

AI can learn from data. But can it learn to reason?

Guy Van den Broeck

Tong Lecture Series at Peking University - Sep 6 2022

The AI Dilemma



Integrate reasoning into modern deep learning algorithms

Outline

1. The paradox of learning to reason from data

~~deep learning~~

2. Tractable deep generative models

probabilistic reasoning + deep learning

3. Learning with symbolic knowledge

logical reasoning + deep learning

Outline

1. The paradox of learning to reason from data

~~deep learning~~

2. Tractable deep generative models

probabilistic reasoning + deep learning

3. Learning with symbolic knowledge

logical reasoning + deep learning

Can Language Models Perform Logical Reasoning?

Language Models achieve high performance on various “reasoning” benchmarks in NLP.

<p>Kristin and her son Justin went to visit her mother Carol on a nice Sunday afternoon. They went out for a movie together and had a good time.</p> 	<p>Q: How is Carol related to Justin ?</p> <p>A: Carol is the grandmother of Justin</p> 
--	---

Reasoning Example
from the CLUTRR
dataset

It is unclear whether they solve the tasks following the rules of logical deduction.

Language Models:

input → ? → *Carol is the grandmother of Justin.*

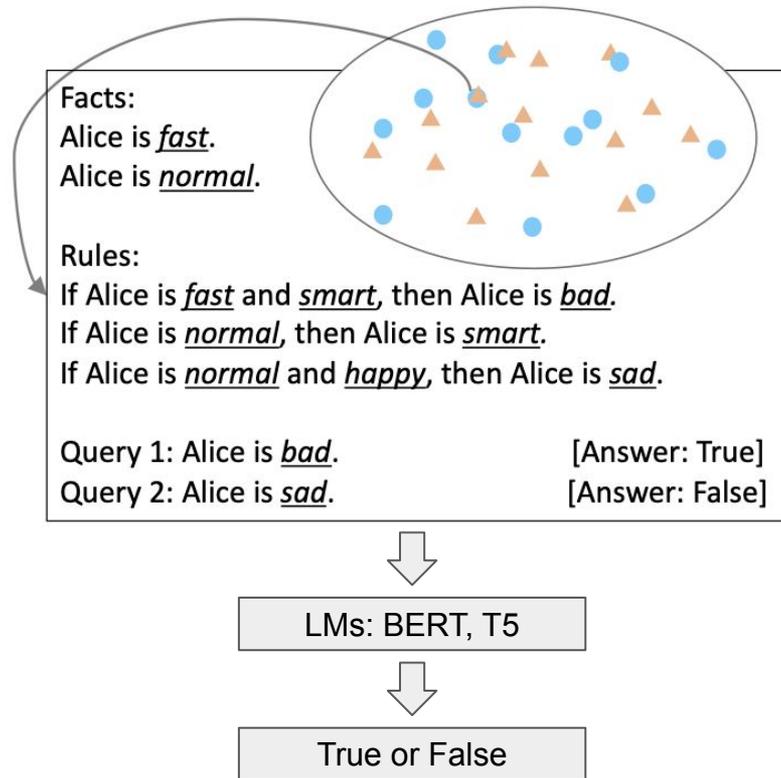
Reasoning:

input → *Justin is Kristin's son; Carol is Kristin's mother;* → *Carol is Justin's mother's mother; if X is Y's mother's mother then X is Y's grandmother* → *Carol is the grandmother of Justin.*

Problem Setting: SimpleLogic

The easiest of reasoning problems:

1. **Propositional logic** fragment
 - a. bounded vocabulary & number of rules
 - b. bounded reasoning depth (≤ 6)
 - c. finite space ($\approx 10^{360}$)
2. **No language variance**: templated language
3. **Self-contained**
No prior knowledge
4. **Purely symbolic** predicates
No shortcuts from word meaning
5. **Tractable** logic (definite clauses)
Can always be solved efficiently

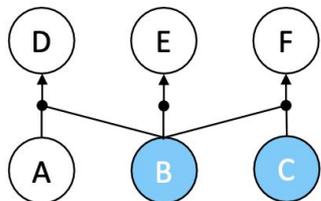


Training a BERT model on SimpleLogic

(1) Randomly sample facts & rules.

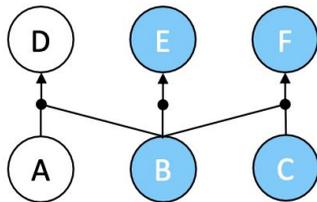
Facts: B, C

Rules: $A, B \rightarrow D$. $B \rightarrow E$. $B, C \rightarrow F$.



Rule-Priority

(2) Compute the correct labels for all predicates given the facts and rules.



Label-Priority



(1) Randomly assign labels to predicates.

True: B, C, E, F.

False: A, D.

(2) Set B, C (randomly chosen among B, C, E, F) as facts and sample rules (randomly) consistent with the label assignments.

Test accuracy for different reasoning depths

Test	0	1	2	3	4	5	6
RP	99.9	99.8	99.7	99.3	98.3	97.5	95.5

Test	0	1	2	3	4	5	6
LP	100.0	100.0	99.9	99.9	99.7	99.7	99.0

Has BERT learned to reason from data?

1. Easiest of reasoning problems (no variance, self-contained, purely symbolic, tractable)
2. RP/LP data covers the whole problem space
3. The learned model has almost 100% test accuracy
4. There exist BERT parameters that compute the ground-truth reasoning function:

Theorem 1: *For a BERT model with n layers and 12 attention heads, by construction, there exists a set of parameters such that the model can correctly solve any reasoning problem in SimpleLogic that requires at most $n - 2$ steps of reasoning.*

**Surely, under these conditions,
BERT has learned the ground-truth reasoning function!**



The Paradox of Learning to Reason from Data

Train	Test	0	1	2	3	4	5	6
RP	RP	99.9	99.8	99.7	99.3	98.3	97.5	95.5
	LP	99.8	99.8	99.3	96.0	90.4	75.0	57.3
LP	RP	97.3	66.9	53.0	54.2	59.5	65.6	69.2
	LP	100.0	100.0	99.9	99.9	99.7	99.7	99.0

The BERT model trained on one distribution fails to generalize to the other distribution within the same problem space.



1. If BERT **has learned** to reason, it should not exhibit such generalization failure.
2. If BERT **has not learned** to reason, it is baffling how it achieves near-perfect in-distribution test accuracy.

Why? Statistical Features

Monotonicity of entailment:

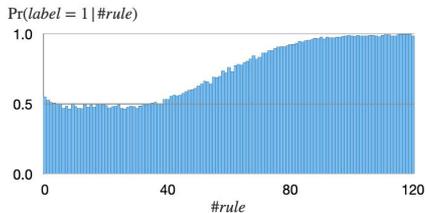
Any rules can be freely added to the hypothesis of any proven fact.



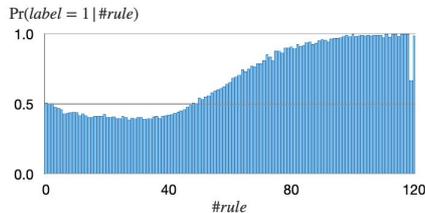
The more rules given, the more likely a predicate will be proved.



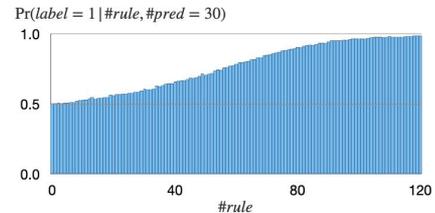
$\Pr(\text{label} = \text{True} \mid \text{Rule \#} = x)$ should increase (roughly) monotonically with x



(a) Statistics for examples generated by Rule-Priority (RP).



(b) Statistics for examples generated by Label-Priority (LP).



(c) Statistics for examples generated by uniform sampling;

BERT leverages statistical features to make predictions

RP_b downsamples from RP such that $\Pr(\text{label} = \text{True} \mid \text{rule\#} = x) = 0.5$ for all x

Train	Test	0	1	2	3	4	5	6
	RP	99.9	99.8	99.7	99.3	98.3	97.5	95.5
RP	RP_b	99.0	99.3	98.5	97.5	96.7	93.5	88.3

1. Accuracy drop from RP to RP_b indicates that **the model is using rule# as a statistical feature to make predictions.**
2. Though removing one statistical feature from training data can help with model generalization, there are potentially countless statistical features and it is computationally infeasible to jointly remove them.

First Conclusion

Experiments unveil the fundamental difference between

1. learning to reason, and
2. learning to achieve high performance on benchmarks using statistical features.

