

**UCLA**

**Computer  
Science**



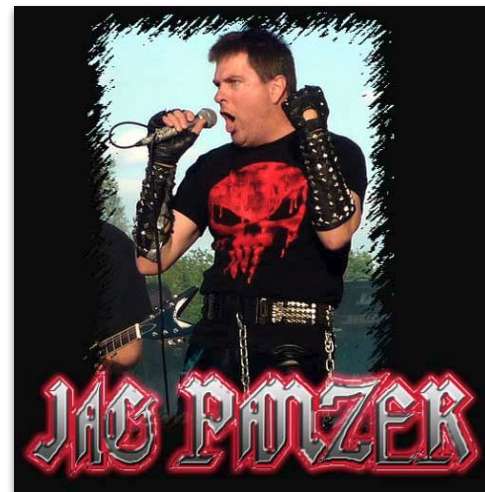
# Computational Abstractions of Probability Distributions

Guy Van den Broeck

PGM - Sep 24, 2020

# Manfred Jaeger Tribute Band

1997-2004-2005



## Relational Bayesian Networks

**Manfred Jaeger\***  
Computer Science Department,  
Stanford University, Stanford CA 94305  
jaeger@robotics.stanford.edu

### Abstract

A new method is developed to represent probabilistic relations on multiple random events. Where previously knowledge bases containing probabilistic rules were used for this purpose,

In case we ind  
served events, e  
 $burglary(\omega_2) =$   
ing a query  $X$   
evidence about  
 $P(X(\omega) = x$

In Int. J. of Uncertainty, Fuzziness

Probabilistic Decision Graphs – Combining  
Verification and AI Techniques for Probabilistic  
Inference

MANFRED JAEGER  
Aalborg Universitet, Institut for Datalogi  
Fredrik Bajers Vej 7, 9220 Aalborg Ø, Denmark  
jaeger@cs.auc.dk



Available online at [www.sciencedirect.com](http://www.sciencedirect.com)

SCIENCE @ DIRECT®

International Journal of Approximate Reasoning  
42 (2006) 4–20

INTERNATIONAL JOURNAL OF  
APPROXIMATE  
REASONING

[www.elsevier.com/locate/ijar](http://www.elsevier.com/locate/ijar)

Compiling relational Bayesian networks  
for exact inference

Mark Chavira <sup>a,\*</sup>, Adnan Darwiche <sup>a</sup>, Manfred Jaeger <sup>b</sup>

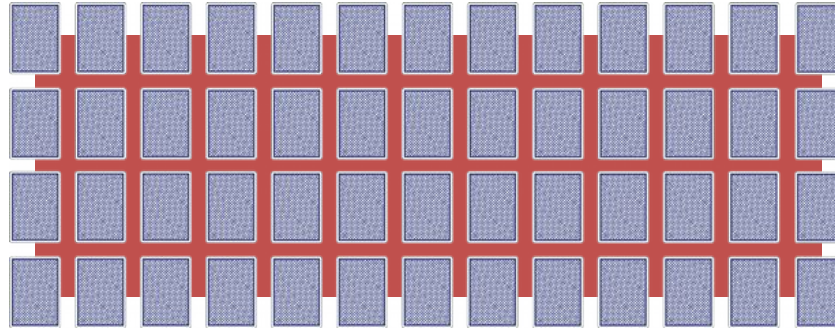
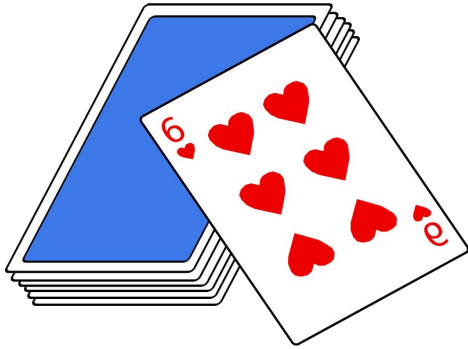
<sup>a</sup> Computer Science Department, UCLA, Los Angeles, CA 90095, United States

<sup>b</sup> Institut for Datalogi, Aalborg Universitet, Fredrik Bajers Vej 7 E, DK-9220 Aalborg Ø, Denmark

Available online 15 November 2005

# Let me be provocative

*Graphical models of variable-level (in)dependence are a broken abstraction.*

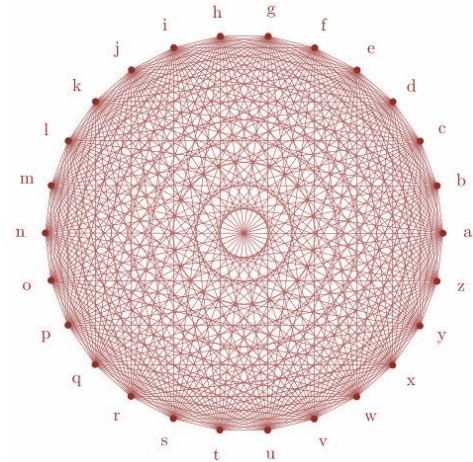


# Let me be provocative

*Graphical models of variable-level (in)dependence are a broken abstraction.*



3.14  $\text{Smokes}(x) \wedge$   
 $\text{Friends}(x,y)$   
 $\Rightarrow \text{Smokes}(y)$



# Let me be provocative

*Graphical models of variable-level (in)dependence are a broken abstraction.*



## Bean Machine

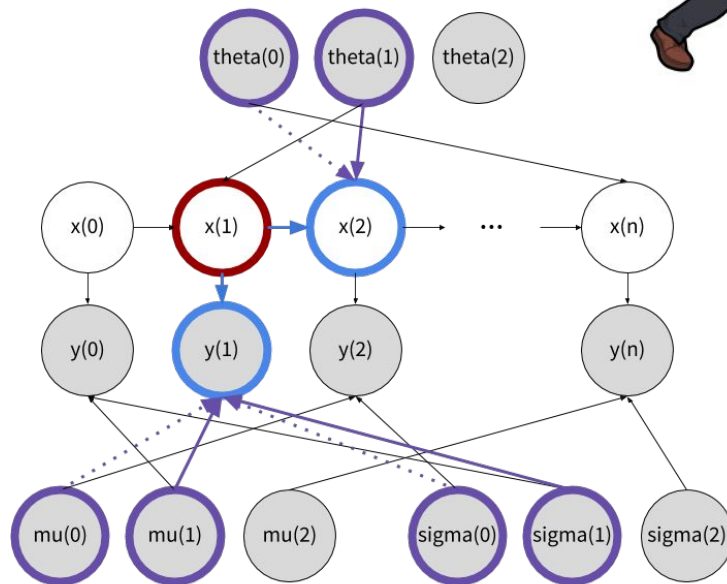
$$\mu_k \sim \text{Normal}(\alpha, \beta)$$

$$\sigma_k \sim \text{Gamma}(\nu, \rho)$$

$$\theta_k \sim \text{Dirichlet}(\kappa)$$

$$x_i \sim \begin{cases} \text{Categorical}(\text{init}) & \text{if } i = 0 \\ \text{Categorical}(\theta_{x_{i-1}}) & \text{if } i > 0 \end{cases}$$

$$y_i \sim \text{Normal}(\mu_{x_i}, \sigma_{x_i})$$



# Let me be **even more** provocative

*Graphical models of variable-level (in)dependence are a broken abstraction.*



We may have gotten stuck in a local optimum?

- Exact probabilistic inference still independence-based
  - Huge effort to extract more local structure from individual tables
- *What do you mean, compute probabilities exactly?*
  - Statistician: inference = Hamiltonian Monte Carlo
  - Machine learner: inference = variational
- Variable-level causality

# Let me be provocative

*Graphical models of variable-level (in)dependence are a broken abstraction.*



The choice of representing a distribution primarily by its variable-level (in)dependencies is a little arbitrary...

*What if we made some different choices?*

# Computational Abstractions

*Let us think of distributions as objects that are computed.*

Abstraction = Structure of Computation

‘closer to the metal’

Two examples:

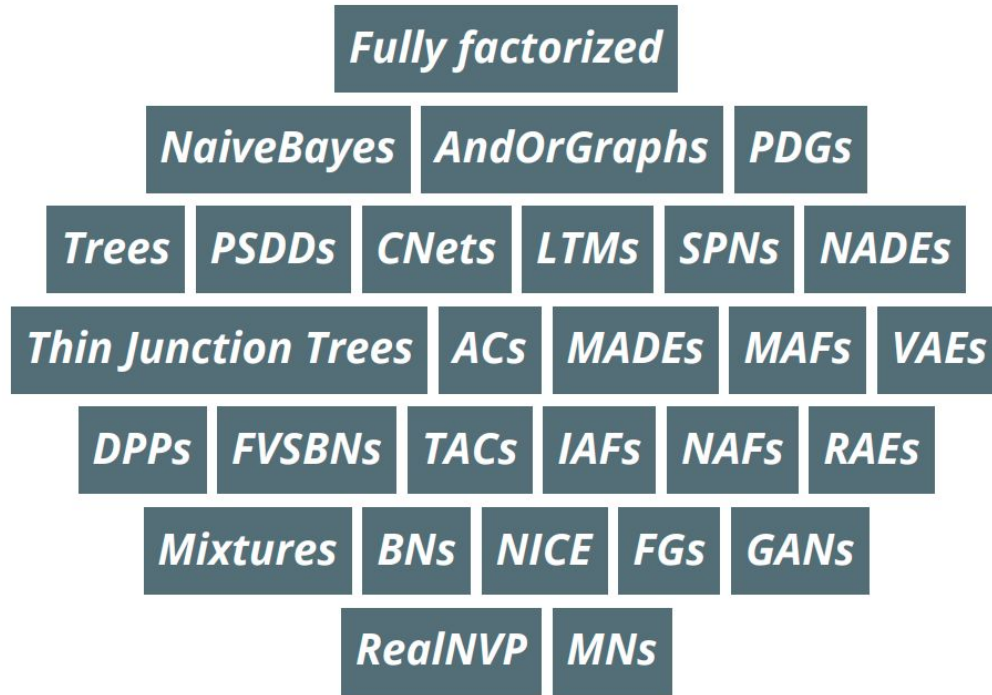
- Probabilistic Circuits
- Probabilistic Programs



