees



Learned **CLT structure** captures strong pairwise dependencies

Anji Liu and Guy Van den Broeck. Tractable Regularization of Probabilistic Circuits, *NeurIPS*, 2021.

# Learning Expressive Probabilistic Circuits

## Hidden Chow-Liu Trees



Learned **CLT structure** captures strong pairwise dependencies

Learned **HCLT structure**

correlations

6 Variables

**Compile** into an equivalent PC

# Learning Expressive Probabilistic Circuits

## Hidden Chow-Liu Trees



Learned **CLT structure** captures strong pairwise dependencies

Learned **HCLT structure**

correlations

6 Variables

**Compile** into an equivalent PC

Mini-batch Stochastic **Expectation Maximization**

Anji Liu and Guy Van den Broeck. Tractable Regularization of Probabilistic Circuits, *NeurIPS*, 2021.

# Lossless Data Compression



**Data**        Encode      **Bitstream**      Decode      **Reconstructed data**

Expressive probabilistic model $p(\boldsymbol{x})$ $\Rightarrow$ Determines the theoretical limit of compression rate

$+$

Efficient coding algorithm $\Rightarrow$ How close we can approach the theoretical limit

Anji Liu, Stephan Mandt and Guy Van den Broeck. Lossless Compression with Probabilistic Circuits, 2021.

# Lossless Neural Compression with Probabilistic Circuits

**Data**　　　　　　　**Bitstream**　　　**Reconstructed data**



Probabilistic Circuits

- Expressive　　$\rightarrow$　SoTA likelihood on MNIST.

- Fast　　$\rightarrow$　Time complexity of en/decoding is **O( |p| log(D) )**, where D is the # variables and |p| is the size of the PC.

Arithmetic Coding:

$$p(X_1 < x_1)$$
$$p(X_1 \leq x_1)$$
$$p(X_2 < x_2 | x_1)$$
$$p(X_2 \leq x_2 | x_1)$$
$$p(X_3 < x_3 | x_1, x_2)$$
$$p(X_3 \leq x_3 | x_1, x_2)$$
$$\vdots$$

Anji Liu, Stephan Mandt and Guy Van den Broeck. Lossless Compression with Probabilistic Circuits, 2021.

# Lossless Neural Compression with Probabilistic Circuits

SoTA compression rates

| Dataset | HCLT (ours) | IDF | BitSwap | BB-ANS | JPEG2000 | WebP | McBits |
|---|---|---|---|---|---|---|---|
| MNIST | **1.24** (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) | **3.35** (3.28) | 3.69 (3.66) | 3.93 | 4.62 | (3.72) |
| EMNIST (Letter) | **1.84** (1.80) | 2.02 (1.95) | 1.90 (1.84) | 2.29 (2.26) | 3.62 | 3.31 | (3.12) |
| EMNIST (ByClass) | **1.89** (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++ | **3.24** | 4.10 | 3.81 |
| PC+IDF | 3.28 | **3.99** | **3.71** |

Anji Liu, Stephan Mandt and Guy Van den Broeck. Lossless Compression with Probabilistic Circuits, 2021.

# PC Learners keep getting better!    *... stay tuned ...*

Table 1: Density estimation performance on MNIST-family datasets in test set bpd.

| Dataset | Sparse PC (ours) | HCLT | RatSPN | IDF | BitSwap | BB-ANS | McBits |
|---|---|---|---|---|---|---|---|
| MNIST | **1.14** | 1.20 | 1.67 | 1.90 | 1.27 | 1.39 | 1.98 |
| EMNIST(MNIST) | **1.52** | 1.77 | 2.56 | 2.07 | 1.88 | 2.04 | 2.19 |
| EMNIST(Letters) | **1.58** | 1.80 | 2.73 | 1.95 | 1.84 | 2.26 | 3.12 |
| EMNIST(Balanced) | **1.60** | 1.82 | 2.78 | 2.15 | 1.96 | 2.23 | 2.88 |
| EMNIST(ByClass) | **1.54** | 1.85 | 2.72 | 1.98 | 1.87 | 2.23 | 3.14 |
| FashionMNIST | **3.27** | 3.34 | 4.29 | 3.47 | 3.28 | 3.66 | 3.72 |

| Dataset | PC | Bipartite flow | AF/SCF | IAF/SCF |
|---|---|---|---|---|
| Penn Treebank | **1.23** | 1.38 | 1.46 | 1.63 |

Meihua Dang, Anji Liu, Guy Van den Broeck, Sparse Probabilistic Circuits via Pruning and Growing, Sparsity in Neural Networks, 2022