Outline

1. The paradox of learning to reason from data

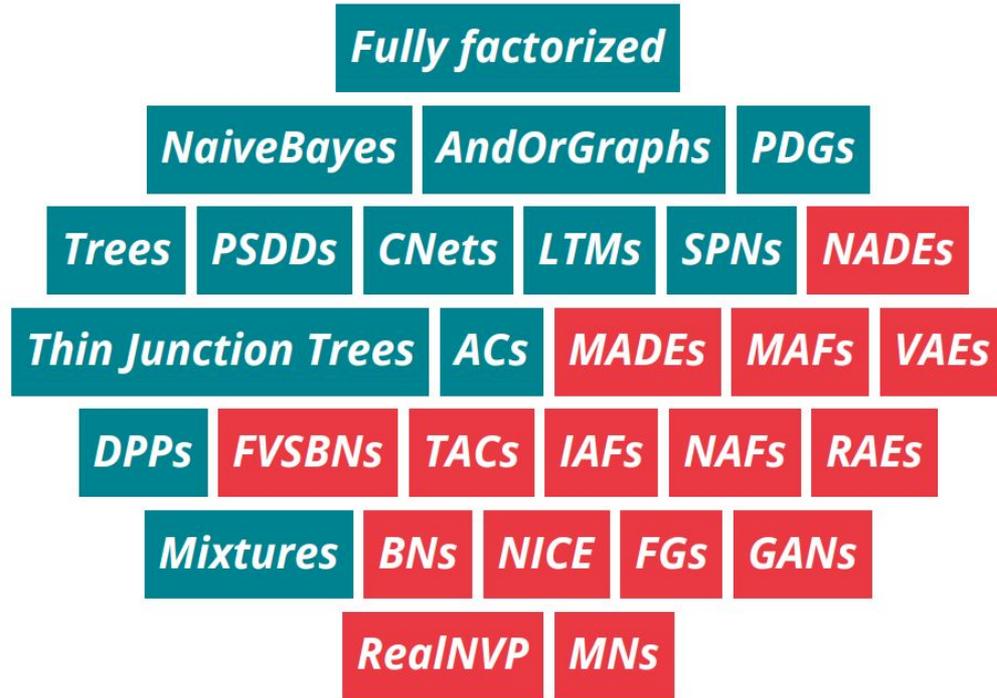
~~deep learning~~

2. **Tractable deep generative models**

probabilistic reasoning + deep learning

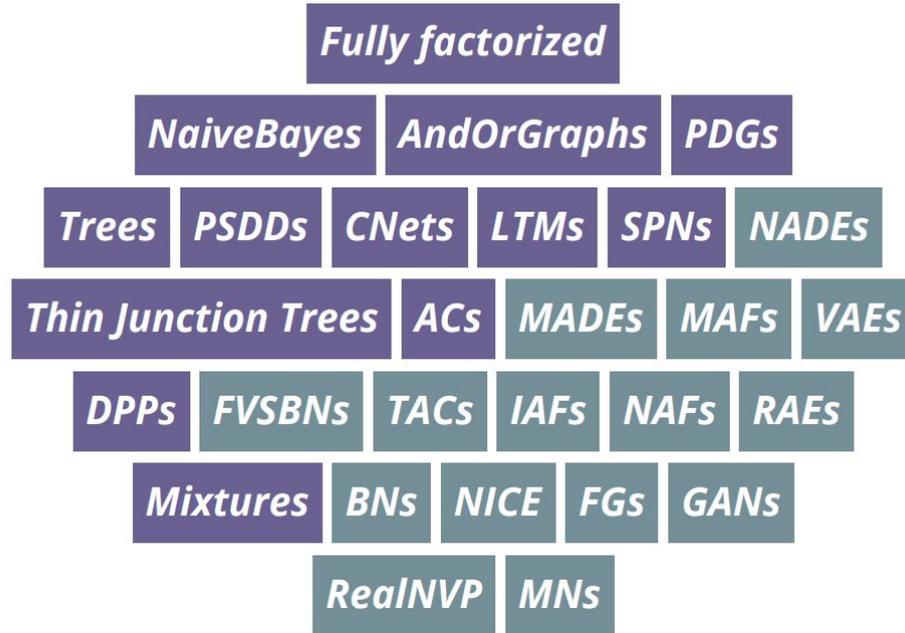
3. Learning with symbolic knowledge

logical reasoning + deep learning



Intractable and ***tractable*** models

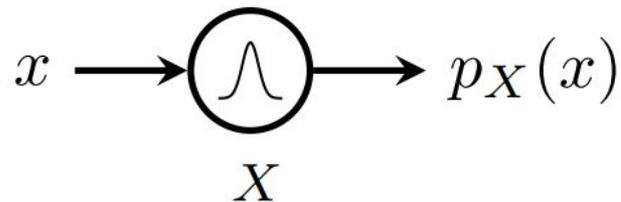
Probabilistic circuits



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

Probabilistic circuits

computational graphs that recursively define distributions



Simple distributions are tractable “black boxes” for:

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

Probabilistic circuits

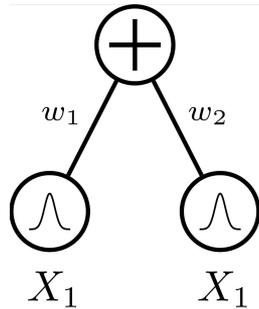
computational graphs that recursively define distributions



$\neg X$



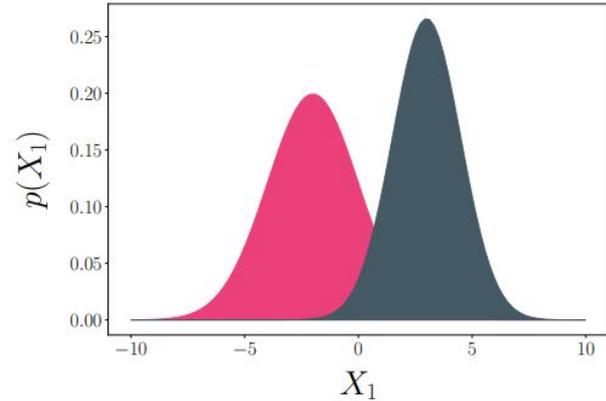
X_1



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

\Rightarrow

mixtures



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

Probabilistic circuits

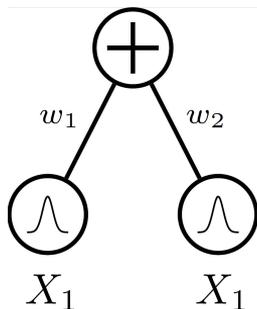
computational graphs that recursively define distributions