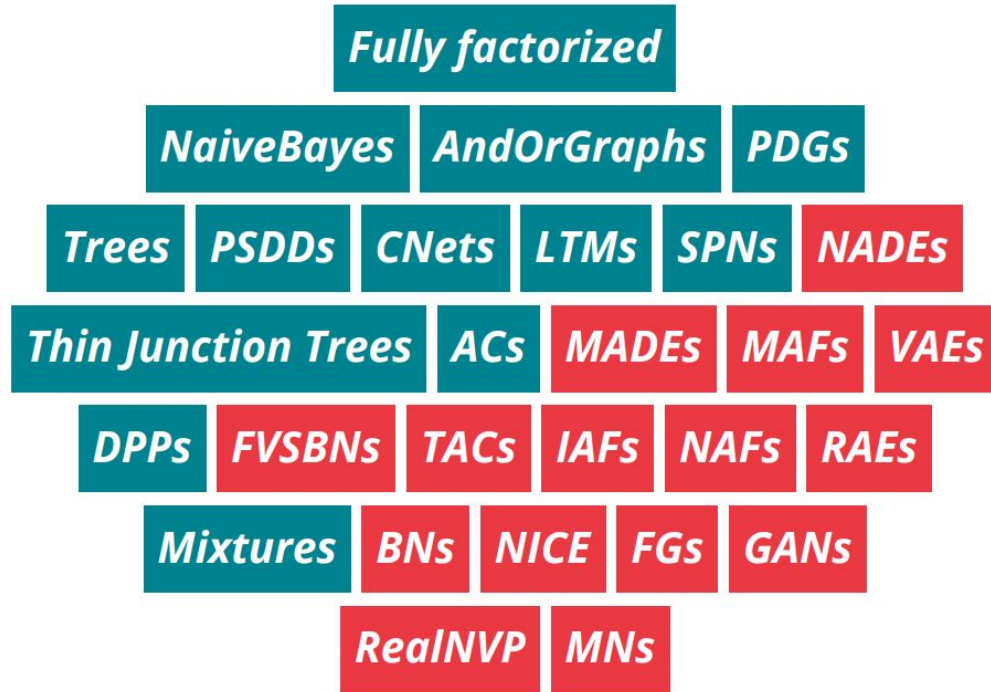


# Probabilistic Circuits

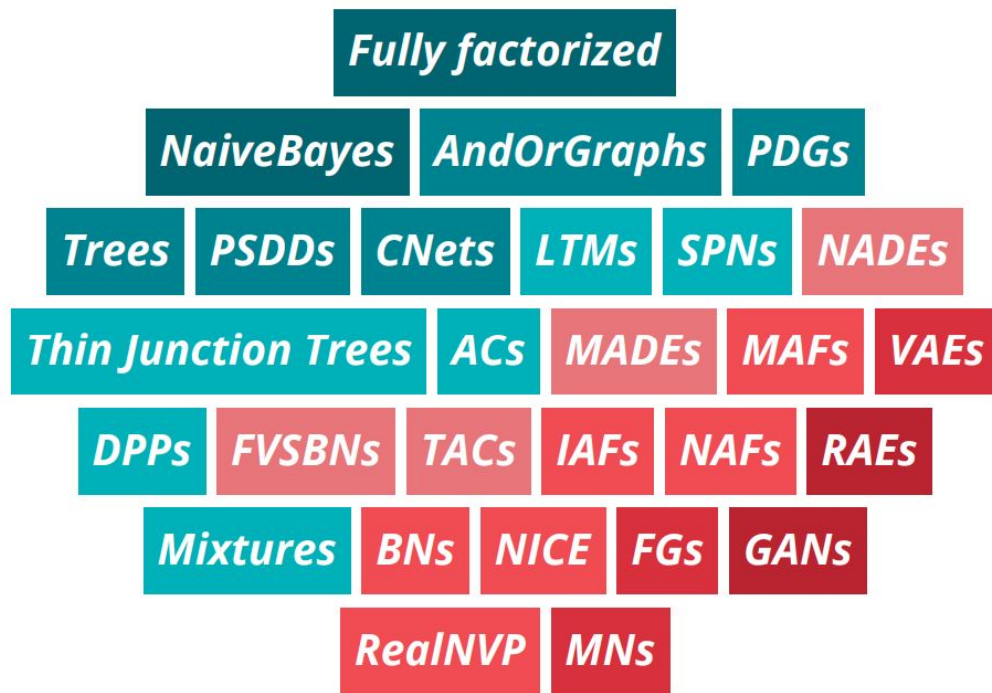




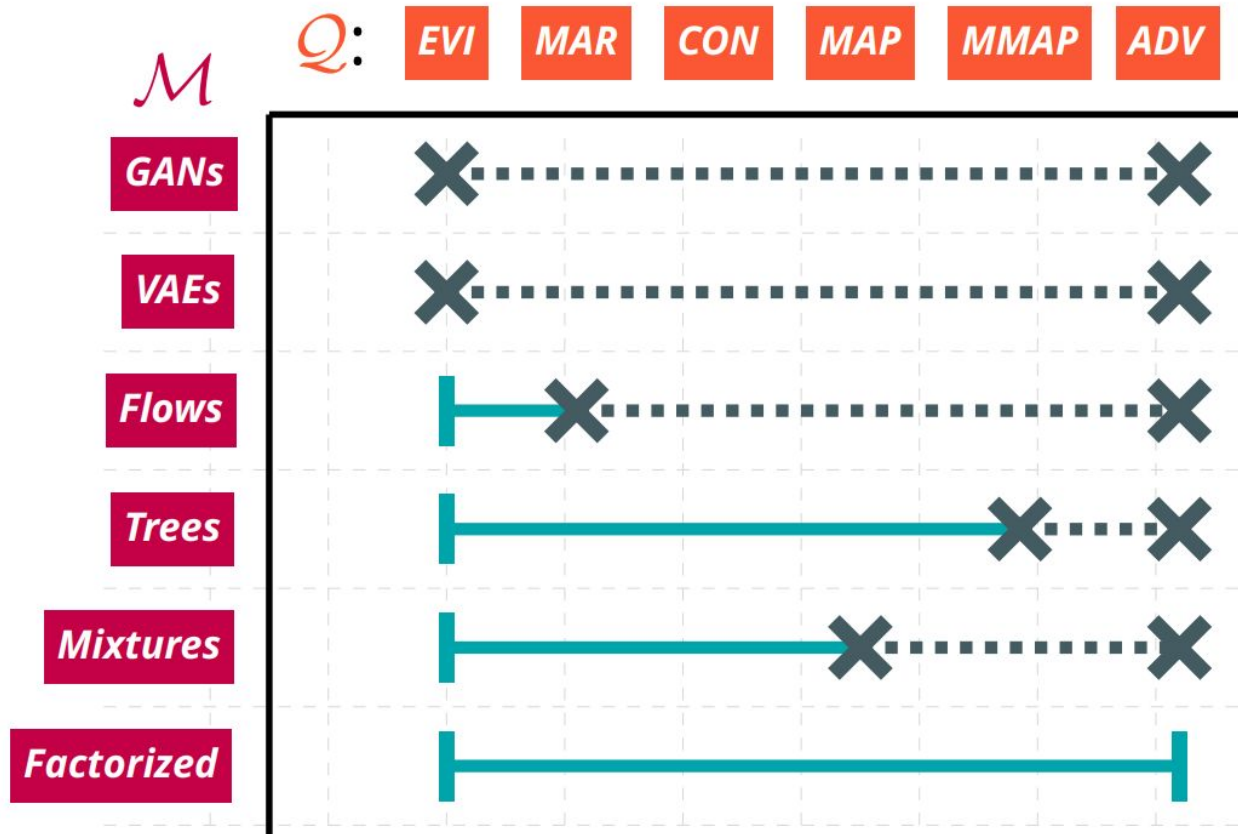
# ***The Alphabet Soup of probabilistic models***

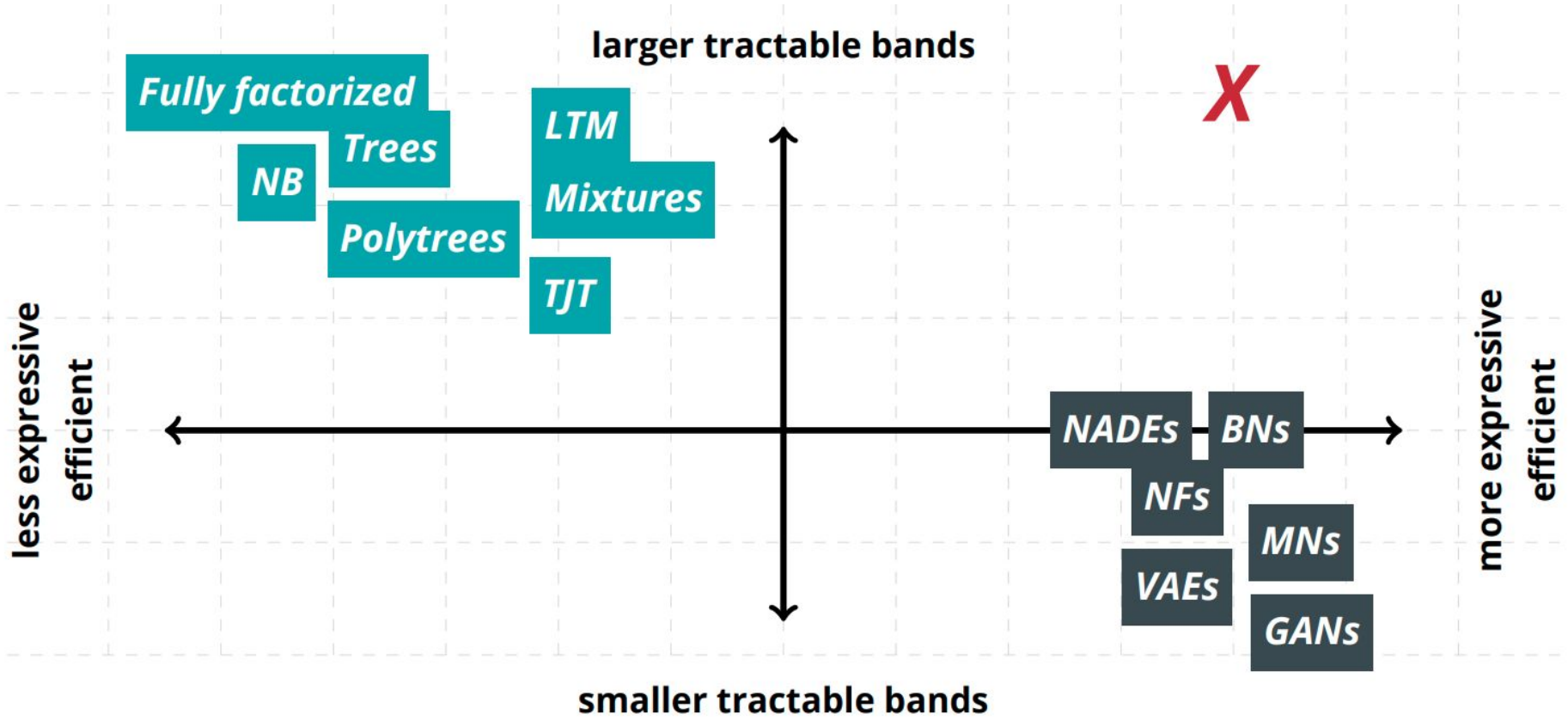


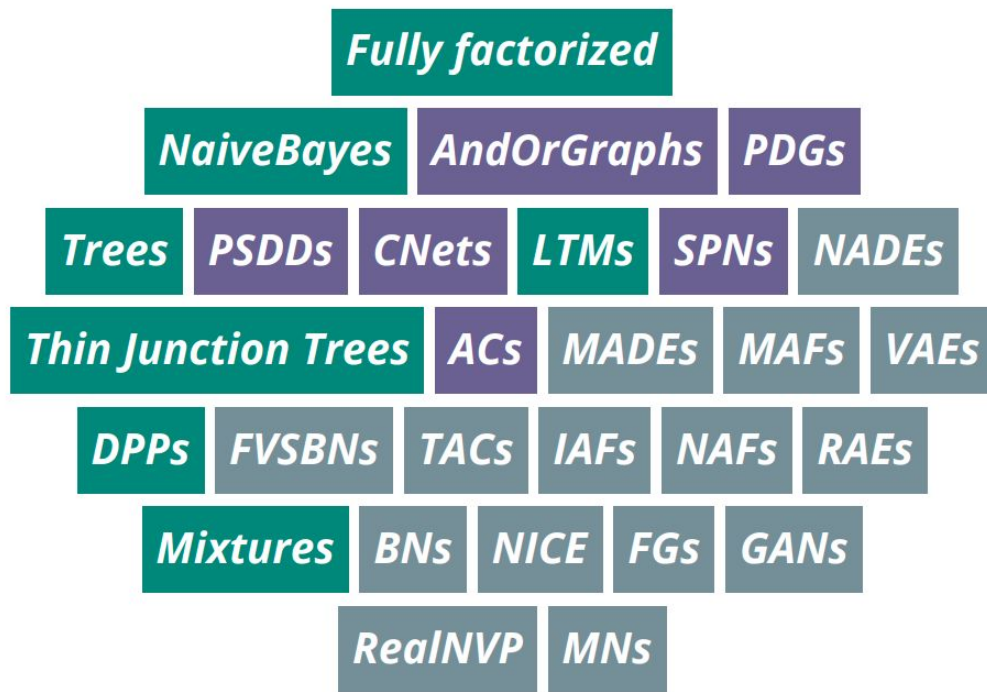
***Intractable*** and ***tractable*** models