Fully factorized

NaiveBayes  AndOrGraphs  PDGs

Trees  PSDDs  CNets  LTMs  SPNs  NADEs

Thin Junction Trees  ACs  MADEs  MAFs  VAEs

DPPs  FVSBNs  TACs  IAFs  NAFs  RAEs

Mixtures  BNs  NICE  FGs  GANs

RealNVP  MNs

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

# Queries as pipelines: KLD

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



Antonio Vergari, YooJung Choi, Anji Liu, Stefano Teso and Guy Van den Broeck. A Compositional Atlas of Tractable Circuit Operations for Probabilistic Inference, *NeurIPS*, 2021.
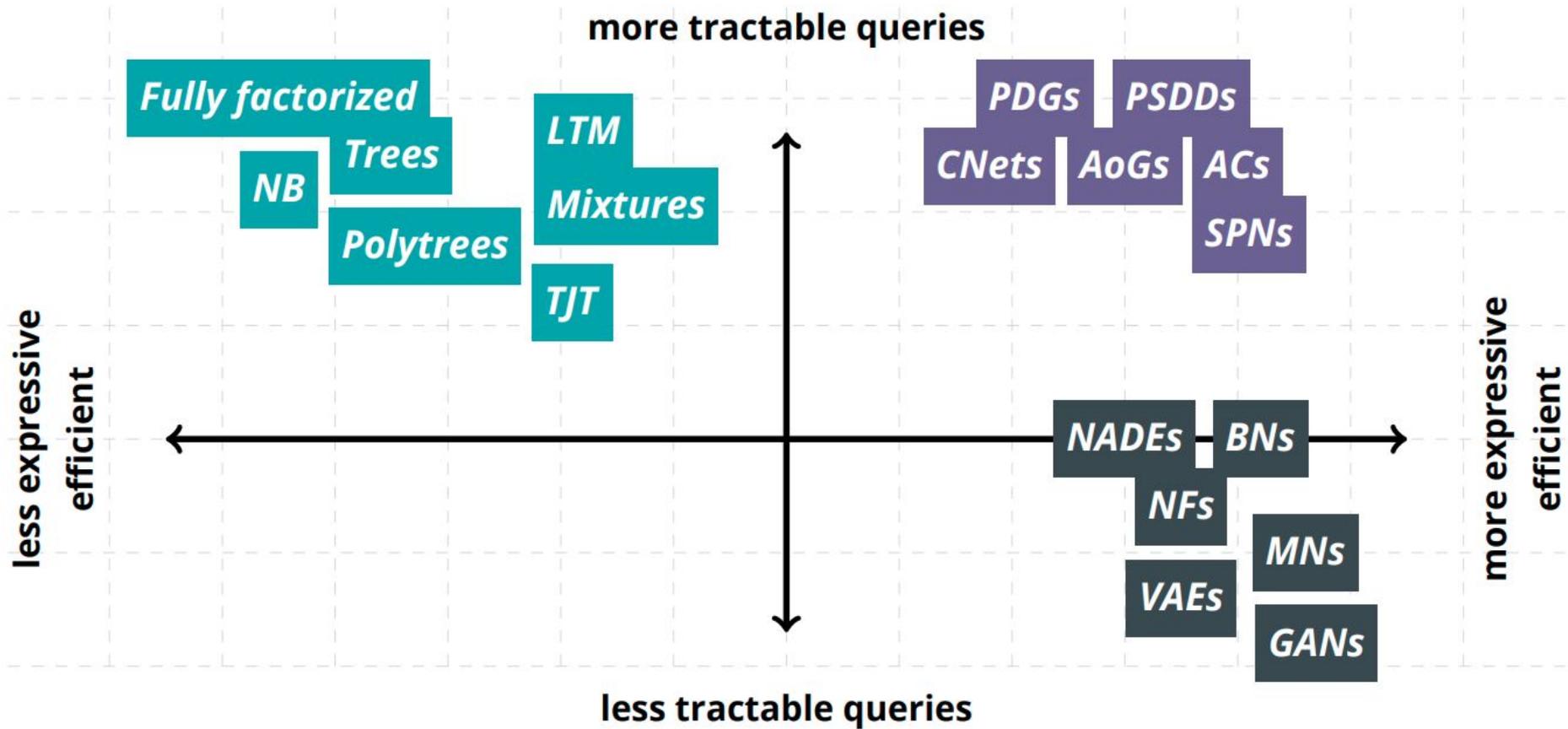
# Inference by tractable operations

***systematically derive*** tractable inference algorithm of complex queries

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

Antonio Vergari, YooJung Choi, Anji Liu, Stefano Teso and Guy Van den Broeck. A Compositional Atlas of Tractable Circuit Operations for Probabilistic Inference, *NeurIPS*, 2021.

**tractability is a spectrum**

more tractable queries

Fully factorized

Trees

NB

Polytrees

LTM

Mixtures

TJT

PDGs · PSDDs

CNets · AoGs · ACs

SPNs

less expressive
efficient

more expressive
efficient

NADEs · BNs

NFs

VAEs

MNs

GANs

less tractable queries

# *Learn more about probabilistic circuits?*

## Tutorial (3h)

## Overview Paper (80p)

# Thanks

*This was the work of many wonderful students/postdocs/collaborators!*

References: http://starai.cs.ucla.edu/publications/