$\neg X$



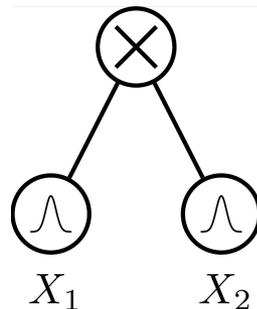
X_1



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

\Rightarrow

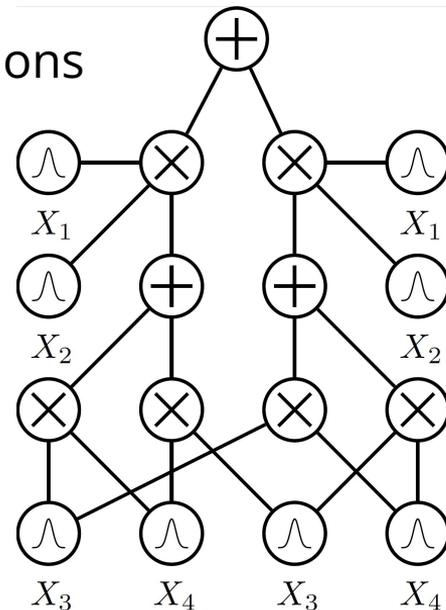
mixtures



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

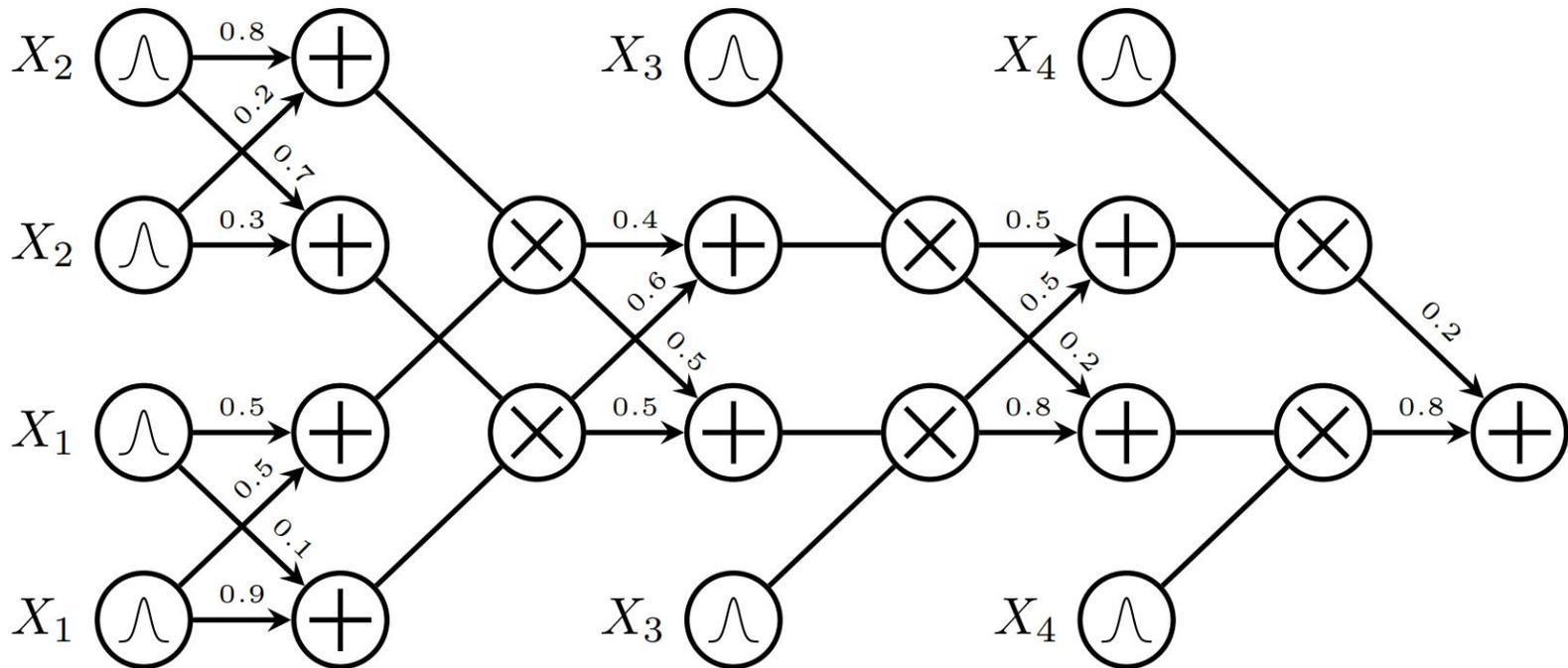
\Rightarrow

factorizations



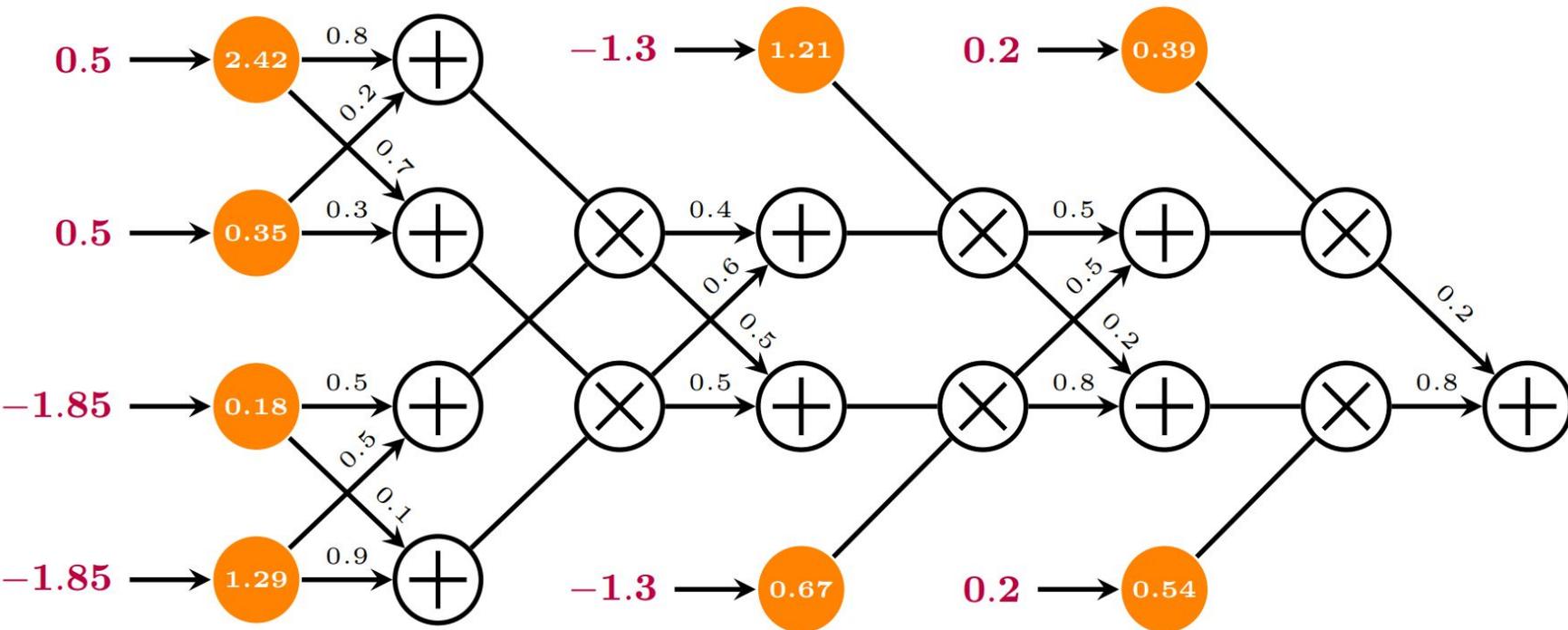
Likelihood

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



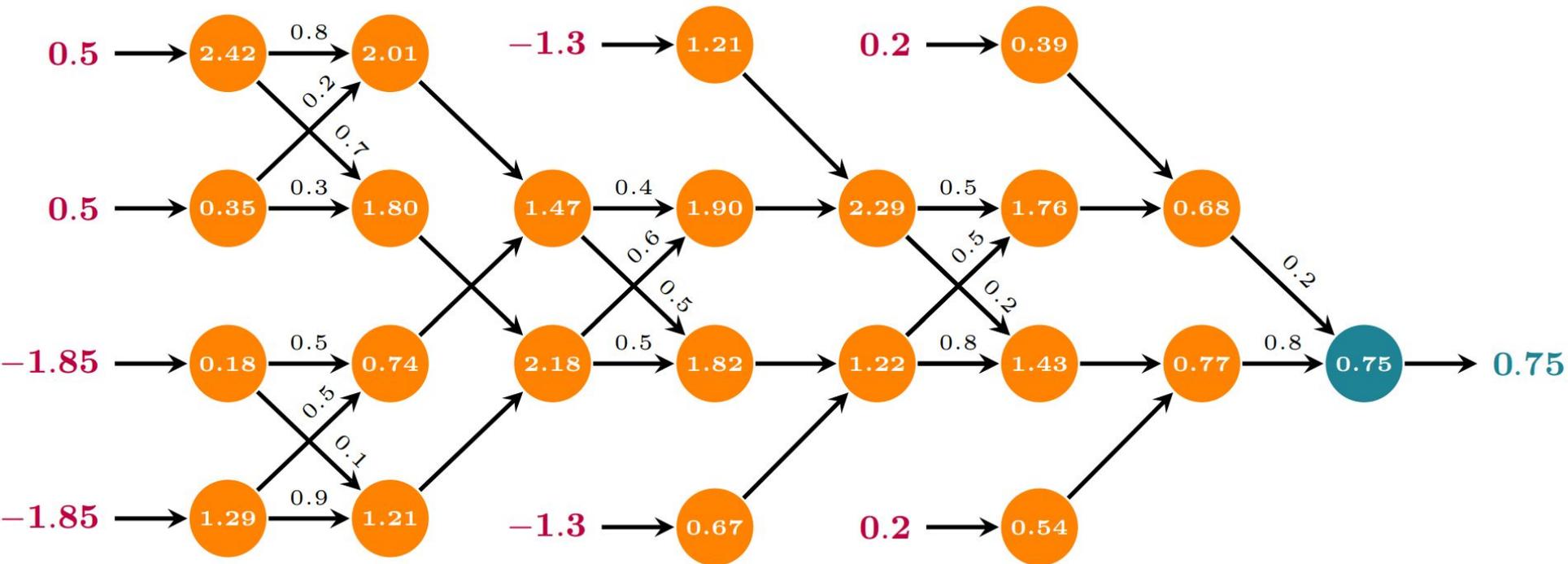
Likelihood

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

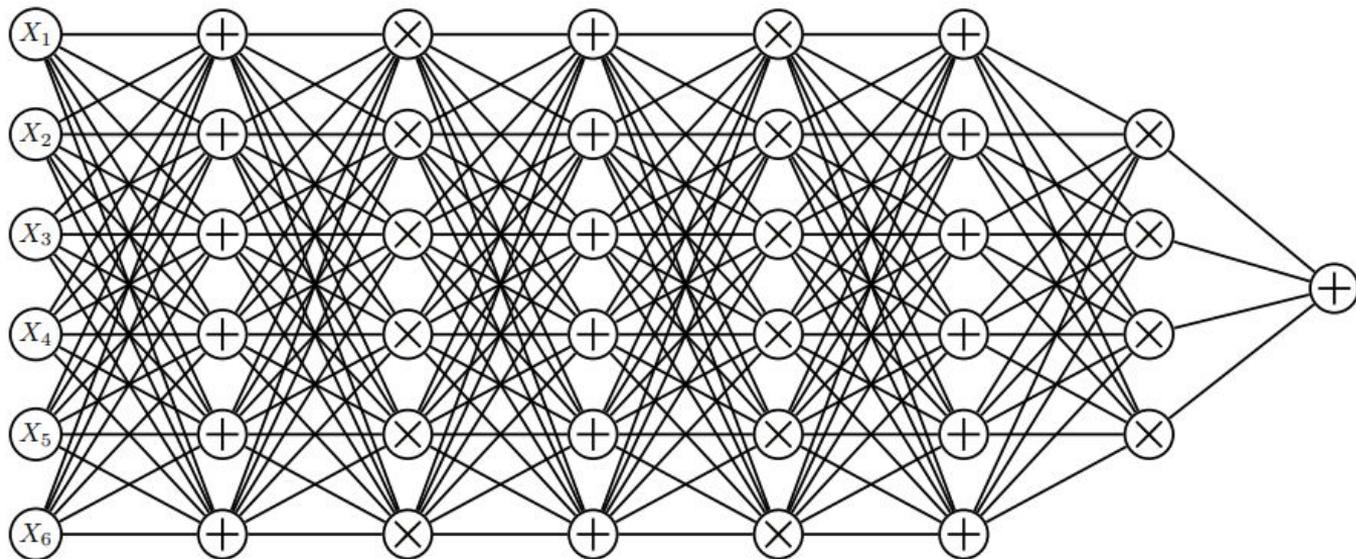


Likelihood

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

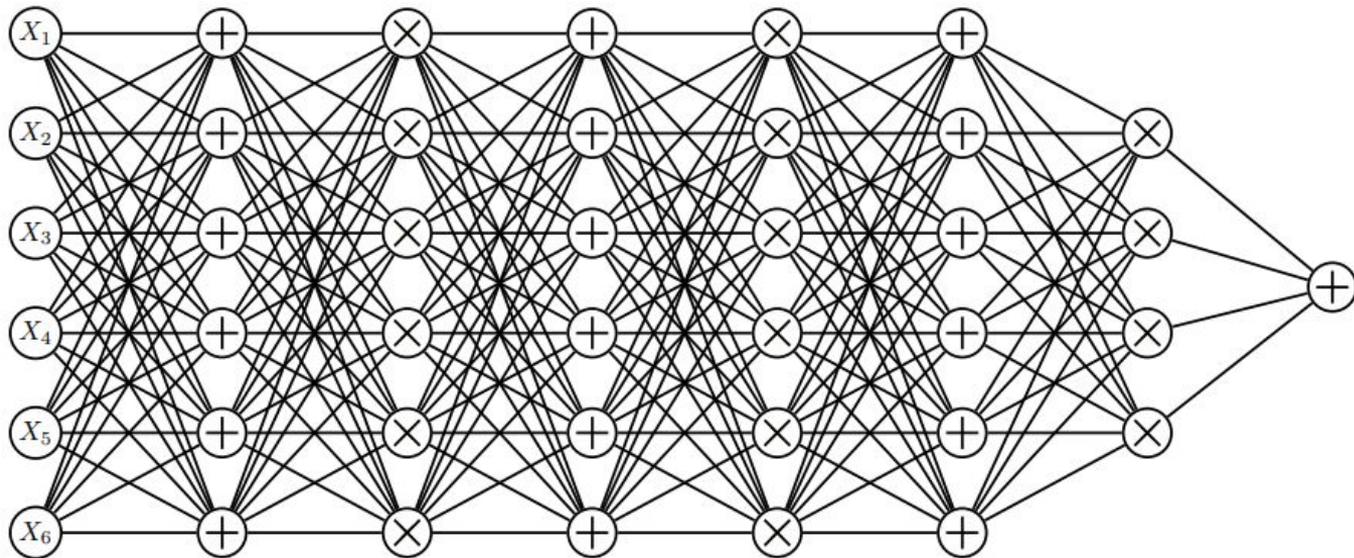


Just sum, products and distributions?



just arbitrarily compose them like a neural network!

Just sum, products and distributions?



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

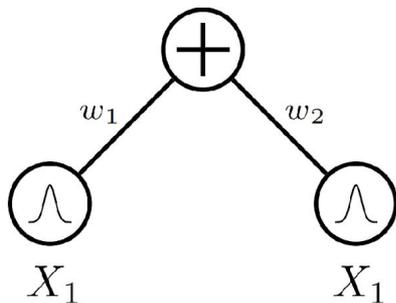


structural constraints needed for tractability

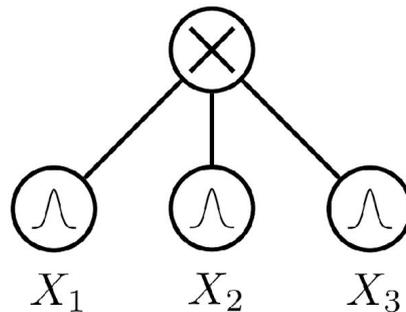
Tractable marginals

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

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



smooth circuit



decomposable circuit

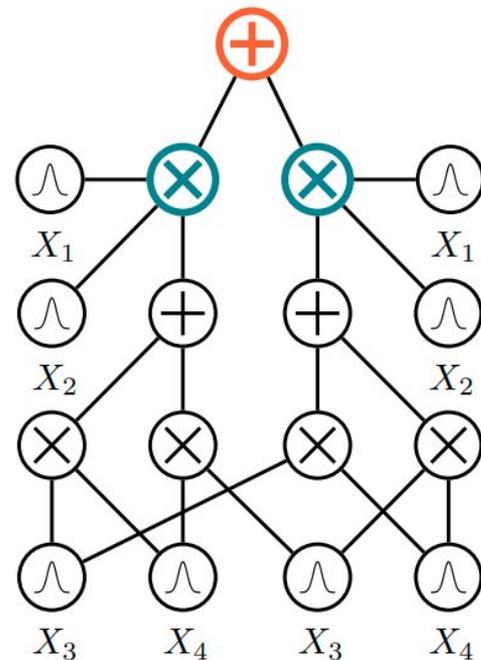
Smoothness + decomposability = tractable MAR