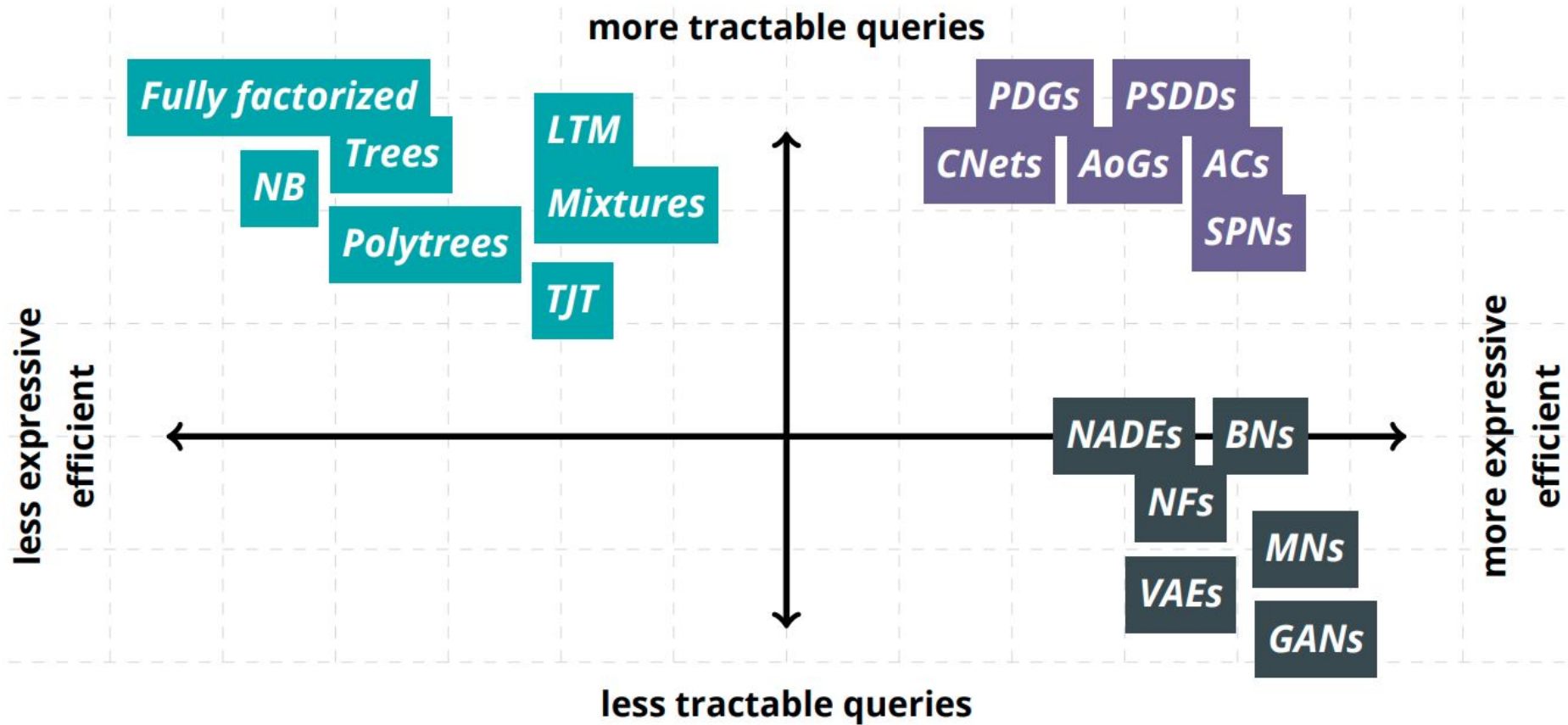
**tractability is a spectrum**





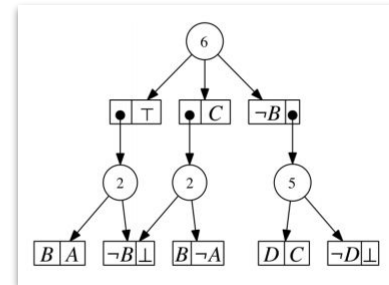
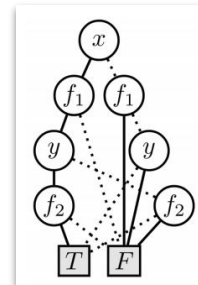
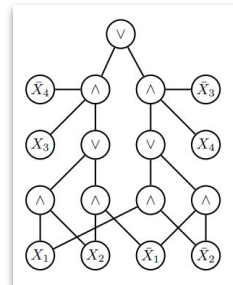
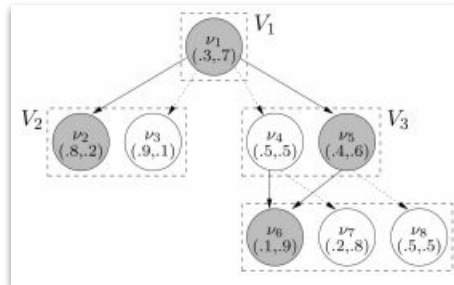
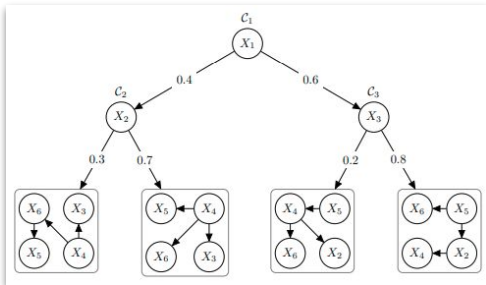
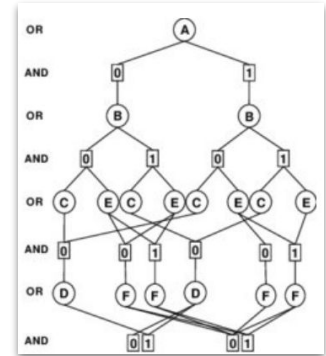
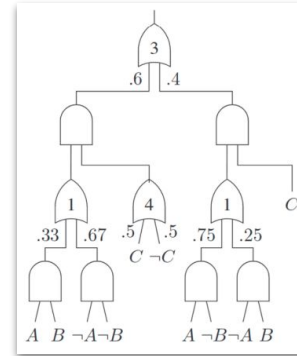
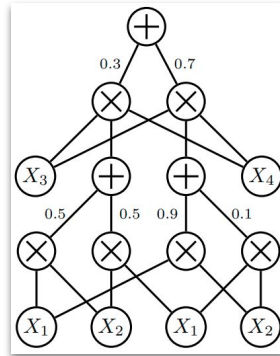
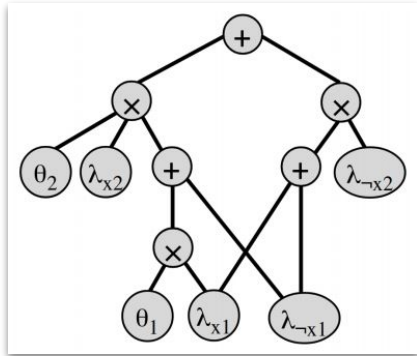
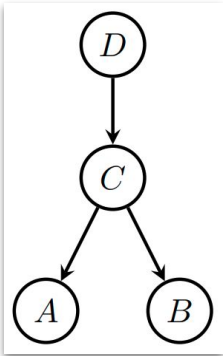


***Expressive* models without *compromises***

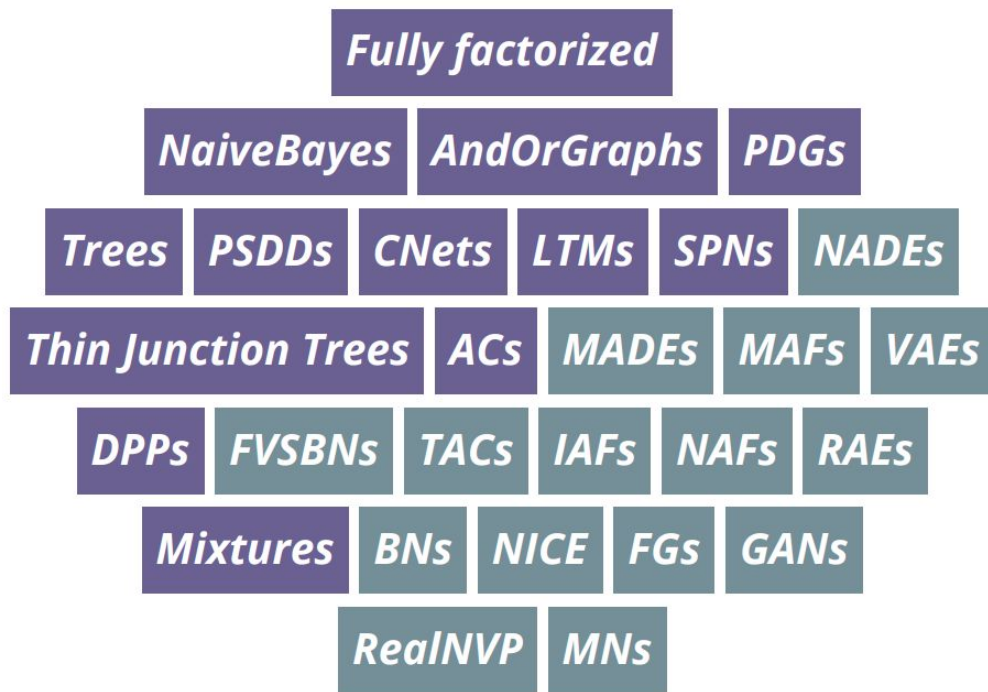




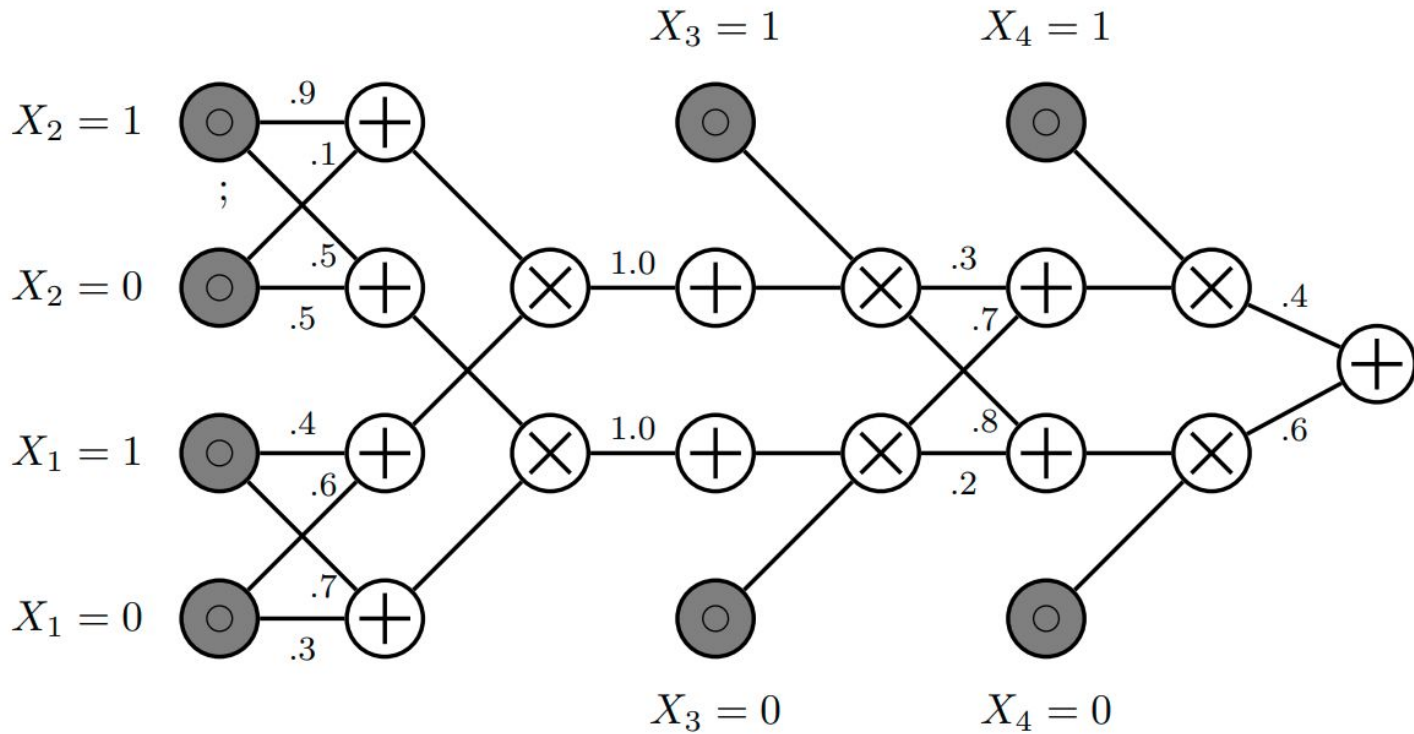
# Tractable Probabilistic Models



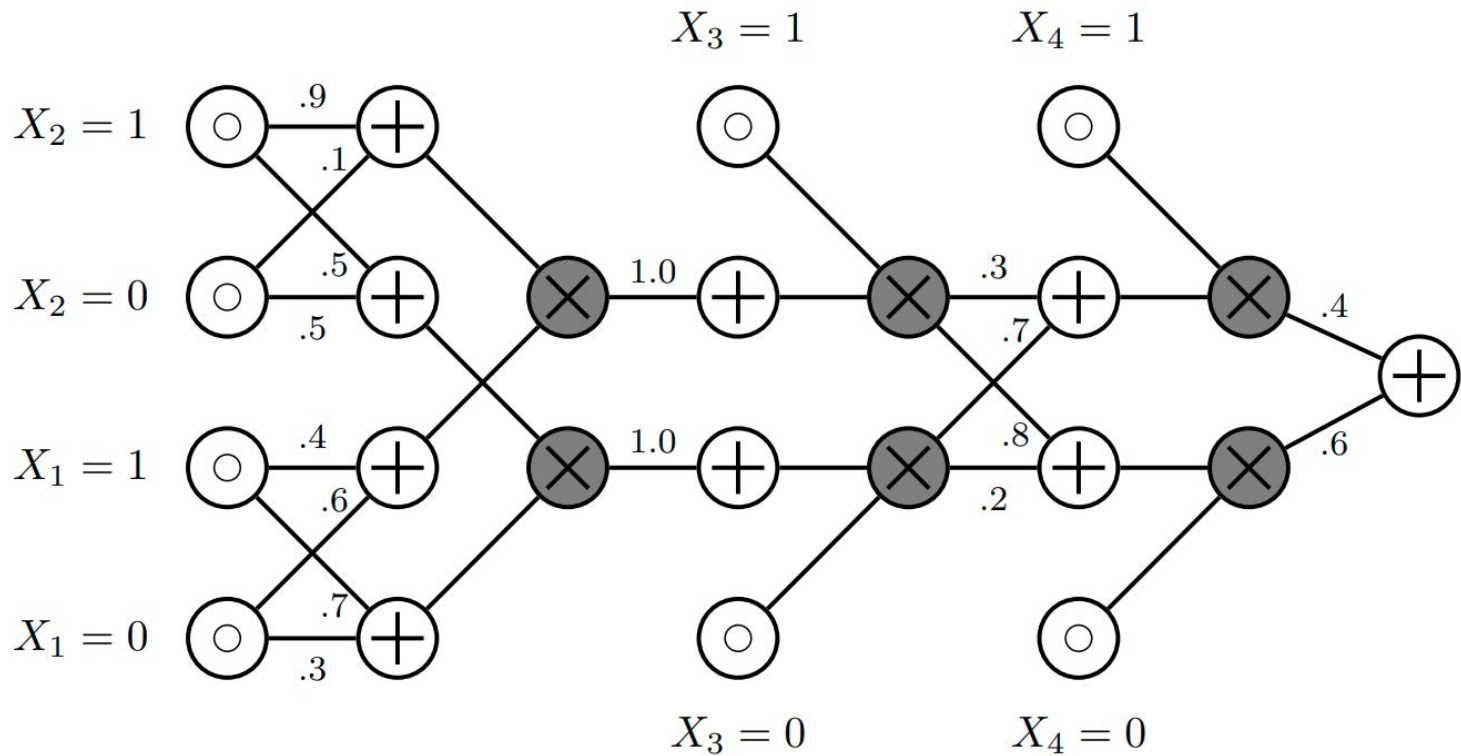
*"Every keynote needs a joke and a literature overview slide, not necessarily distinct"*  
 - after Ron Graham



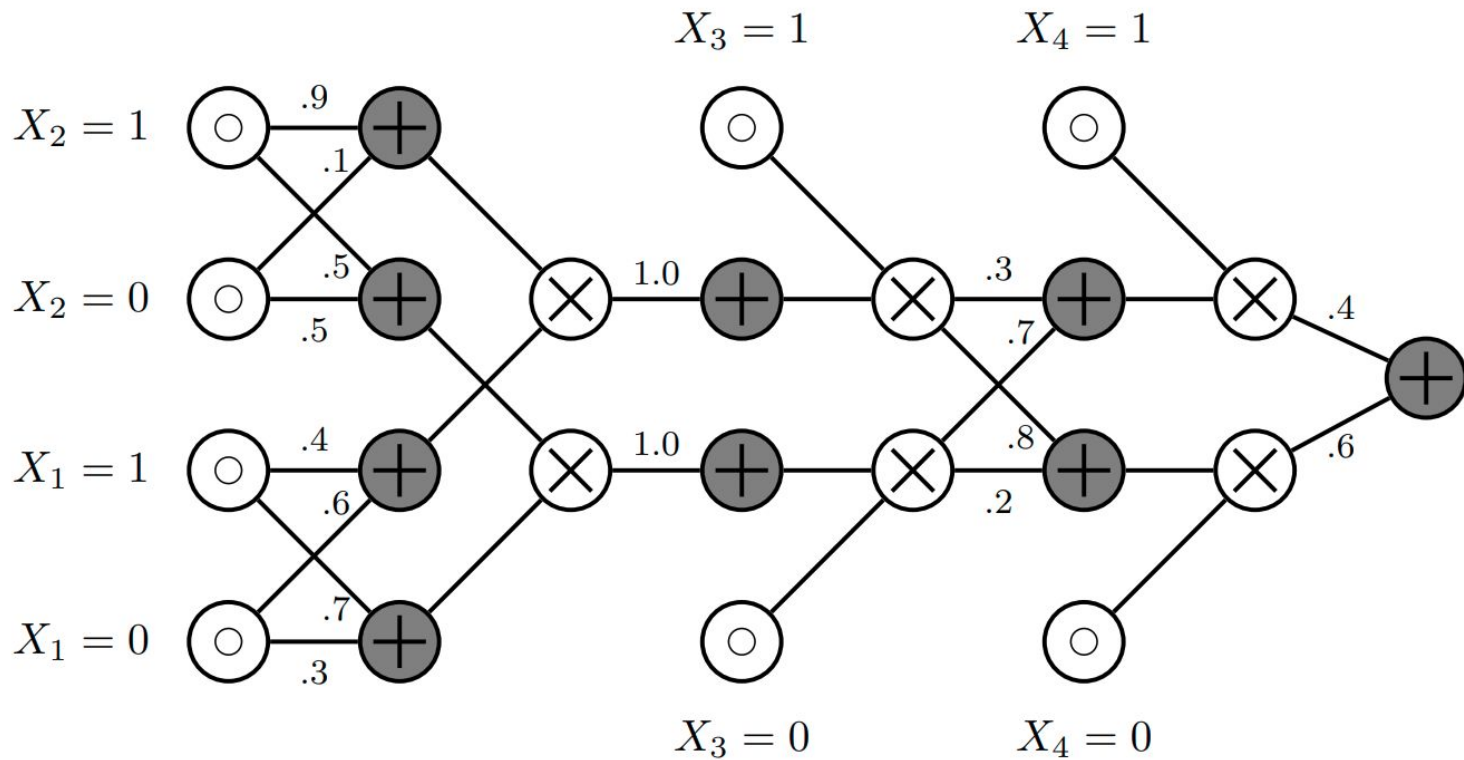
***a unifying framework* for tractable models**