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

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

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

\Rightarrow integrals are "pushed down" to children

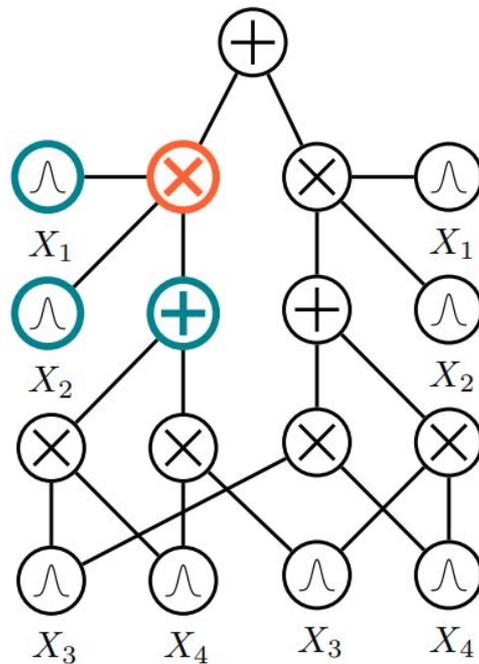


Smoothness + decomposability = tractable MAR

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

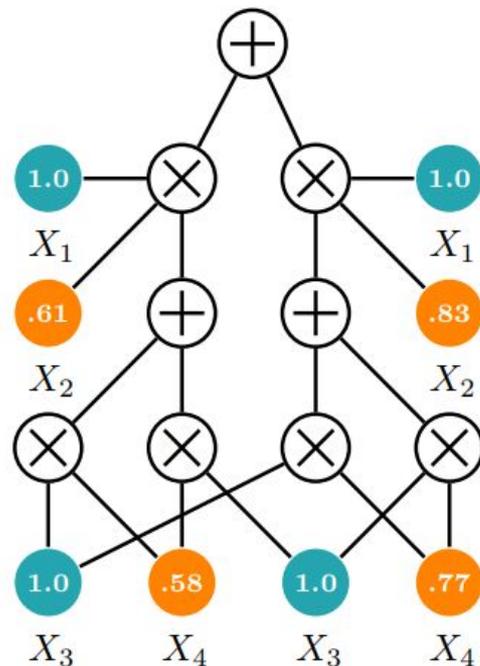
E.g. to compute $p(x_2, x_4)$:

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)



Smoothness + decomposability = tractable MAR

Forward pass evaluation for MAR

\Rightarrow linear in circuit size!

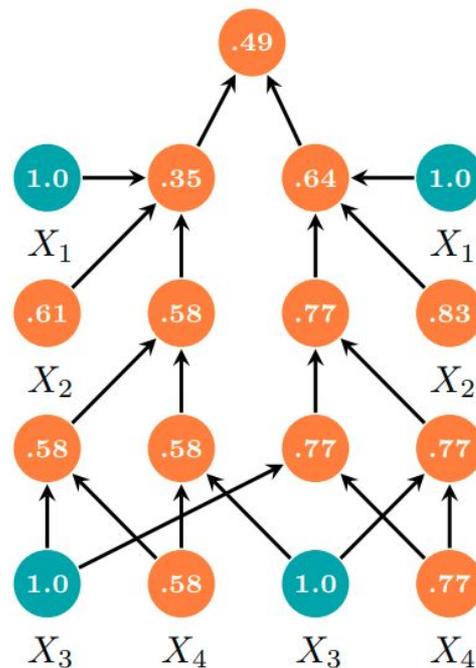
E.g. to compute $p(x_2, x_4)$:

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

\Rightarrow for normalized leaf distributions: 1.0

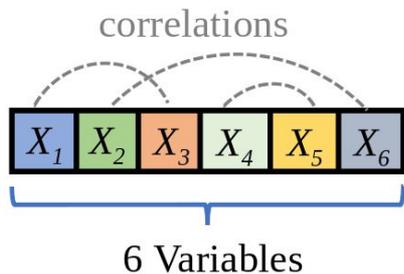
leaves over X_2 and X_4 output **EVI**

feedforward evaluation (bottom-up)

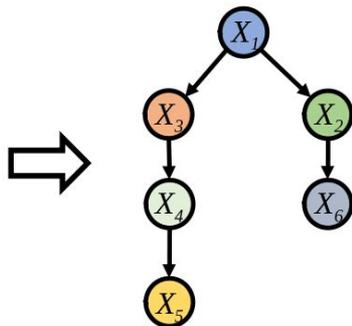


Learning Expressive Probabilistic Circuits

Hidden Chow-Liu Trees

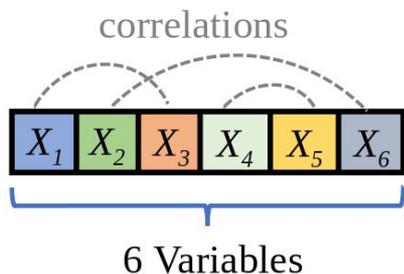


Learned **CLT structure**
captures strong pairwise
dependencies

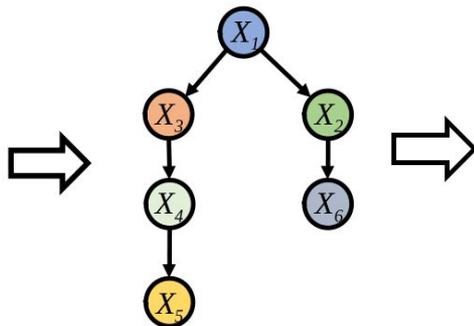


Learning Expressive Probabilistic Circuits

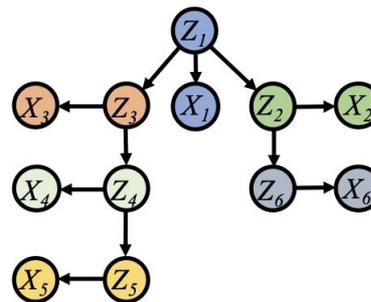
Hidden Chow-Liu Trees



Learned **CLT** structure captures strong pairwise dependencies



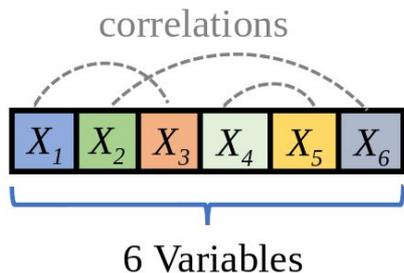
Learned **HCLT** structure



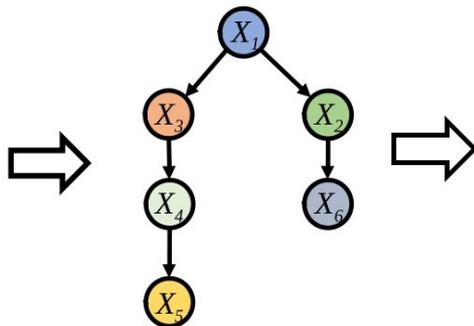
⇒ **Compile into an equivalent PC**

Learning Expressive Probabilistic Circuits

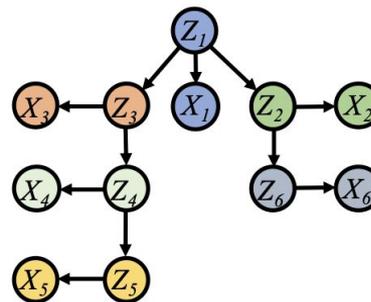
Hidden Chow-Liu Trees



Learned **CLT structure**
captures strong pairwise
dependencies



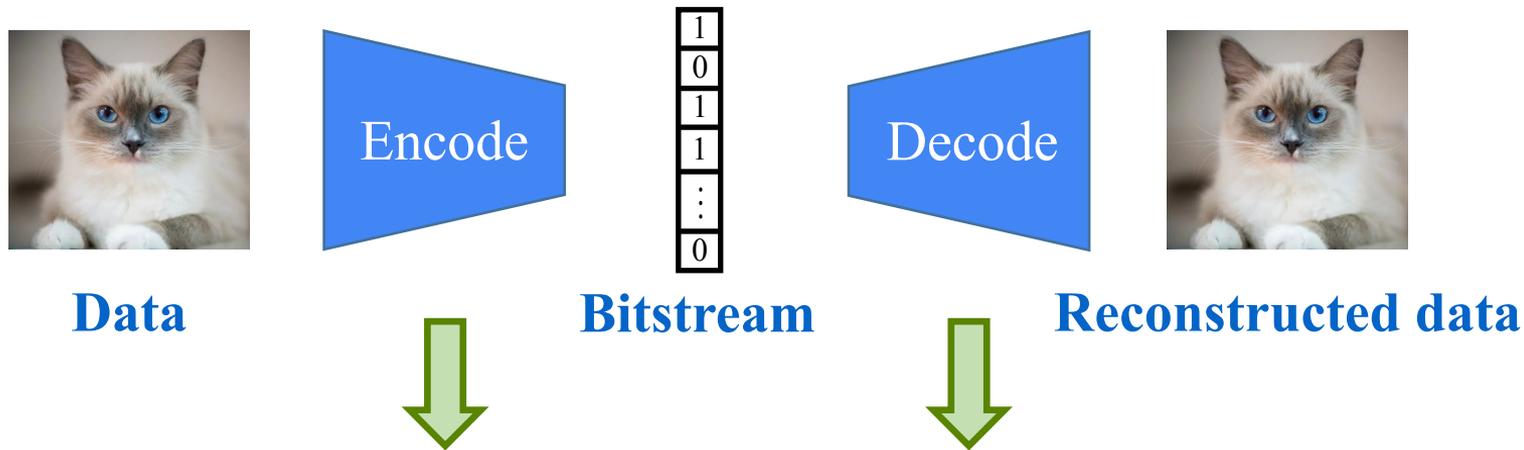
Learned **HCLT structure**



⇒ **Compile** into an
equivalent PC

⇒ Mini-batch Stochastic
Expectation Maximization

Lossless Data Compression



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

+

Efficient coding algorithm



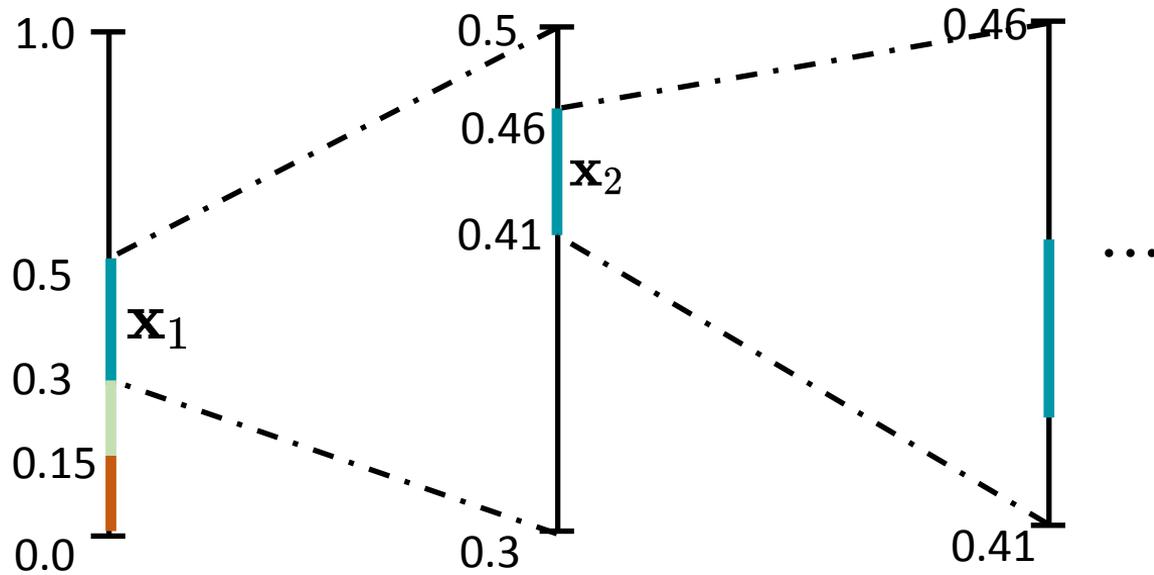
Determines the theoretical limit of compression rate



How close we can approach the theoretical limit

A Typical Streaming Code – Arithmetic Coding

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



Compress \mathbf{x}_1 with
 $-\log p(\mathbf{x}_1)$ bits

Compress \mathbf{x}_2 with
 $-\log p(\mathbf{x}_2|\mathbf{x}_1)$ bits

Compress \mathbf{x}_3 with
 $-\log p(\mathbf{x}_3|\mathbf{x}_1, \mathbf{x}_2)$ bits

Need to compute

$$p(X_1 < x_1)$$

$$p(X_1 \leq x_1)$$

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

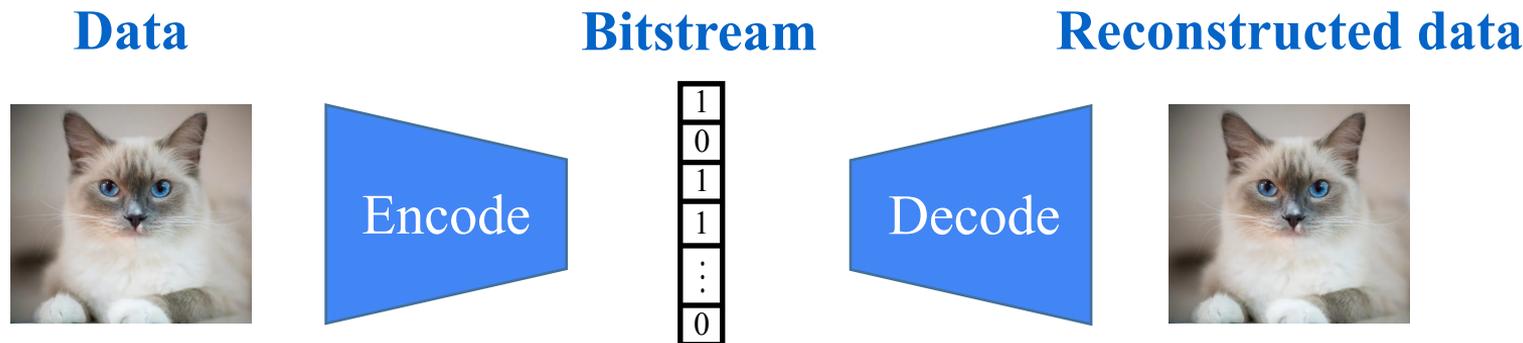
$$p(X_2 \leq x_2 | x_1)$$

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

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

\vdots

Lossless Neural Compression with Probabilistic Circuits



Probabilistic Circuits

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

Arithmetic Coding:

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

Lossless Neural Compression with Probabilistic Circuits

SoTA compression rates

Dataset	HCLT (ours)	IDF	BitSwap	BB-ANS	JPEG2000	WebP	McBits
MNIST	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

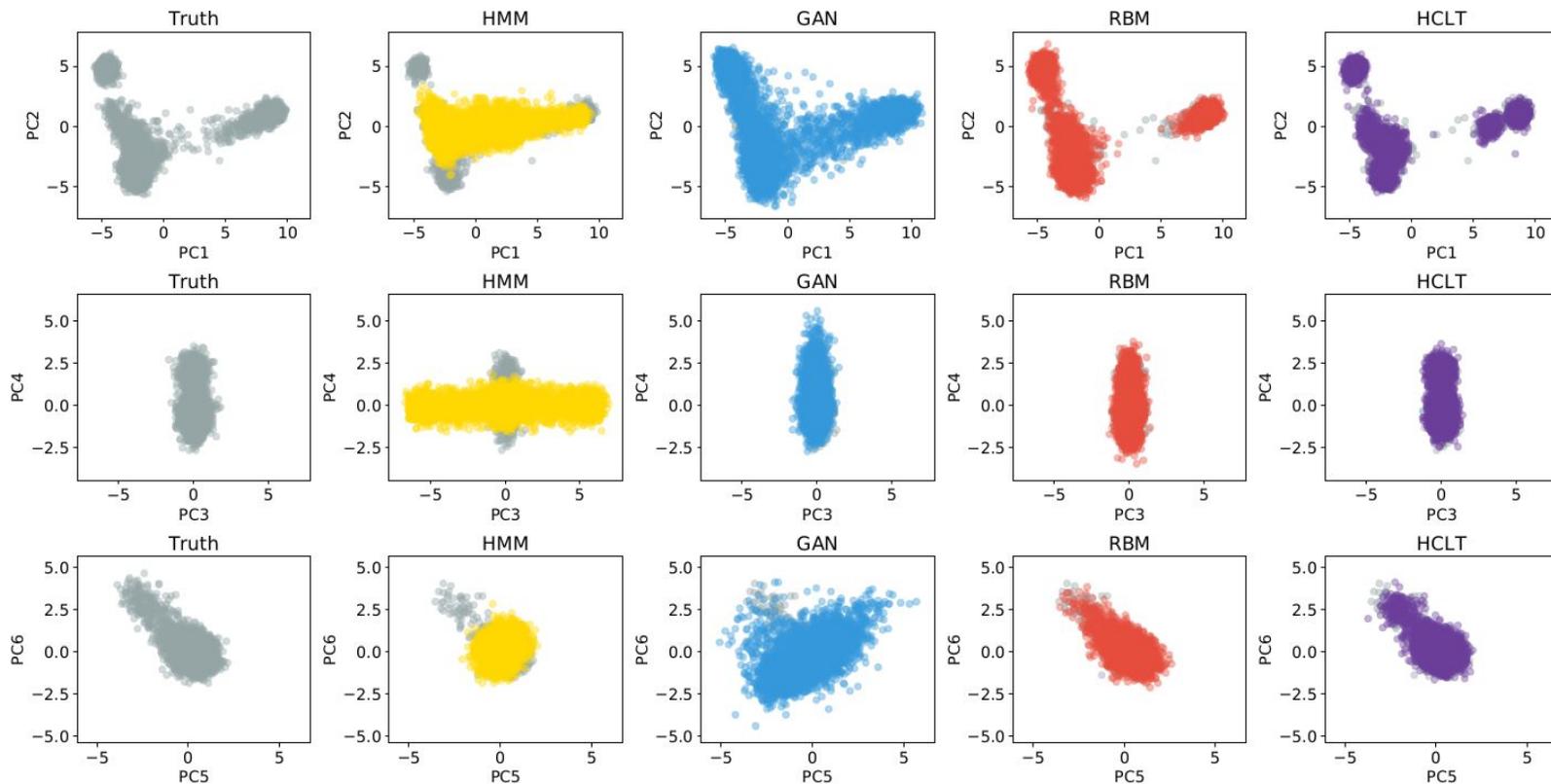
Method	# parameters	Theoretical bpd	Codeword bpd	Comp. time (s)	Decomp. time (s)
PC (HCLT, $M=16$)	3.3M	1.26	1.30	9	44
PC (HCLT, $M=24$)	5.1M	1.22	1.26	15	86
PC (HCLT, $M=32$)	7.0M	1.20	1.24	26	142
IDF	24.1M	1.90	1.96	288	592
BitSwap	2.8M	1.27	1.31	578	326

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

Tractable and expressive generative models of genetic variation data

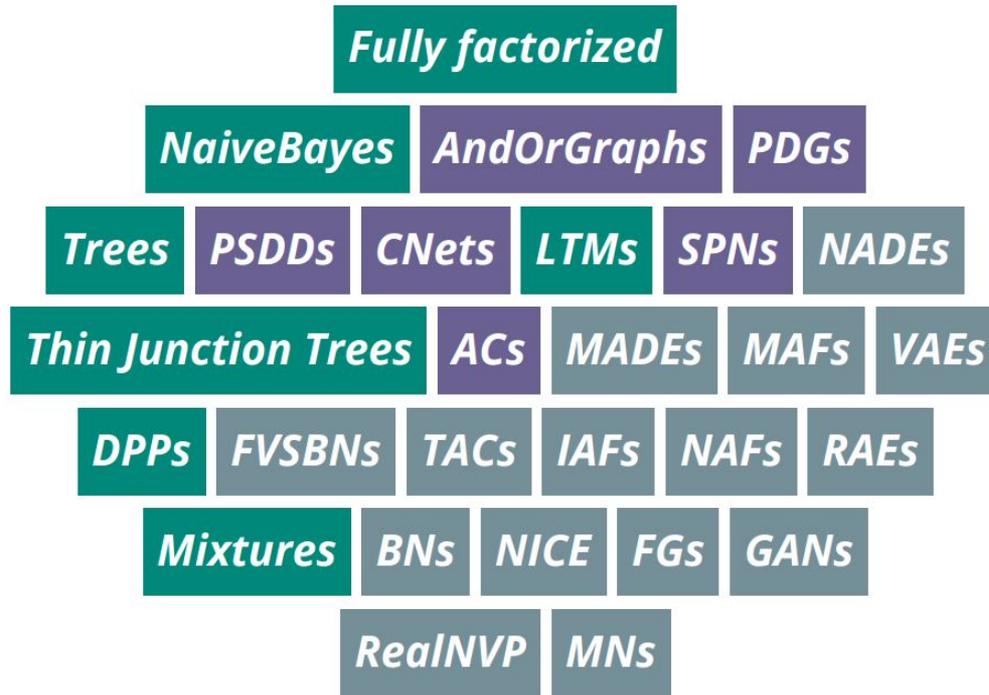


PC Learners keep getting better! ... stay tuned ...