Input nodes are tractable (simple) distributions, e.g., indicator functions  $p_n(X=1) = [X=1]$



Product nodes are factorizations  $\prod_{c \in \text{in}(n)} p_c(\mathbf{x})$



Sum nodes are mixture models  $\sum_{c \in \text{in}(n)} \theta_{n,c} p_c(\mathbf{x})$

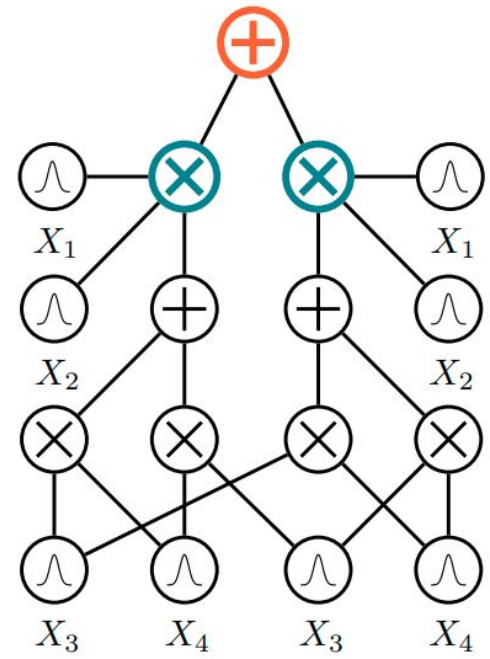
# Smoothness + decomposability = tractable MAR

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

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

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

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

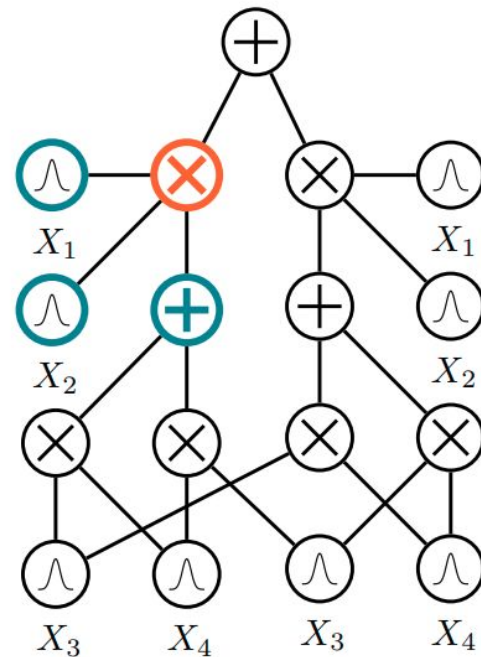


# Smoothness + decomposability = tractable MAR

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

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

$\Rightarrow$  integrals decompose into easier ones





# Smoothness + decomposability = tractable MAR

Forward pass evaluation for MAR

$\Rightarrow$  linear in circuit size!

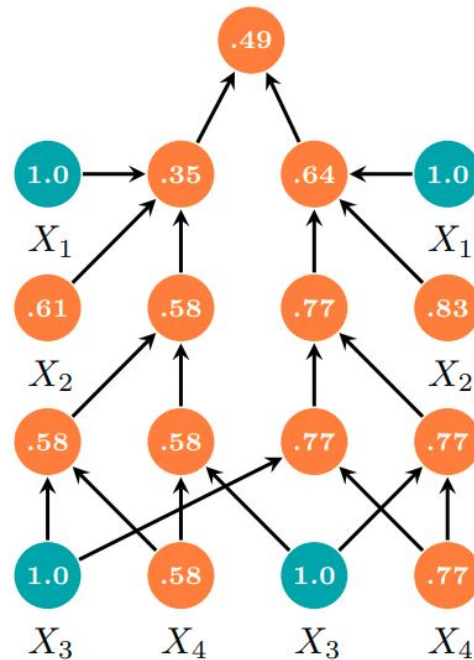
E.g. to compute  $p(x_2, x_4)$ :

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

$\Rightarrow$  for normalized leaf distributions: 1.0

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

feedforward evaluation (bottom-up)



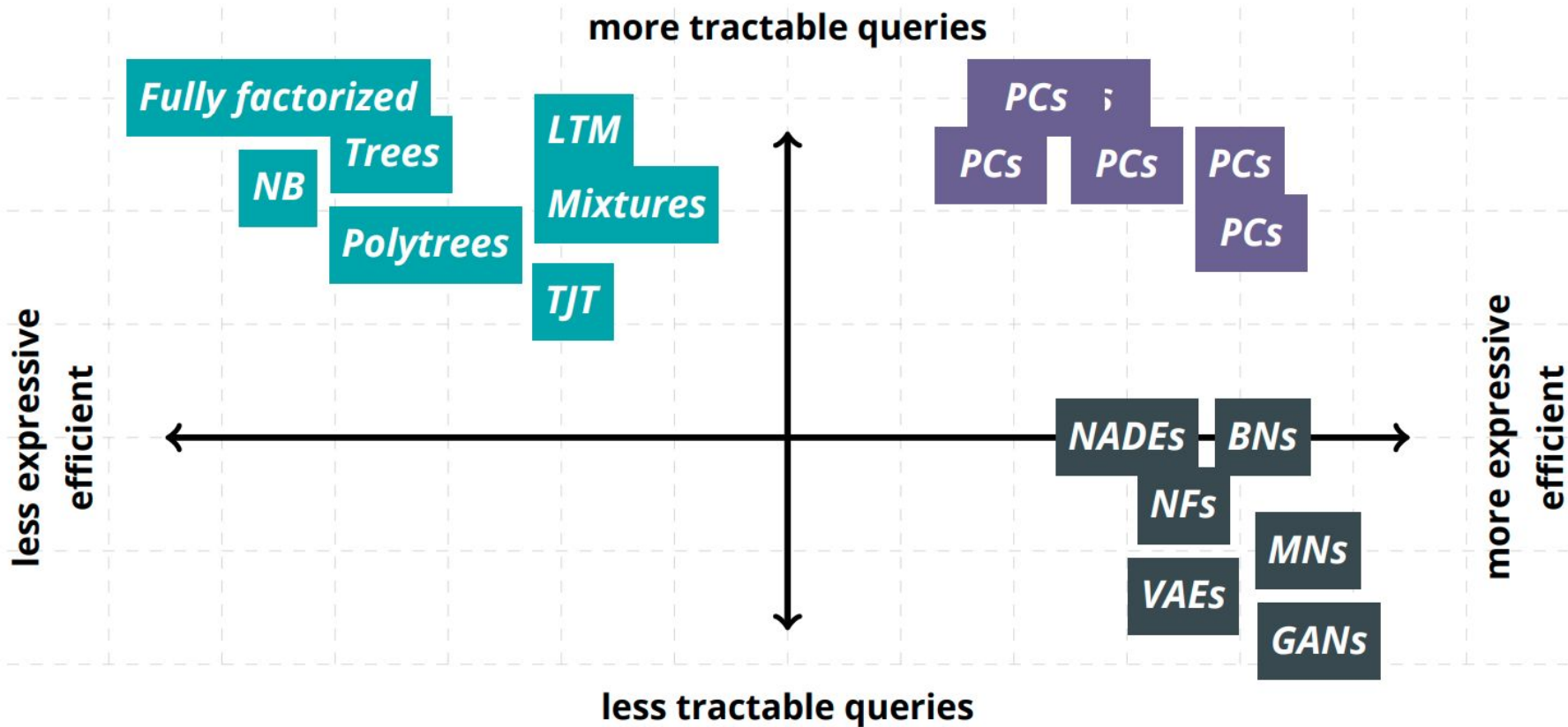


$\mathcal{P}$		$\mathcal{Q}$ :								
		MAR	CON	MOM	MAP	MMAP	ENT	DIV	EXP	
		marginal queries	conditional queries	moments (mean,...)	maximum a posteriori	marginal MAP	entropy	divergences (KLD,...)	expected predictions	
smoothness	<b>SMO</b>	✓	✓	✓	✗	✗	✓	✓	✓	
decomposability	<b>DEC</b>	✓	✓	✓	✗	✗	✗	✗	✗	
consistency	<b>CON</b>	✗	✗	✗	✓	✓	✗	✗	✗	
determinism	<b>DET</b>	✗	✗	✗	✓	✗	✗	✗	✗	
marginal determinism	<b>MAR-DET</b>	✗	✗	✗	✗	✓	✓	✓	✗	
structured decomposability	<b>STR-DEC</b>	✗	✗	✗	✗	✗	✓	✗	✗	
paired str. decomposability	<b>P-STR-DEC</b>	✗	✗	✗	✗	✗	✗	✓	✓	

---

	<i>smooth</i>	<i>dec.</i>	<i>det.</i>	<i>str.dec.</i>
Arithmetic Circuits (ACs) <i>[Darwiche 2003]</i>	✓	✓	✓	✗
Sum-Product Networks (SPNs) <i>[Poon et al. 2011]</i>	✓	✓	✗	✗
Cutset Networks (CNets) <i>[Rahman et al. 2014]</i>	✓	✓	✓	✗
Probabilistic Decision Graphs <i>[Jaeger 2004]</i>	✓	✓	✓	✓
PSDDs <i>[Kisa et al. 2014a]</i>	✓	✓	✓	✓
AndOrGraphs <i>[Dechter et al. 2007]</i>	✓	✓	✓	✓

---



# How expressive are probabilistic circuits?

*density estimation benchmarks*

dataset	best circuit	BN	MADE	VAE	dataset	best circuit	BN	MADE	VAE
<i>nlcs</i>	<b>-5.99</b>	-6.02	-6.04	<b>-5.99</b>	<i>dna</i>	<b>-79.88</b>	-80.65	-82.77	-94.56
<i>msnbc</i>	<b>-6.04</b>	<b>-6.04</b>	-6.06	-6.09	<i>kosarek</i>	<b>-10.52</b>	-10.83	-	-10.64
<i>kdd</i>	-2.12	-2.19	<b>-2.07</b>	-2.12	<i>msweb</i>	-9.62	-9.70	<b>-9.59</b>	-9.73
<i>plants</i>	<b>-11.84</b>	-12.65	-12.32	-12.34	<i>book</i>	-33.82	-36.41	-33.95	<b>-33.19</b>
<i>audio</i>	-39.39	-40.50	-38.95	<b>-38.67</b>	<i>movie</i>	-50.34	-54.37	-48.7	<b>-47.43</b>
<i>jester</i>	-51.29	<b>-51.07</b>	-52.23	-51.54	<i>webkb</i>	-149.20	-157.43	-149.59	<b>-146.9</b>
<i>netflix</i>	-55.71	-57.02	-55.16	<b>-54.73</b>	<i>cr52</i>	-81.87	-87.56	-82.80	<b>-81.33</b>
<i>accidents</i>	-26.89	<b>-26.32</b>	-26.42	-29.11	<i>c20ng</i>	-151.02	-158.95	-153.18	<b>-146.9</b>
<i>retail</i>	<b>-10.72</b>	-10.87	-10.81	-10.83	<i>bbc</i>	<b>-229.21</b>	-257.86	-242.40	-240.94
<i>pumbs*</i>	-22.15	<b>-21.72</b>	-22.3	-25.16	<i>ad</i>	-14.00	-18.35	<b>-13.65</b>	-18.81

# Want to learn more?

## Tutorial (3h)

**Probabilistic Circuits**

**Inference  
Representations  
Learning  
Theory**

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

**Robert Peharz**  
TU Eindhoven

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

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

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

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

## Overview Paper (80p)

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

YooJung Choi

Antonio Vergari

Guy Van den Broeck

*Computer Science Department*

*University of California*

*Los Angeles, CA, USA*

#### Contents

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

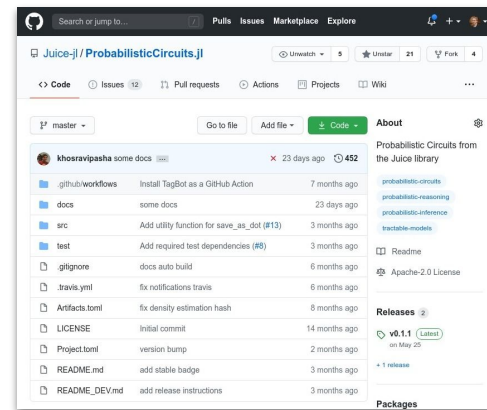
<http://starai.cs.ucla.edu/papers/ProbCirc20.pdf>

# Training PCs in Julia with Juice.jl



Training maximum likelihood parameters of probabilistic circuits

```
julia> using ProbabilisticCircuits;
julia> data, structure = load(...);
julia> num_examples(data)
17412
julia> num_edges(structure)
270448
julia> @btime estimate_parameters(structure , data);
63 ms
```



Custom SIMD and CUDA kernels to parallelize over layers and training examples.

<https://github.com/Juice-jl/>

Probabilistic circuits seem awfully general.