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

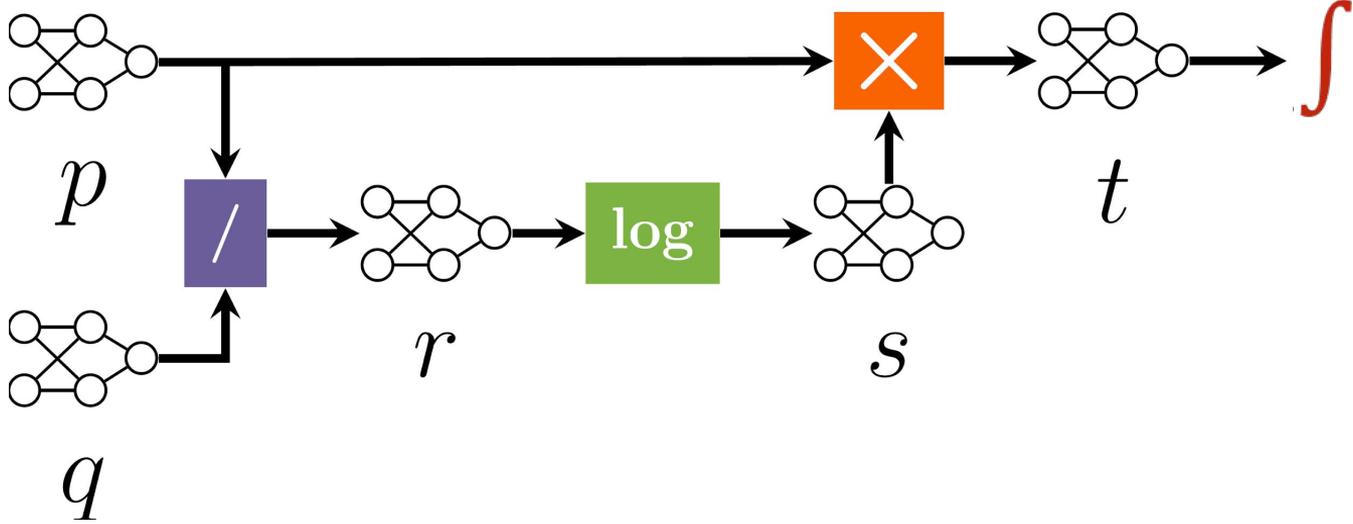
Dataset	PC	Bipartite flow	AF/SCF	IAF/SCF
Penn Treebank	1.23	1.38	1.46	1.63