*Are all tractable probabilistic models  
probabilistic circuits?*



# Determinantal Point Processes (DPPs)

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

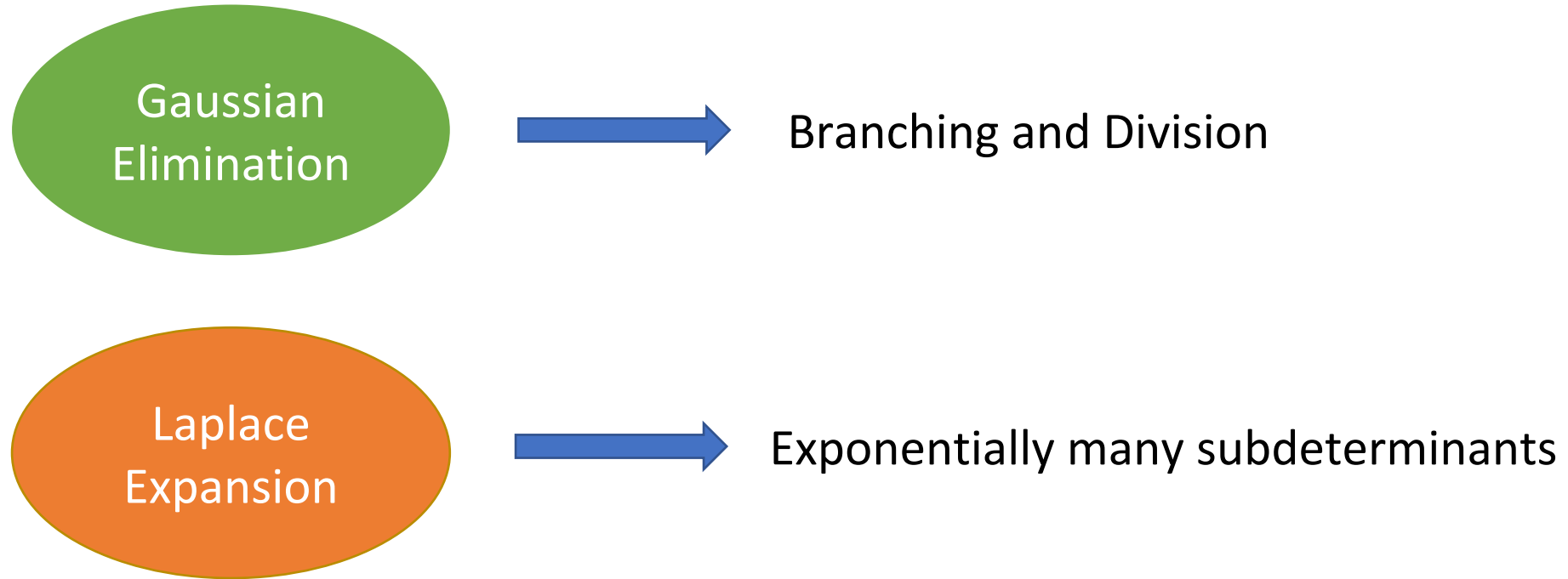
$$L = \begin{bmatrix} \mathbf{1} & 0.9 & \mathbf{0.8} & 0 \\ 0.9 & 0.97 & 0.96 & 0 \\ \mathbf{0.8} & 0.96 & \mathbf{1} & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\Pr_L(\mathbf{X}_1 = \mathbf{1}, X_2 = 0, \mathbf{X}_3 = \mathbf{1}, X_4 = 0) = \frac{1}{\det(L + I)} \det(L_{\{1,2\}})$$

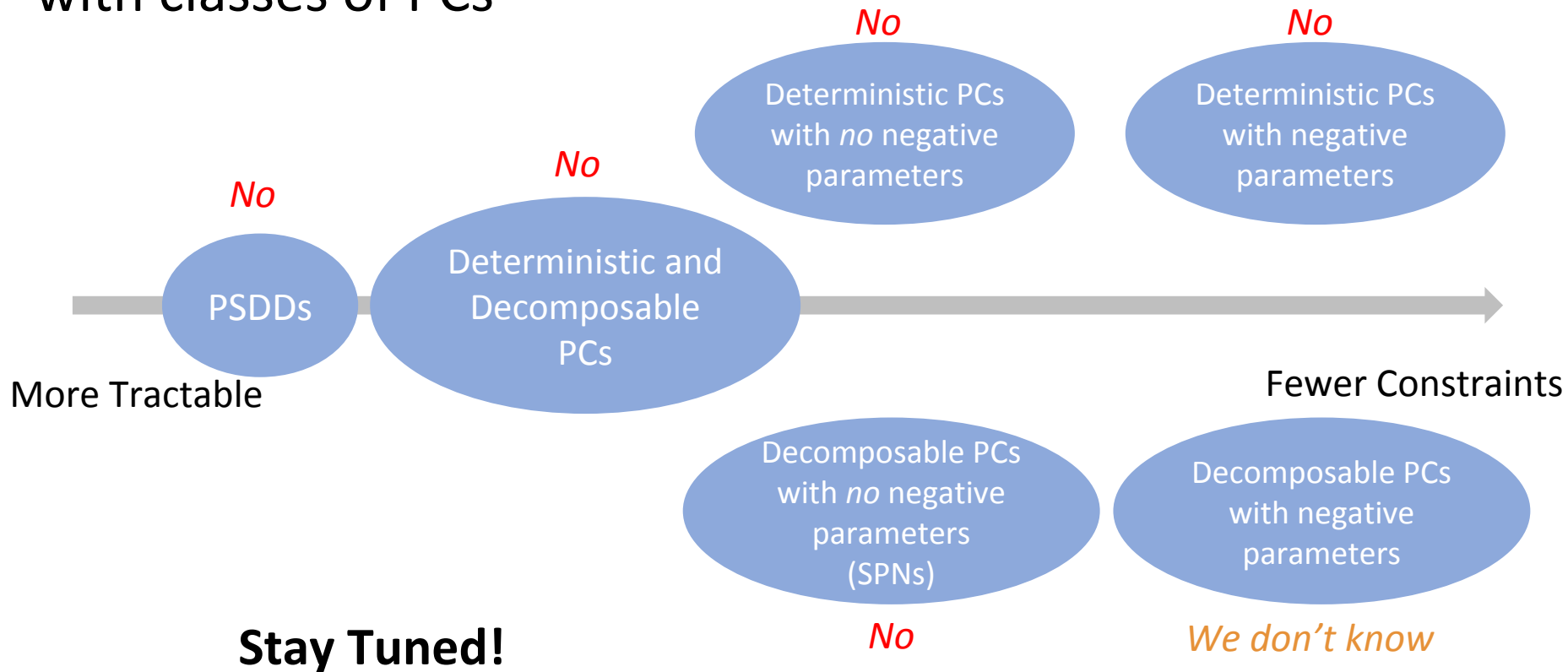
Computing marginal probabilities is *tractable*.



# Representing the Determinant as a PC is *not easy*



# We cannot tractably represent DPPs with classes of PCs



# The AI Dilemma



**Pure Logic**

**Pure Learning**

# The AI Dilemma



## Pure Logic

## Pure Learning

- Slow thinking: deliberative, cognitive, model-based, extrapolation
- Amazing achievements until this day
- “*Pure logic is brittle*”  
noise, uncertainty, incomplete knowledge, ...



# The AI Dilemma



## Pure Logic

## Pure Learning

- Fast thinking: instinctive, perceptive, model-free, interpolation
- Amazing achievements recently
- “*Pure learning is brittle*”
  - bias, algorithmic fairness, interpretability, explainability, adversarial attacks, unknown unknowns, calibration, verification, missing features, missing labels, data efficiency, shift in distribution, general robustness and safety
  - fails to incorporate a sensible model of the world





Pure Logic   Probabilistic World Models   Pure Learning



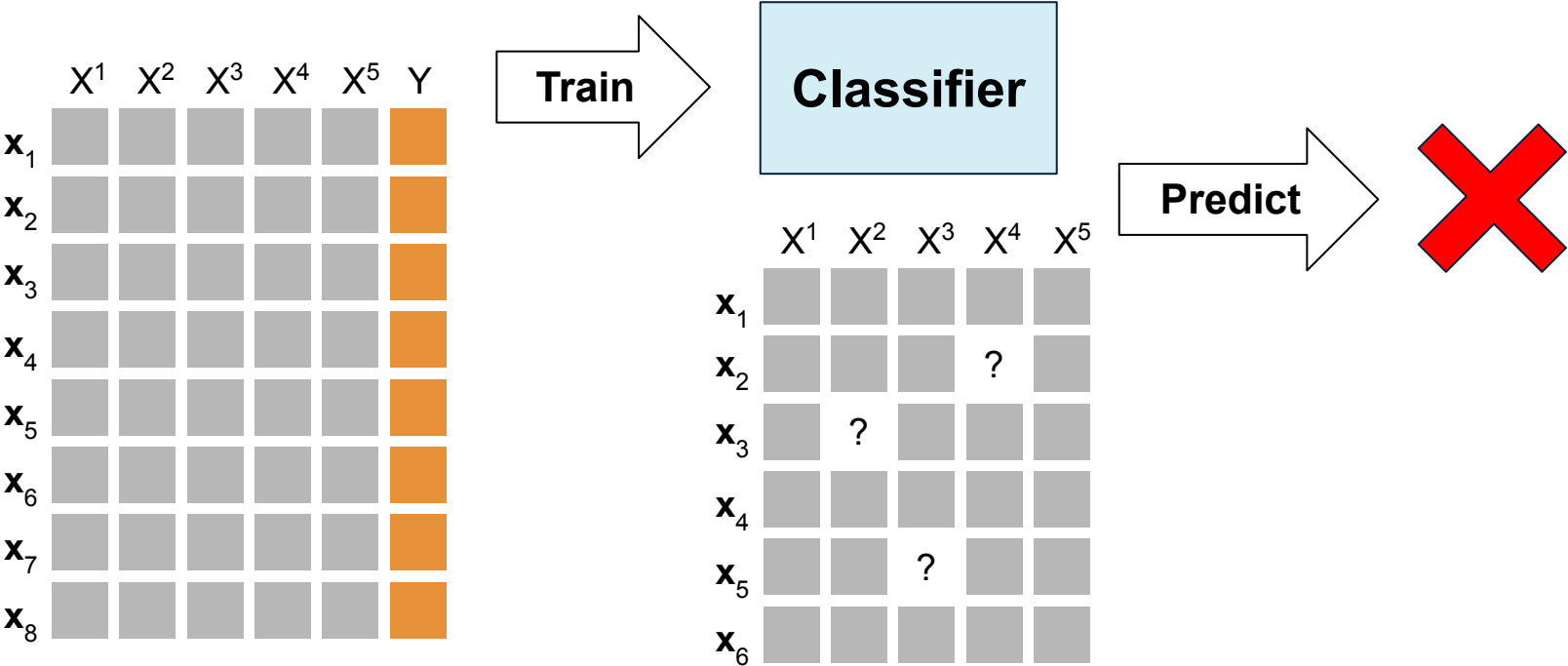
**A New Synthesis of  
Learning and Reasoning**

*“Pure learning is brittle”*

bias, **algorithmic fairness**, interpretability, **explainability**, adversarial attacks, unknown unknowns, calibration, verification, **missing features**, missing labels, data efficiency, shift in distribution, general robustness and safety

We need to incorporate a sensible probabilistic model of the world

# Prediction with Missing Features



Test with missing features

# Expected Predictions

Consider **all possible complete inputs** and **reason** about the expected behavior of the classifier

$$\mathbb{E}_{\mathbf{x}^m \sim p(\mathbf{x}^m | \mathbf{x}^o)} [f(\mathbf{x}^m \mathbf{x}^o)]$$

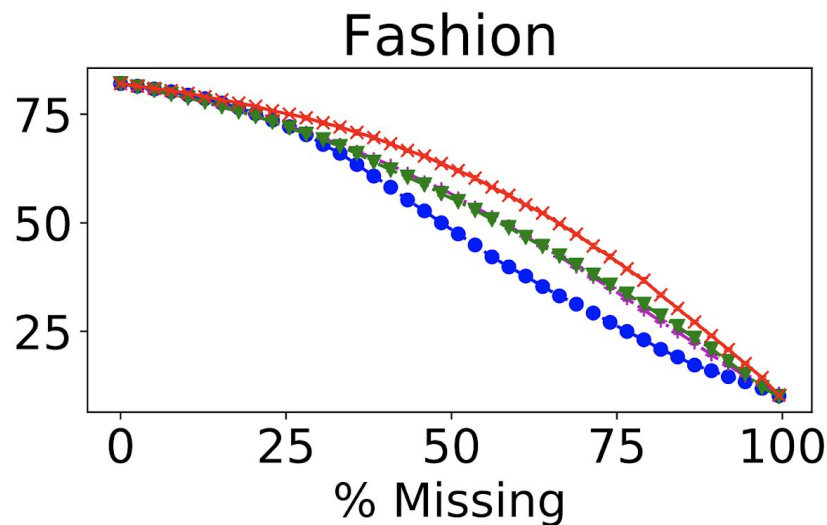
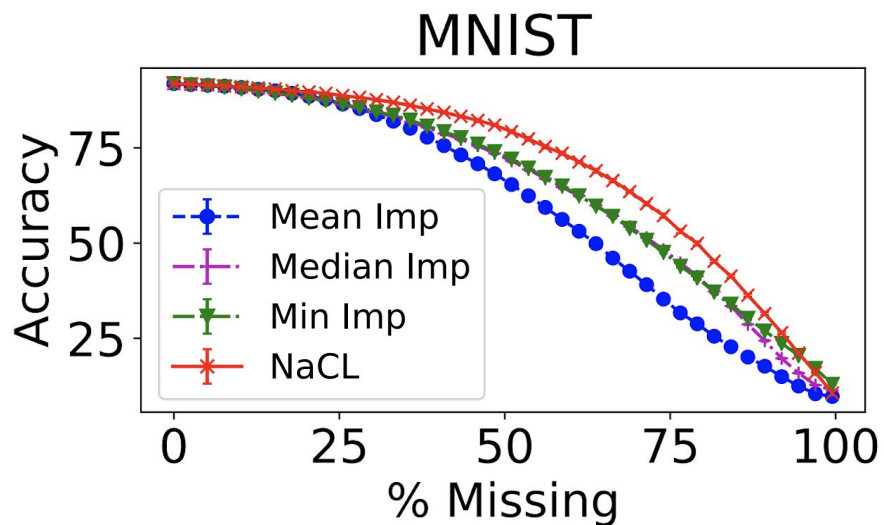
$\mathbf{x}^o$  = observed features  
 $\mathbf{x}^m$  = missing features