Expressive* models without *compromises

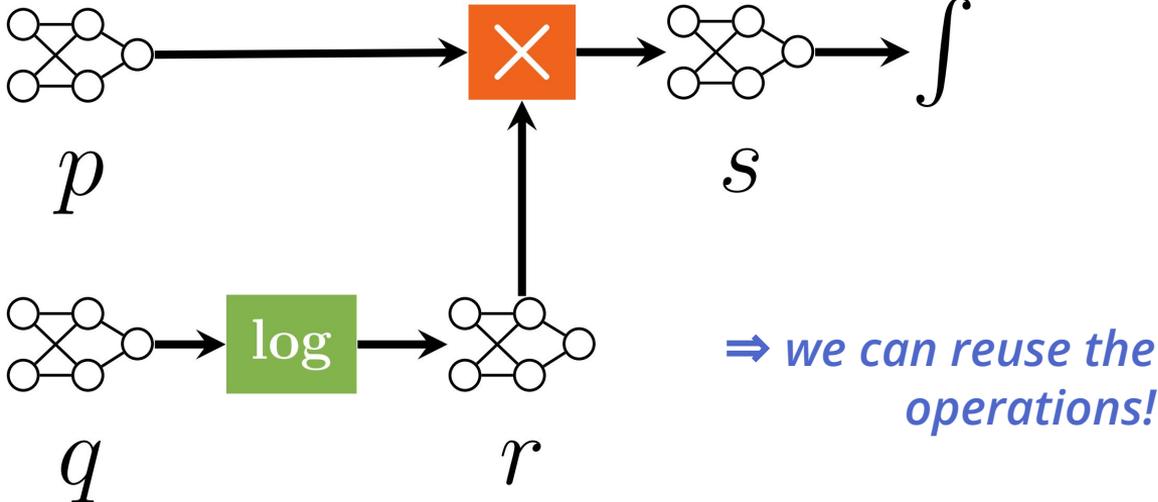
Queries as pipelines: KLD

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



Queries as pipelines: Cross Entropy

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



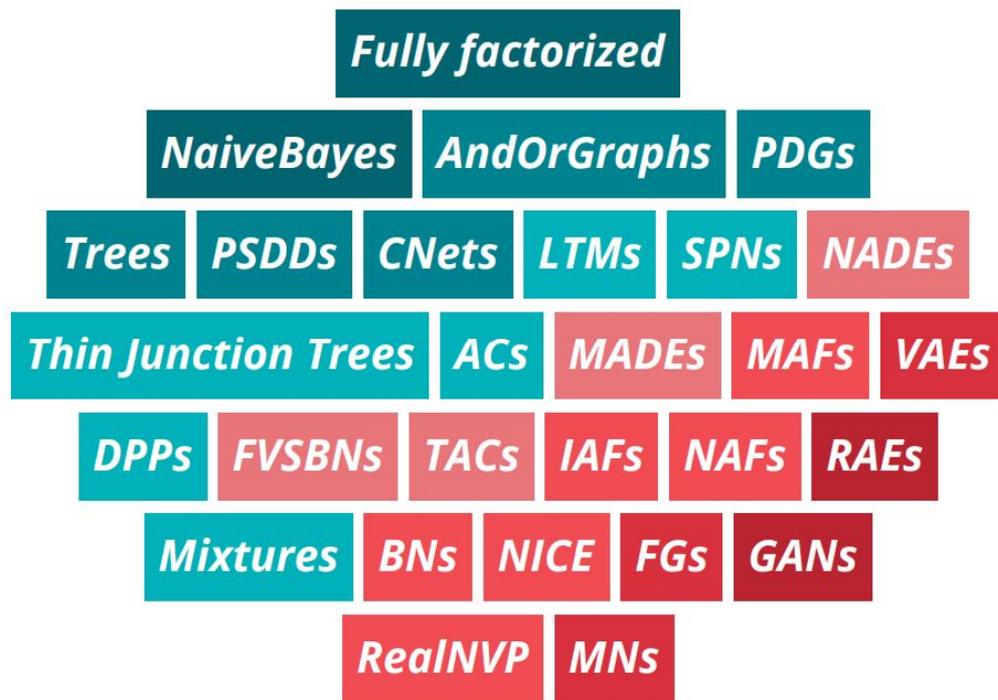
Tractable circuit operations

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

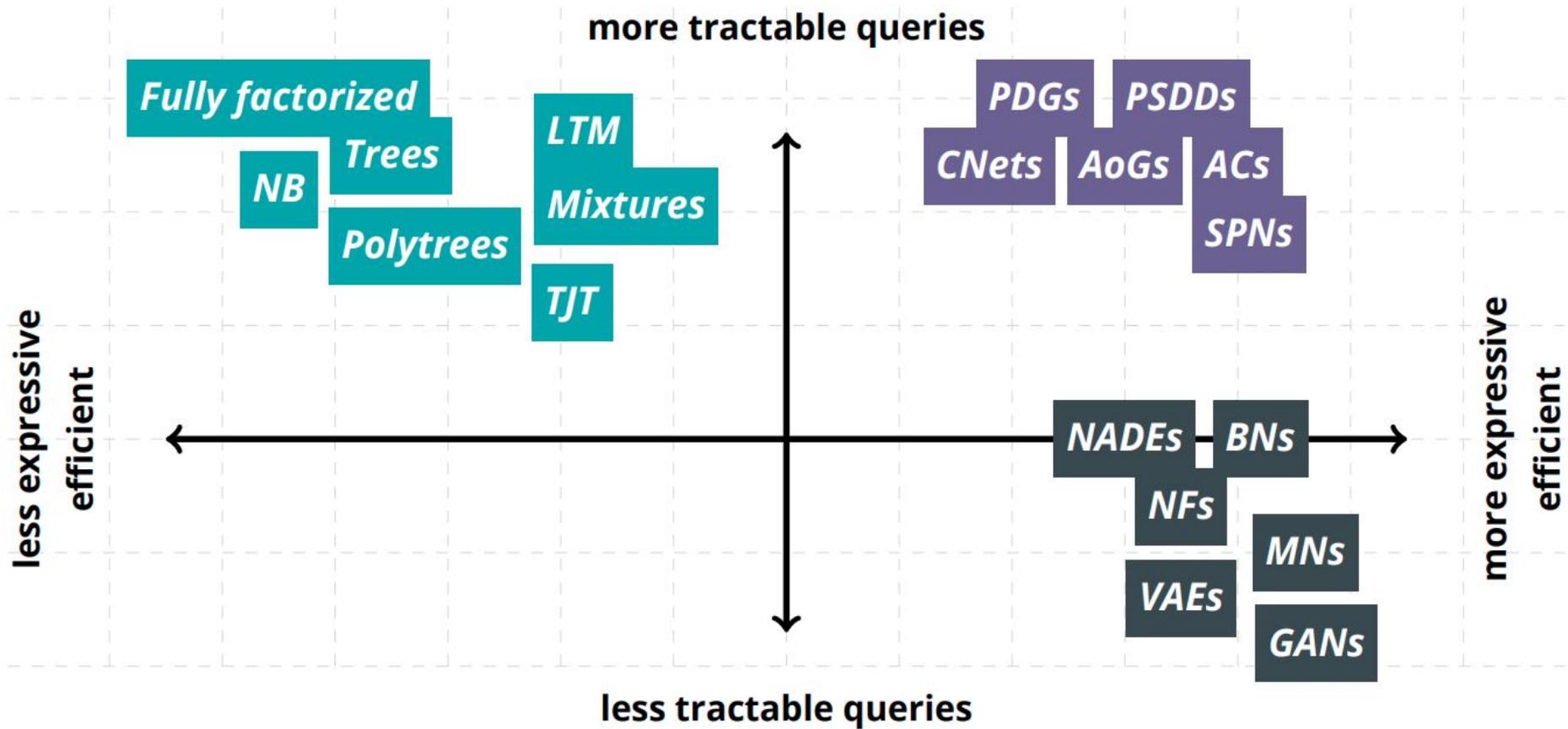
Inference by tractable operations

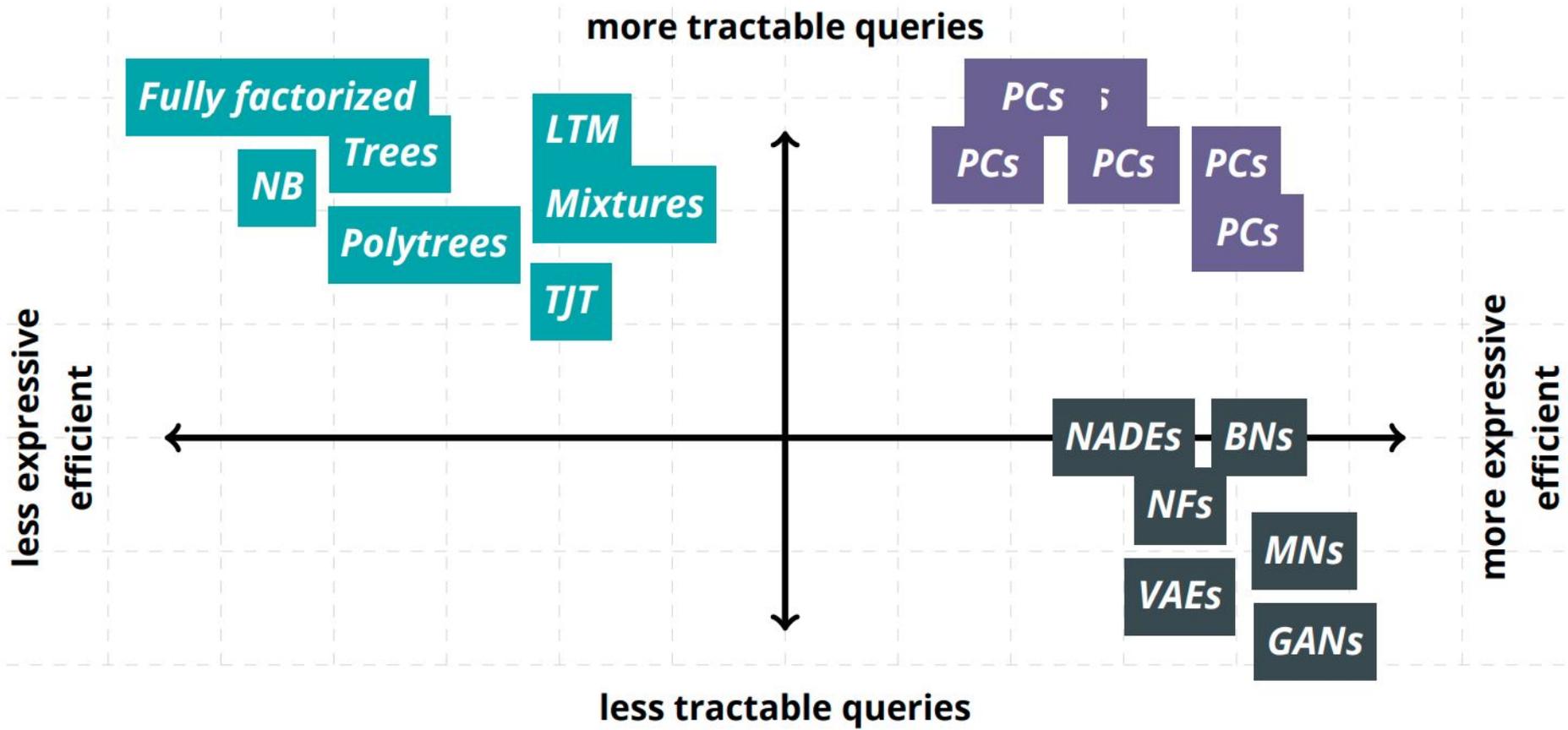
systematically derive tractable inference algorithm of complex queries

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



tractability is a spectrum





Learn more about probabilistic circuits?



Tutorial (3h)

Probabilistic Circuits

**Inference
Representations
Learning
Theory**

Antonio Vergari
University of California, Los Angeles

Robert Peharz
TU Eindhoven

YooJung Choi
University of California, Los Angeles

Guy Van den Broeck
University of California, Los Angeles

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

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

Overview Paper (80p)

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

YooJung Choi
Antonio Vergari
Guy Van den Broeck
*Computer Science Department
University of California
Los Angeles, CA, USA*

Contents

1	Introduction	3
2	Probabilistic Inference: Models, Queries, and Tractability	4
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 SotA likelihood full MNIST probabilistic circuit model in ~7 minutes on GPU:

https://github.com/Juice-jl/ProbabilisticCircuits.jl/blob/master/examples/train_mnist_hclt.ipynb

Public

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

master ProbabilisticCircuits.jl / examples / train_mnist_hclt.ipynb Go to file ...

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

1 contributor

609 Lines (609 sloc) | 26.5 KB

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

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

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

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

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

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

Outline

1. The paradox of learning to reason from data
deep learning
2. Tractable deep generative models
probabilistic reasoning + deep learning
3. **Learning with symbolic knowledge**
logical reasoning + deep learning

The AI Dilemma



Pure (Logic) Reasoning

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) Reasoning

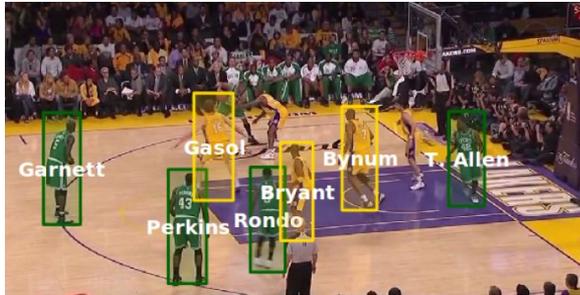
Pure Learning

- 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



Knowledge in Vision, Robotics, NLP, Activity Recognition

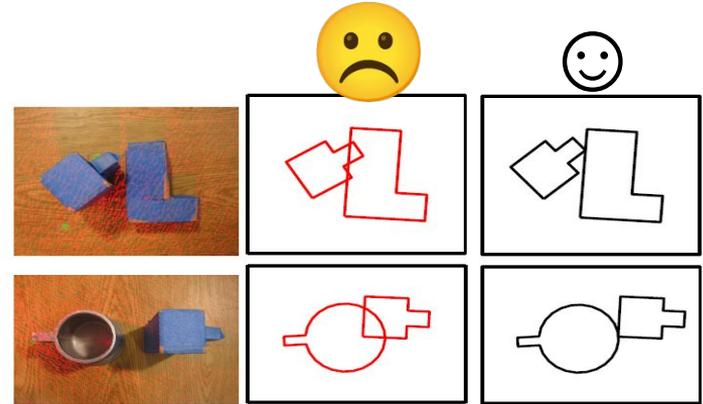


People appear at most once in a frame

A= At least one verb
in each sentence.
If X and Y are married,
then they are people.



Cut the orange before squeezing the orange

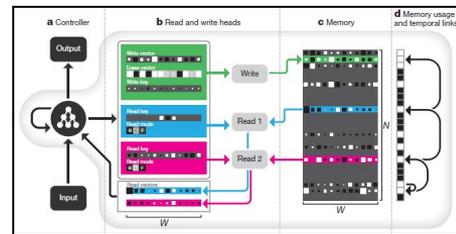


Rigid objects don't overlap

Motivation: Deep Learning

The image shows a screenshot of a New Scientist article. The header features the 'New Scientist' logo and navigation links for HOME, NEWS, TECHNOLOGY, SPACE, PHYSICS, HEALTH, EARTH, HUMANS, LIFE, TOPICS, EVENTS, and JOBS. The main headline reads 'Meet The People Shaping The Future Of Energy: Reinventing Energy Summit - 25 November in London'. Below this, there is a social media sharing bar with icons for Google+, Facebook, Twitter, and a plus sign, along with a '26' notification. The article title is 'DeepMind's AI has learned to navigate the Tube using memory'. The main image shows a person with long blonde hair looking at a large London Underground map.

The image shows a screenshot of a Nature article. The header features the 'nature' logo and the text 'International weekly journal of science'. Navigation links include Home, News & Comment, Research, Careers & Jobs, Current Issue, Archive, Audio & Video, and For. The article title is 'Google's AI reasons its way around the London Underground'. Below the title, there is a sub-headline: 'DeepMind's latest technique uses external memory to solve tasks that require logic and reasoning — a step toward more human-like AI.' There are also social media sharing icons for a share symbol and a printer icon.



[Graves, A., Wayne, G., Reynolds, M., Harley, T., Danihelka, I., Grabska-Barwińska, A., et al.. (2016). Hybrid computing using a neural network with dynamic external memory. *Nature*, 538(7626), 471-476.]

Motivation: Deep Learning

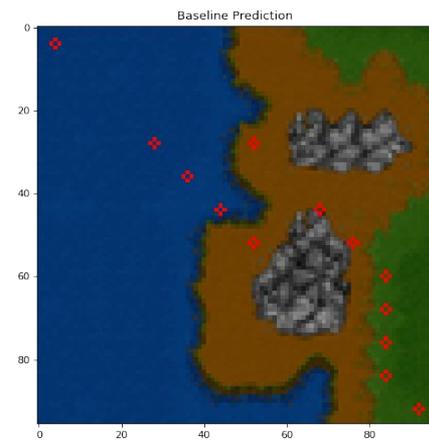
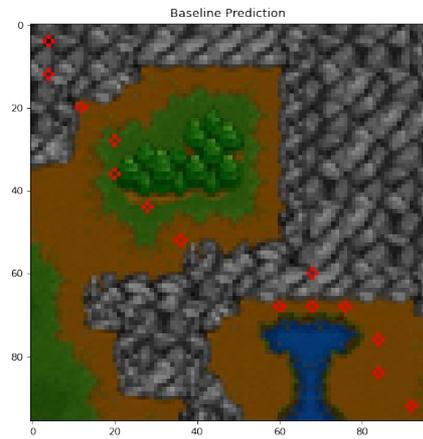
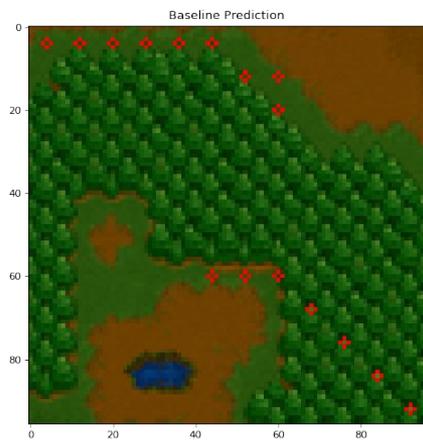
DeepMind's latest technique uses external memory to solve tasks that require **logic** and **reasoning** — a step toward more human-like AI.

... but ...

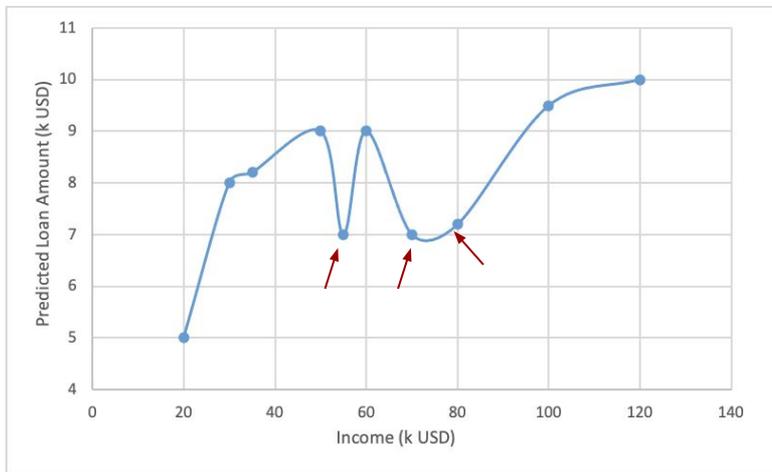


optimal planner recalculating a shortest path to the end node. To ensure that the network always moved to a valid node, the output distribution was renormalized over the set of possible triples outgoing from the current node. The performance

it also received input triples during the answer phase, indicating the actions chosen on the previous time-step. This makes the problem a 'structured prediction'



Predict Loan Amount



Neural Network Model: **Increasing income can decrease the approved loan amount**

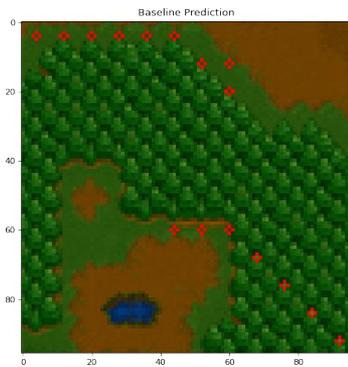
Monotonicity (Prior Knowledge):

Increasing income should increase the approved loan amount

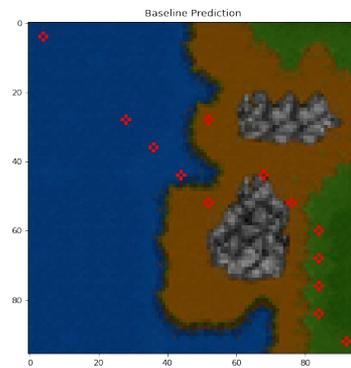
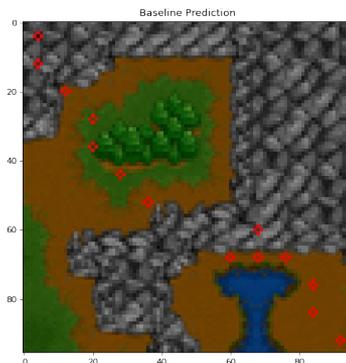
Knowledge vs. Data

- Where did the world knowledge go?
 - Python scripts
 - Decode/encode/search cleverly
 - Fix inconsistent beliefs
 - Rule-based decision systems
 - Dataset design
 - “a big hack” (with author’s permission)
- In some sense we went backwards
 - Less principled, scientific, and intellectually satisfying ways of incorporating knowledge

without constraint



without constraint



Warcraft min-cost simple-path prediction results

Test accuracy %	Coherent	Incoherent	Constraint
ResNet-18	44.8	97.7	56.9

*Is prediction
the shortest path?*
This is the real task!