Generalizes what we've been doing all along...

$$\begin{aligned} P(C|\mathbf{y}) &= \sum_{\mathbf{m}} P(C, \mathbf{m}|\mathbf{y}) \\ &= \sum_{\mathbf{m}} P(C|\mathbf{m}, \mathbf{y}) P(\mathbf{m}|\mathbf{y}) \\ &= \mathbb{E}_{\mathbf{m} \sim P(\mathbf{M}|\mathbf{y})} P(C|\mathbf{m}, \mathbf{y}) \end{aligned}$$



# Experiments with simple distributions (Naive Bayes) to reason about missing data in logistic regression



*“Conformant learning”*

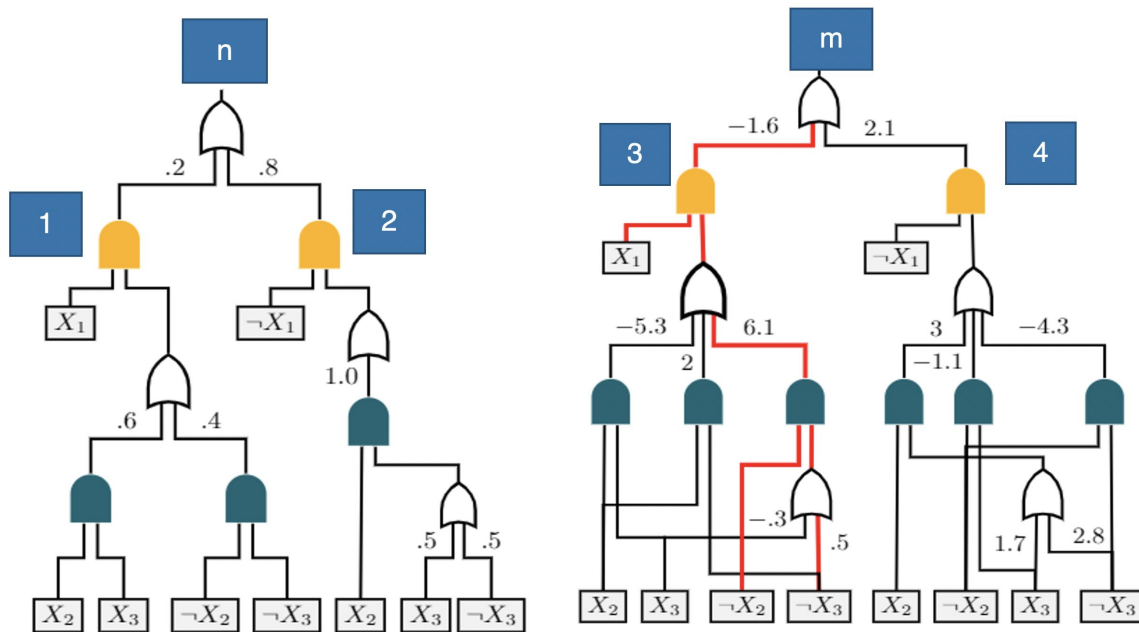
# What about complex classifiers and distributions?

Tractable expected predictions if the classifier is a regression circuit, and the feature distribution is a compatible probabilistic circuits

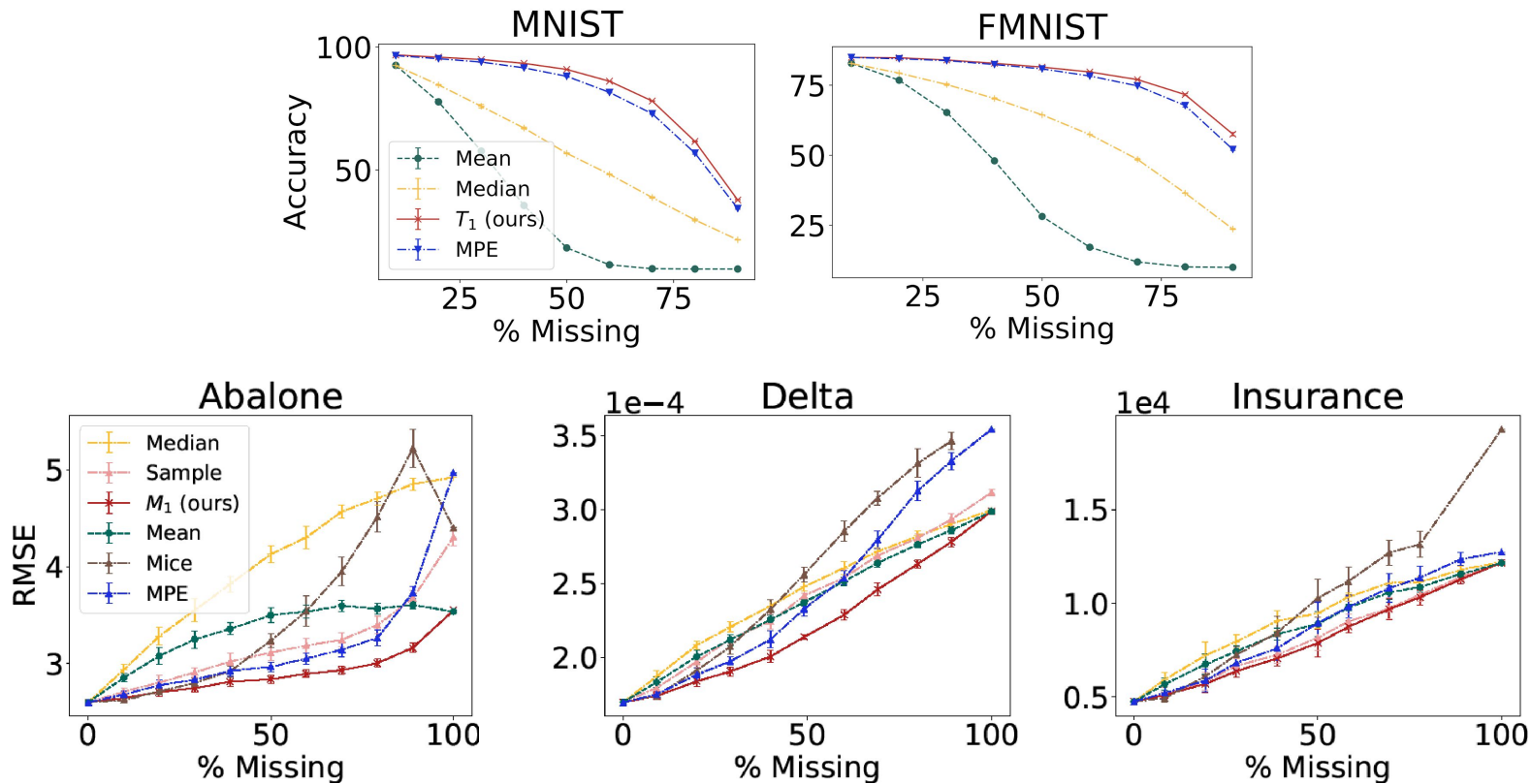


Recursion that “breaks down” the computation.

For + nodes (n,m), look at subproblems (1,3), (1,4), (2,3), (2,4)



# Experiments with Probabilistic Circuits



# ADV inference in Julia with Juice.jl



```
using ProbabilisticCircuits
pc = load_prob_circuit(zoo_psdd_file("insurance.psdd"));
rc = load_logistic_circuit(zoo_lc_file("insurance.circuit"), 1);
```

**q<sub>8</sub>**: *How different is the insurance costs between smokers and non smokers?*

```
groups = make_observations([["!smoker"], ["smoker"]])
exps, _ = Expectation(pc, rc, groups);
println("Smoker      : \$ $(exps[2])");
println("Non-Smoker: \$ $(exps[1])");
println("Difference: \$ $(exps[2] - exps[1])");
Smoker      : $ 31355.32630488978
Non-Smoker: $ 8741.747258310648
Difference: $ 22613.57904657913
```

# What If Training Also Has Missingness

This time we consider decision trees as the classifier

$$\mathcal{L}(\Theta; \mathbf{D}_{\text{train}}) = \frac{1}{|\mathbf{D}_{\text{train}}|} \sum_{\mathbf{x}^o, y \in \mathbf{D}_{\text{train}}} \mathbb{E}_{p_{\Phi}(\mathbf{X}^m | \mathbf{x}^o)} [l(y, f_{\Theta}(\mathbf{x}))]$$

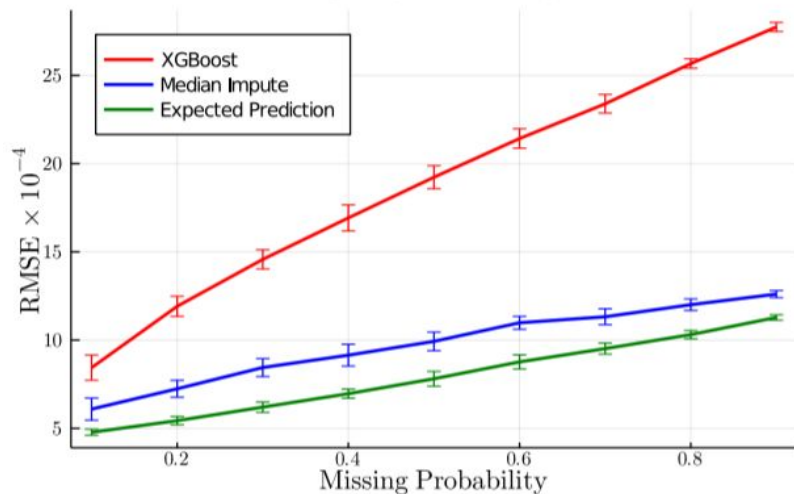
For one decision tree and using MSE loss, can be computed exactly

$$\theta_{\ell}^* = \frac{\sum_{\mathbf{x}^o, y \in \mathbf{D}_{\text{train}}} y \cdot p_{\ell}(\mathbf{x}^o) / p(\mathbf{x}^o)}{\sum_{\mathbf{x}^o, y \in \mathbf{D}_{\text{train}}} p_{\ell}(\mathbf{x}^o) / p(\mathbf{x}^o)}$$

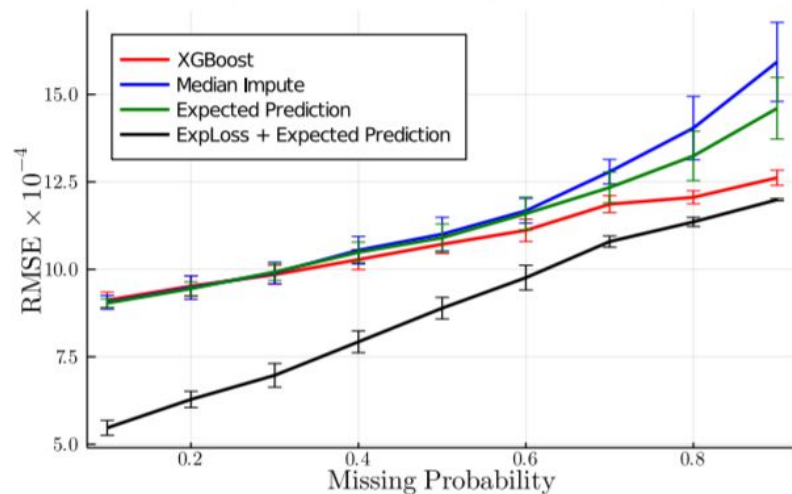
More scenarios such as bagging/boosting in the paper.

# Preliminary Experiments

Missing only at Deployment



Missing at both Learning and Deployment



# ADV inference in Julia with Juice.jl



```
using ProbabilisticCircuits
pc = load_prob_circuit(zoo_psdd_file("insurance.psdd"));
rc = load_logistic_circuit(zoo_lc_file("insurance.circuit"), 1);
```

**Q9:** *Is the predictive model biased by gender?*

```
groups = make_observations([[ "male" ], [ "female" ]])
exps, _ = Expectation(pc, rc, groups);
println("Female   : \$ $(exps[2])");
println("Male     : \$ $(exps[1])");
println("Diff      : \$ $(exps[2] - exps[1])");
Female   : $ 14170.125469335406
Male     : $ 13196.548926381849
Diff     : $ 973.5765429535568
```

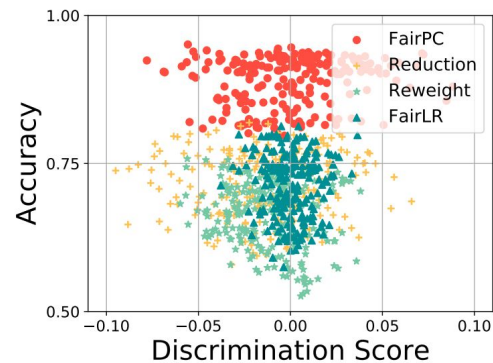
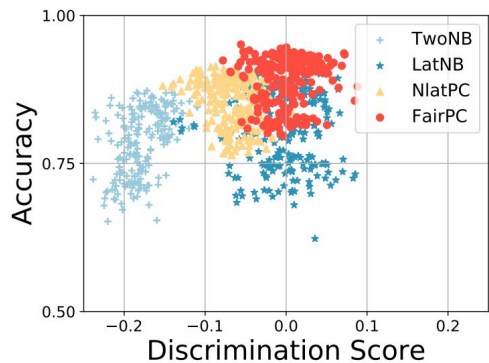
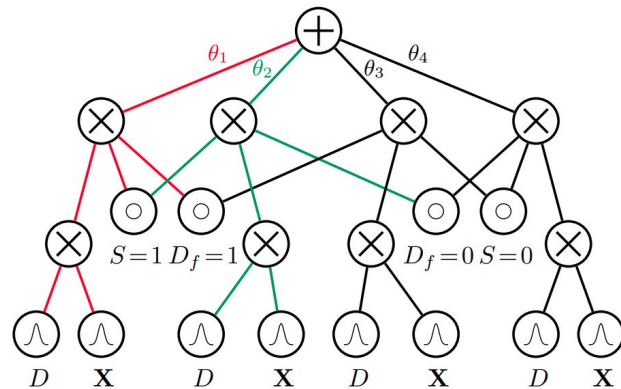
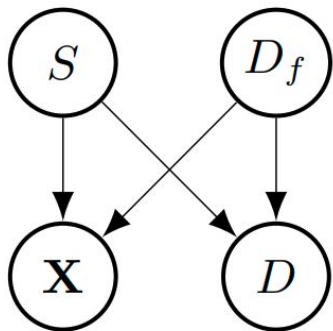


# Model-Based Algorithmic Fairness: FairPC

Learn classifier given

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

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



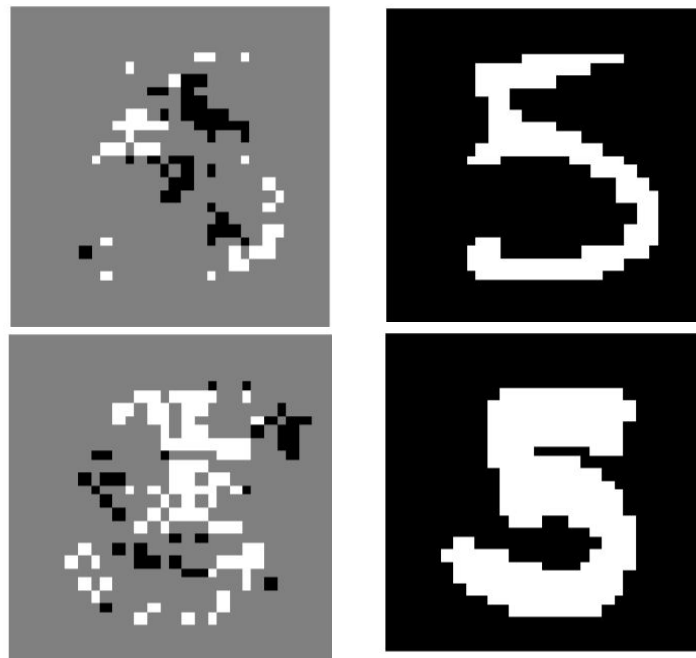


# Probabilistic Sufficient Explanations

Goal: explain an instance of classification

Choose a **subset** of features s.t.

1. Given only the explanation it is “probabilistically sufficient”  
*Under the feature distribution, it is likely to make the prediction to be explained*
2. It is minimal and “simple”





**Pure Logic**    **Probabilistic World Models**    **Pure Learning**

## A New Synthesis of Learning and Reasoning

*“Pure learning is brittle”*

bias, **algorithmic fairness**, interpretability, **explainability**, adversarial attacks, unknown unknowns, calibration, verification, **missing features**, missing labels, data efficiency, shift in distribution, general robustness and safety

We need to incorporate a sensible probabilistic model of the world

# Probabilistic Programs



# What are probabilistic programs?

```
let x = flip 0.5 in
let y = flip 0.7 in
let z = x || y in
let w = if z then
  my_func(x,y)
else
  ...
in
observe(z);
```

means “flip a coin, and output true with probability  $\frac{1}{2}$ ”

Standard (functional) programming constructs: let, if, ...

means “reject this execution if z is not true”

# Why Probabilistic Programming?

PPLs are proliferating



Pyro

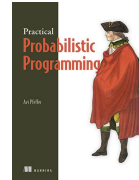
Edward



HackPPL



Stan




Figaro

Venture, Church, IBAL, WebPPL, Infer.NET, Tensorflow Probability, ProbLog, PRISM, LPADs, CProlog, CLP(BN), ICL, PHA, Primula, Storm, Gen, PRISM, PSI, Bean Machine, etc. ... *and many many more*

Programming languages are humanity's biggest knowledge representation achievement!

# Dice probabilistic programming language

<http://dicelang.cs.ucla.edu/>



## Dice

The dice probabilistic programming language

About GitHub

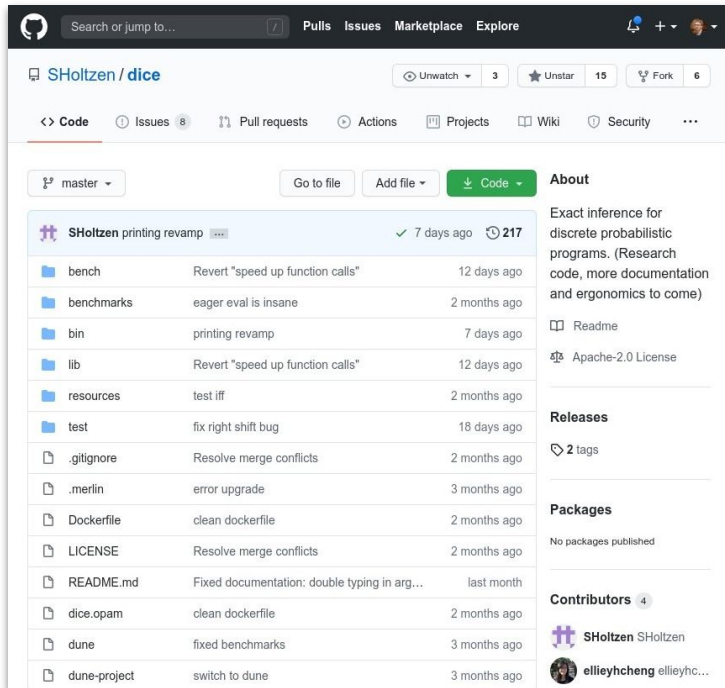
`dice` is a probabilistic programming language focused on fast exact inference for discrete probabilistic programs. For more information on `dice`, see the [about page](#).

Below is an online `dice` code demo. To run the example code, press the "Run" button.

```
1 fun sendChar(key: int(2), observation: int(2)) {
2   let gen = discrete(0.5, 0.25, 0.125, 0.125) in // sample a FooLang character
3   let enc = key + gen in // encrypt the character
4   observe observation == enc
5 }
6
7 // sample a uniform random key: A=0, B=1, C=2, D=3
8
9 let key = discrete(0.25, 0.25, 0.25, 0.25) in
10
11 // observe the ciphertext CCCC
12 let tmp = sendChar(key, int(2, 2)) in
13 let tmp = sendChar(key, int(2, 2)) in
14 let tmp = sendChar(key, int(2, 2)) in
15 let tmp = sendChar(key, int(2, 2)) in
16
17 key
```

Run

<https://github.com/SHoltzen/dice>



Search or jump to... Pulls Issues Marketplace Explore

SHoltzen / dice Unwatch 3 Unstar 15 Fork 6

<> Code Issues 8 Pull requests Actions Projects Wiki Security

master Go to file Add file Code

SHoltzen printing revamp 7 days ago 217

bench	Revert "speed up function calls"	12 days ago
benchmarks	eager eval is insane	2 months ago
bin	printing revamp	7 days ago
lib	Revert "speed up function calls"	12 days ago
resources	test iff	2 months ago
test	fix right shift bug	18 days ago
.gitignore	Resolve merge conflicts	2 months ago
.merlin	error upgrade	3 months ago
Dockerfile	clean dockerfile	2 months ago
LICENSE	Resolve merge conflicts	2 months ago
README.md	Fixed documentation: double typing in arg...	last month
dice.opam	clean dockerfile	2 months ago
dune	fixed benchmarks	3 months ago
dune-project	switch to dune	3 months ago

About

Exact inference for discrete probabilistic programs. (Research code, more documentation and ergonomics to come)

Readme Apache-2.0 License

Releases 2 tags

Packages No packages published

Contributors 4

SHoltzen SHoltzen ellieyhcheng ellieyh...

# What is a possible world?

```
let x = flip 0.4 in
let y = flip 0.7 in
let z = x || y in
let x = if z then
  x
else
  1
in (x,y)
```

Execution A

```
x=1
x=1, y=1
x=1, y=1, z=1

x=1, y=1, z=1

(1, 1)
```

$$P = 0.4 * 0.7$$

Execution B

```
x=1
x=1, y=0
x=1, y=0, z=1

x=1, y=0, z=1

(1,0)
```

$$P = 0.4 * 0.3$$

Execution C

```
x=0
x=0, y=1
x=0, y=1, z=1

x=0, y=1, z=1

(0,1)
```

$$P = 0.6 * 0.7$$

Execution D

```
x=0
x=0, y=0
x=0, y=0, z=0

x=1, y=0, z=0

(1,0)
```

$$P = 0.6 * 0.3$$

# Why should I care? I like PGMs

- Better abstraction:
  - Beyond variable-level dependencies
  - modularity through functions  
reuse (cf. templative graphical models)
  - intuitive language for local structure; arithmetic
  - data structures
  - first-class observations



# First-Class Observations, Functions

---

```
1 fun EncryptChar(key:int, obs:char):Bool {
2   let randomChar = ChooseChar() in
3   let ciphertext = (randomChar + key) % 26 in
4   let _ = observe ciphertext = obs in
5   true}
6 let k = UniformInt(0, 25) in
7 let _ = EncryptChar(k, 'H') in ...
8 let _ = EncryptChar(k, 'D') in k
```

---

Frequency Analyzer for a Caesar cipher in Dice

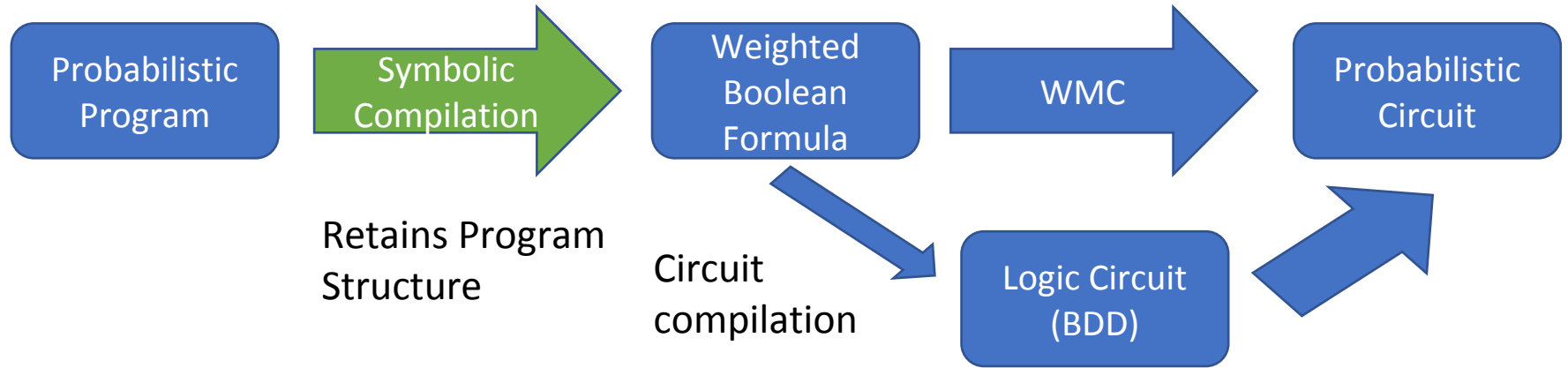
# *What do PGMs bring to the table?*

1. Real programs have inherently discrete structure (e.g. if-statements)
2. Discrete structure is inherent in many domains (graphs, text/topic models, ranking, etc.)
3. Many existing PPLs assume smooth and differentiable densities and do not handle these programs correctly.

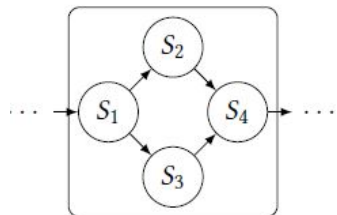
**Discrete probabilistic programming is the important unsolved open problem!**

**PGM community knows how to solve this!**

# Symbolic Compilation to Probabilistic Circuits



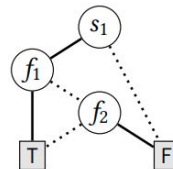
# Inference in Dice



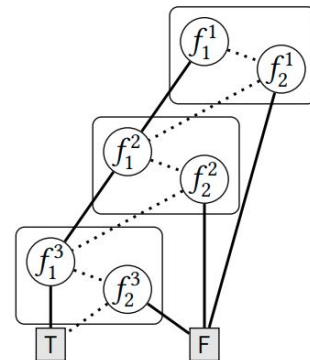
(a) Network diagram.

```
1 fun diamond( $s_1$ :Bool):Bool {  
2   let route = flip1 0.5 in  
3   let  $s_2$  = if route then  $s_1$  else F in  
4   let  $s_3$  = if route then F else  $s_1$  in  
5   let drop = flip2 0.0001 in  
6    $s_2 \vee (s_3 \wedge \neg \text{drop})$ }  
7 diamond(diamond(diamond(T)))
```

(b) Probabilistic program defining the network.



(c) diamond function.



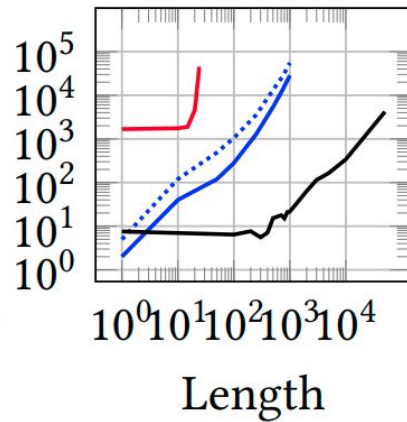
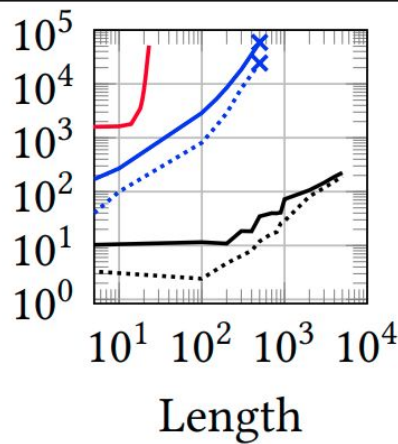
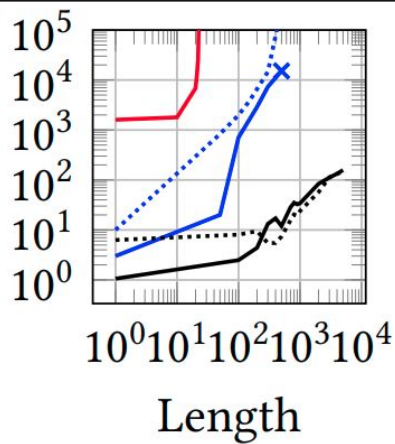
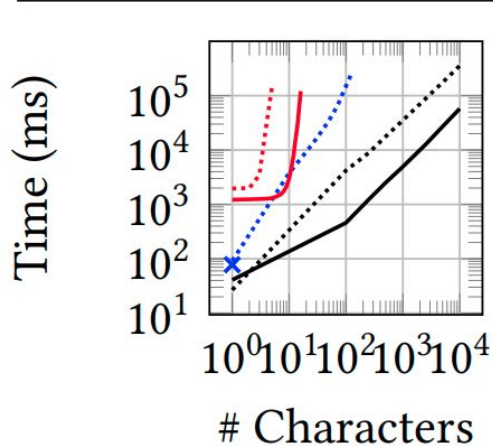
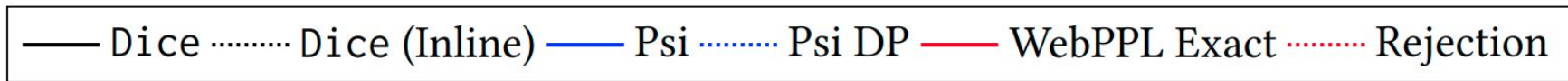
(d) Final BDD.

## Network Verification

# PPL benchmarks from PL community

Benchmark	Psi (ms)	DP (ms)	Dice (ms)	# Paths	BDD Size
Grass	167	57	<b>1.0</b>	$9.5 \times 10^1$	15
Burglar Alarm	98	10	<b>1.1</b>	$2.5 \times 10^2$	11
Coin Bias	94	23	<b>1.0</b>	4	13
Noisy Or	81	152	<b>1.0</b>	$1.6 \times 10^4$	35
Evidence1	48	32	<b>1.0</b>	9	5
Evidence2	59	28	<b>1.0</b>	9	6
Murder Mystery	193	75	<b>1.0</b>	$1.6 \times 10^1$	6

# Scalable Inference

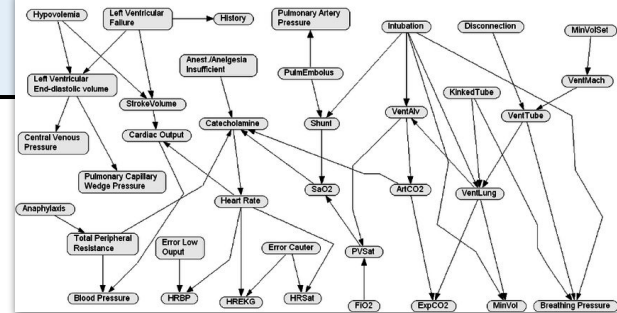


# Scalable Inference

Benchmark	Psi (ms)	DP (ms)	Dice (ms)	# Parameters	# Paths	BDD Size
Cancer [48]	772	46	<b>1.0</b>	10	$1.1 \times 10^3$	28
Survey [73]	2477	152	<b>2.0</b>	21	$1.3 \times 10^4$	73
Alarm [5]	<b>X</b>	<b>X</b>	<b>9.0</b>	509	$1.0 \times 10^{36}$	$1.3 \times 10^3$
Insurance [7]	<b>X</b>	<b>X</b>	<b>75.0</b>	984	$1.2 \times 10^{40}$	$1.0 \times 10^5$
Hepar2 [63]	<b>X</b>	<b>X</b>	<b>54.0</b>	1453	$2.9 \times 10^{69}$	$1.3 \times 10^3$
Hailfinder [1]	<b>X</b>	<b>X</b>	<b>526.0</b>	2656	$2.0 \times 10^{76}$	$6.5 \times 10^4$
Pigs	<b>X</b>	<b>X</b>	<b>32.0</b>	5618	$7.3 \times 10^{492}$	35
Water [43]	<b>X</b>	<b>X</b>	<b>2926.0</b>	$1.0 \times 10^4$	$3.2 \times 10^{54}$	$5.1 \times 10^4$
Munin [3]	<b>X</b>	<b>X</b>	<b>1945.0</b>	$8.1 \times 10^5$	$2.1 \times 10^{1622}$	$1.1 \times 10^4$

# Alarm Bayesian Network

```
let HYPOVOLEMIA = flip 0.2 in
let LVFAILURE = flip 0.05 in
let STROKEVOLUME =
  if (HYPOVOLEMIA) then
    (if (LVFAILURE) then (discrete(0.98,0.01,0.01)) else (discrete(0.50,0.49,0.01)))
  else
    (if (LVFAILURE) then (discrete(0.95,0.04,0.01)) else (discrete(0.05,0.90,0.05)))
in
let LVEDVOLUME =
  if (HYPOVOLEMIA) then
    (if (LVFAILURE) then (discrete(0.95,0.04,0.01)) else (discrete(0.01,0.09,0.90)))
  else
    (if (LVFAILURE) then (discrete(0.98,0.01,0.01)) else (discrete(0.05,0.90,0.05)))
in
...
```





# Why should I care? I like PGMs

- Better abstraction:
  - Beyond variable-level dependencies
  - modularity through functions  
reuse (cf. templative graphical models)
  - intuitive language for local structure; arithmetic
  - data structures
  - first-class observations
- Better inference? correctness? analysis?  
**import PL.\***

# Denotational Semantics

- Goal: associate with every expression “e” a semantic object.
- Notation: semantic bracket:  $\llbracket \cdot \rrbracket$ 
  - In Bayesian network:  $\llbracket \text{BN} \rrbracket = \text{Pr}_{\text{BN}}(\cdot)$
  - In probabilistic programs:  $\llbracket e \rrbracket(\cdot)$  for all expressions
  - Accepting and distributional semantics:

$$\llbracket e \rrbracket_A \triangleq \sum_v \llbracket e \rrbracket(v), \quad \llbracket e \rrbracket_D(v) \triangleq \frac{1}{\llbracket e \rrbracket_A} \llbracket e \rrbracket(v)$$

- Idea: don't need to run ‘flip 0.4’ infinite times to know meaning

# Denotational Semantics + Formal Inference Rules

$$\llbracket v_1 \rrbracket (v) \triangleq (\delta(v_1))(v) \quad \llbracket \text{fst } (v_1, v_2) \rrbracket (v) \triangleq (\delta(v_1))(v) \quad \llbracket \text{snd } (v_1, v_2) \rrbracket (v) \triangleq (\delta(v_2))(v)$$

$$\llbracket \text{if } v_g \text{ then } e_1 \text{ else } e_2 \rrbracket (v) \triangleq \begin{cases} \llbracket e_1 \rrbracket (v) & \text{if } v_g = \top \\ \llbracket e_2 \rrbracket (v) & \text{if } v_g = \text{F} \\ 0 & \text{otherwise} \end{cases} \quad \llbracket \text{flip } \theta \rrbracket (v) \triangleq \begin{cases} \theta & \text{if } v = \top \\ 1 - \theta & \text{if } v = \text{F} \\ 0 & \text{otherwise} \end{cases}$$

$$\llbracket \text{observe } v_1 \rrbracket (v) \triangleq \begin{cases} 1 & \text{if } v_1 = \top \text{ and } v = \top, \\ 0 & \text{otherwise} \end{cases} \quad \llbracket f(v_1) \rrbracket (v) \triangleq ((T(f))(v_1))(v)$$

$$\llbracket \text{let } x = e_1 \text{ in } e_2 \rrbracket (v) \triangleq \sum_{v'} \llbracket e_1 \rrbracket (v') \times \llbracket e_2[x \mapsto v'] \rrbracket (v)$$

$$\frac{}{\top \rightsquigarrow (\top, \top, \emptyset)} \text{ (C-TRUE)} \quad \frac{}{\text{F} \rightsquigarrow (\text{F}, \top, \emptyset)} \text{ (C-FALSE)} \quad \frac{}{x \rightsquigarrow (\mathbf{x}, \top, \emptyset)} \text{ (C-IDENT)}$$

$$\frac{\text{fresh } \mathbf{f}}{\text{flip } \theta \rightsquigarrow (\mathbf{f}, \top, (\mathbf{f} \mapsto \theta, \top, \bar{\mathbf{f}} \mapsto 1 - \theta))} \text{ (C-FLIP)} \quad \frac{\text{aexp } \rightsquigarrow (\varphi, \top, \emptyset)}{\text{observe aexp } \rightsquigarrow (\top, \varphi, \emptyset)} \text{ (C-OBS)}$$

$$\frac{\text{aexp } \rightsquigarrow (\varphi_g, \top, \emptyset) \quad e_T \rightsquigarrow (\varphi_T, \gamma_T, w_T) \quad e_E \rightsquigarrow (\varphi_E, \gamma_E, w_E)}{\text{if aexp then } e_T \text{ else } e_E \rightsquigarrow (((\varphi_g \wedge \varphi_T) \vee ((\bar{\varphi}_g \wedge \varphi_E), ((\varphi_g \wedge \gamma_T) \vee ((\bar{\varphi}_g \wedge \gamma_E), w_T \cup w_E)))} \text{ (C-ITE)}$$

$$\frac{e_1 \rightsquigarrow (\varphi_1, \gamma_1, w_1) \quad e_2 \rightsquigarrow (\varphi_2, \gamma_2, w_2)}{\text{let } x = e_1 \text{ in } e_2 \rightsquigarrow (\varphi_2[\mathbf{x} \mapsto \varphi_1], \gamma_1 \wedge \gamma_2[\mathbf{x} \mapsto \varphi_1], w_1 \cup w_2)} \text{ (C-LET)}$$

# Provably Correct Inference!

$$\begin{array}{c}
 \text{fresh } f_1 \\
 \hline
 \text{flip 0.1} \rightsquigarrow \left( \begin{array}{c} \text{f}_1 \\ \text{T} \text{---} \text{---} \text{F} \end{array}, \text{T}, w_1 \right) \\
 \hline
 \text{fresh } f_2 \\
 \hline
 \text{flip 0.4} \rightsquigarrow \left( \begin{array}{c} \text{x} \\ \text{T} \text{---} \text{---} \text{F} \end{array}, \text{T}, \emptyset \right) \quad \text{fresh } f_2 \\
 \hline
 \text{flip 0.4} \rightsquigarrow \left( \begin{array}{c} \text{f}_2 \\ \text{T} \text{---} \text{---} \text{F} \end{array}, \text{T}, w_2 \right) \\
 \hline
 \text{flip 0.4} \vee \text{flip 0.1} \rightsquigarrow \left( \begin{array}{c} \text{x} \\ \text{f}_2 \text{---} \text{---} \text{T} \\ \text{F} \text{---} \end{array}, \text{T}, w_2 \right) \\
 \hline
 \text{let } x = \text{flip 0.1} \text{ in flip 0.4} \vee \text{flip 0.1} \rightsquigarrow \left( \begin{array}{c} \text{f}_1 \\ \text{f}_2 \text{---} \text{---} \text{T} \\ \text{F} \text{---} \end{array}, \text{T}, w_1 \cup w_2 \right)
 \end{array}$$

# Better Inference?

Exploit modularity

1. AI modularity:  
Discover contextual independencies and factorize
2. PL modularity:  
Compile procedure summaries and reuse at each call site

Reason about programs! Compiler optimizations.

Quick preview:

3. Flip hoisting optimization
4. Eager compilation

# From programs to circuits directly:

---

```
1 let z = flip1 0.5 in
2 let x = if z then flip2 0.6 else flip3 0.7 in
3 let y = if z then flip4 0.7 else x in (x, y)
```

---

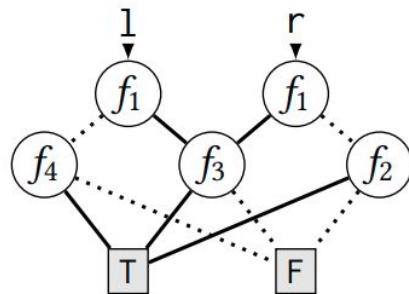
(a) Context-specific independence.

---

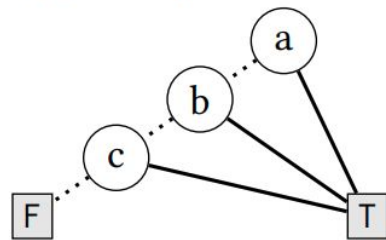
```
1 fun foo(a:Bool, b:Bool, c:Bool):Bool {
2   a ∨ b ∨ c
3 }
```

---

(c) Structure without independence.



(b) Compiled BDD.



(d) Compiled BDD.

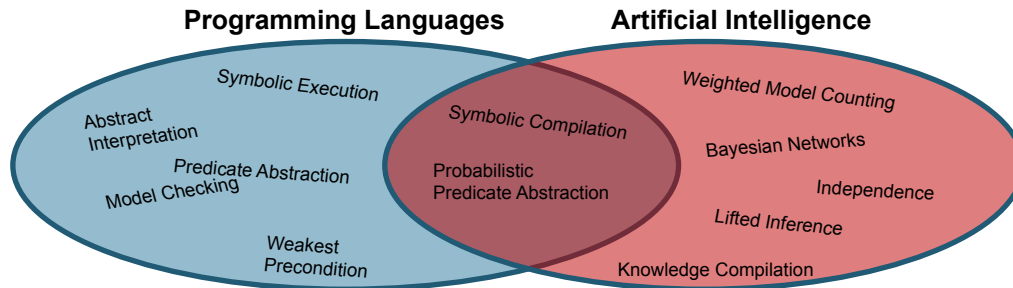
# Compiler Optimizations (sneak preview)

Benchmark	Naive compilation	determinism	flip hoisting + determinism	Eager + flip lifting	<i>Ace baseline</i>
alarm	156	140	83	69	422
water	56,267	65,975	1509	941	605
insurance	140	100	100	128	492
hepar2	95	80	80	80	495
pigs	3,772	2490	2112	186	985
munin	>1,000,000	>1,000,000	109,687	16,536	3,500

Inference time in milliseconds

# Conclusions

- Are we already in the age of computational abstractions?
- Probabilistic circuits for learning deep tractable probabilistic models
- Probabilistic programs as the new probabilistic knowledge representation language





Thanks