*Are individual
edge predictions
correct?*

*Is output
a path?*

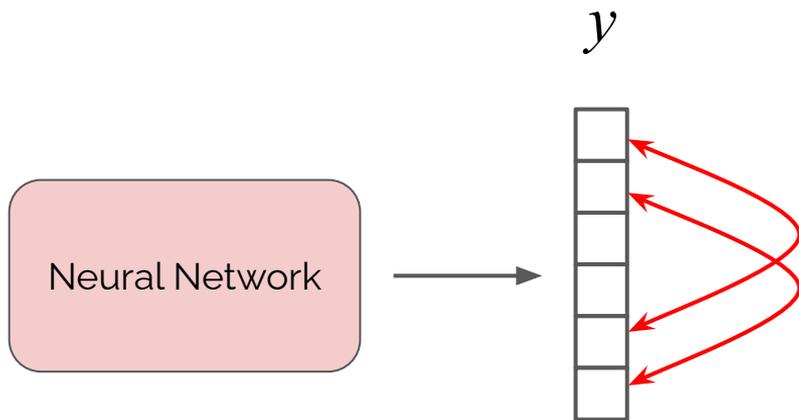
pylon

A PyTorch Framework for Learning with Constraints

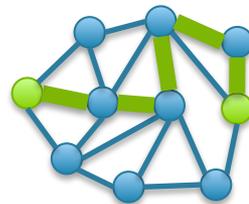
Kareem Ahmed Tao Li Thy Ton Quan Guo,
Kai-Wei Chang Parisa Kordjamshidi Vivek Srikumar
Guy Van den Broeck Sameer Singh

<http://pylon-lib.github.io>

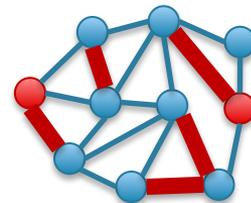
Declarative Knowledge of the Output



How is the output structured?
Are all possible outputs valid?



vs.



How are the outputs related to each other?

Learning this from data is inefficient
Much easier to express this declaratively

pylon

Library that extends PyTorch to allow injection of declarative knowledge

- **Easy to Express Knowledge:** users write **arbitrary constraints** on the output
- **Integrates with PyTorch:** **minimal change** to existing code
- **Efficient Training:** compiles into loss that can be **efficiently optimized**
 - Exact semantic loss (see later)
 - Monte-carlo estimate of loss
 - T-norm approximation
 - *your solver?*

pylon

PyTorch Code

```
for i in range(train_iters):  
    ...  
    py = model(x)  
    ...  
    loss = CrossEntropy(py, ...)
```

1

Specify knowledge as a predicate

```
def check(y):  
    ...  
    return isValid
```

pylon

PyTorch Code

```
for i in range(train_iters):  
    ...  
    py = model(x)  
    ...  
    loss = CrossEntropy(py, ...)  
    loss += constraint_loss(check)(py)
```

1

Specify knowledge as a predicate

```
def check(y):  
    ...  
    return isValid
```

2

Add as loss to training

```
loss += constraint_loss(check)
```

pylon

PyTorch Code

```
for i in range(train_iters):  
    ...  
    py = model(x)  
    ...  
    loss = CrossEntropy(py, ...)  
    loss += constraint_loss(check)(py)
```

1 Specify knowledge as a predicate

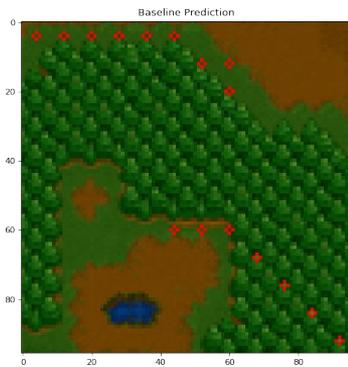
```
def check(y):  
    ...  
    return isValid
```

2 Add as loss to training

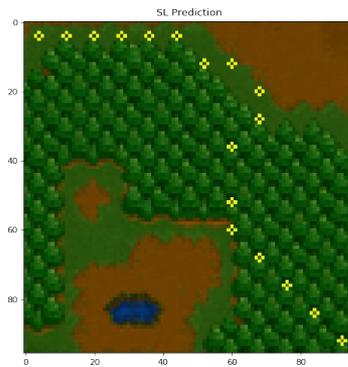
```
loss += constraint_loss(check)
```

3 pylon derives the gradients
(solves a combinatorial problem)

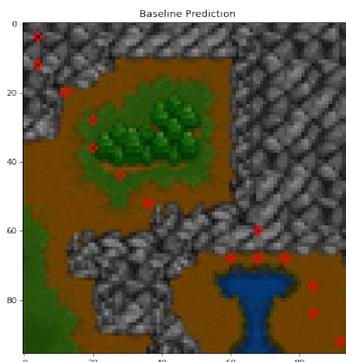
without constraint



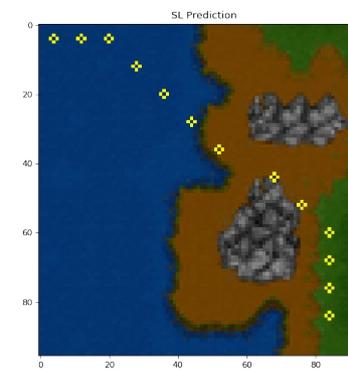
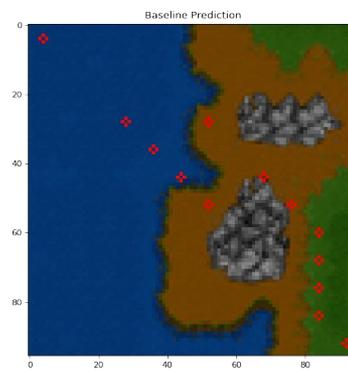
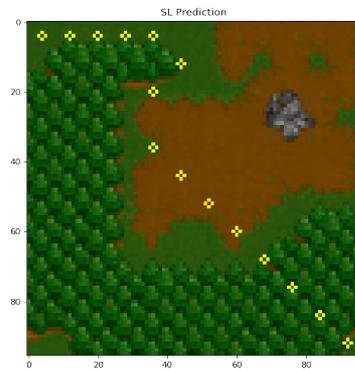
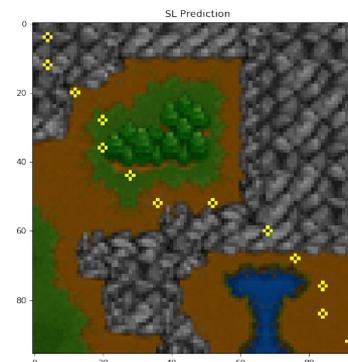
with constraint



without constraint



with constraint



Warcraft min-cost simple-path prediction results

Test accuracy %	Coherent	Incoherent	Constraint
ResNet-18	44.8	97.7	56.9
+ Semantic loss	50.9	97.7	67.4

Semantic Loss

Q: How close is output \mathbf{p} to satisfying constraint α ?

A: Semantic loss function $L(\alpha, \mathbf{p})$

- Axioms, for example:
 - If α constrains to one label, $L(\alpha, \mathbf{p})$ is cross-entropy
 - If α implies β then $L(\alpha, \mathbf{p}) \geq L(\beta, \mathbf{p})$ (α more strict)
- Implied Properties:
 - If α is equivalent to β then $L(\alpha, \mathbf{p}) = L(\beta, \mathbf{p})$
 - If \mathbf{p} is Boolean and satisfies α then $L(\alpha, \mathbf{p}) = 0$

 **SEMANTIC**
Loss!

Axioms imply unique semantic loss:

$$L^s(\alpha, \mathbf{p}) \propto -\log \underbrace{\sum_{\mathbf{x} \models \alpha} \prod_{i: \mathbf{x} \models X_i} p_i \prod_{i: \mathbf{x} \models \neg X_i} (1 - p_i)}_{\text{Probability of satisfying constraint } \alpha \text{ after sampling from neural net output layer } \mathbf{p}}$$

Probability of satisfying constraint α after sampling from neural net output layer \mathbf{p}

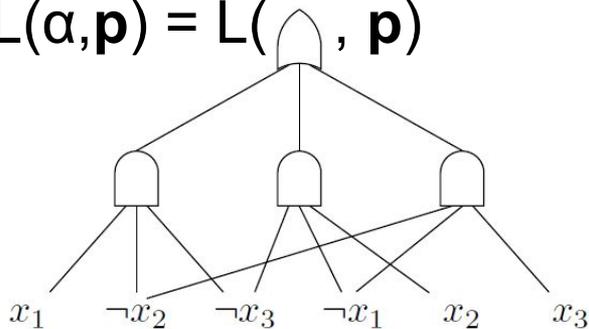
In general: #P-hard 😞

Do this probabilistic-logical reasoning during learning in a computation graph

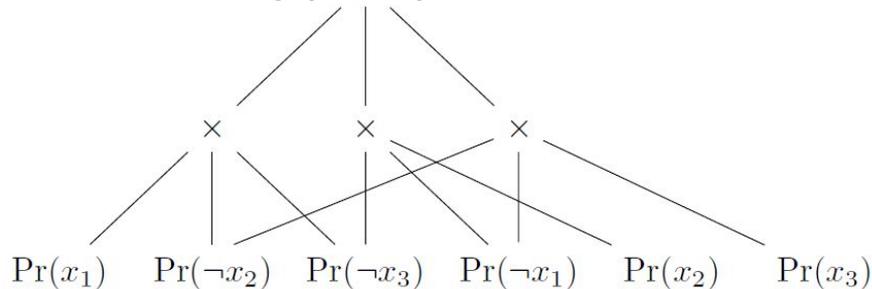
Circuits = Computation Graphs

- Logical circuits that can count solutions (#SAT)
also compute semantic loss efficiently in size of circuit

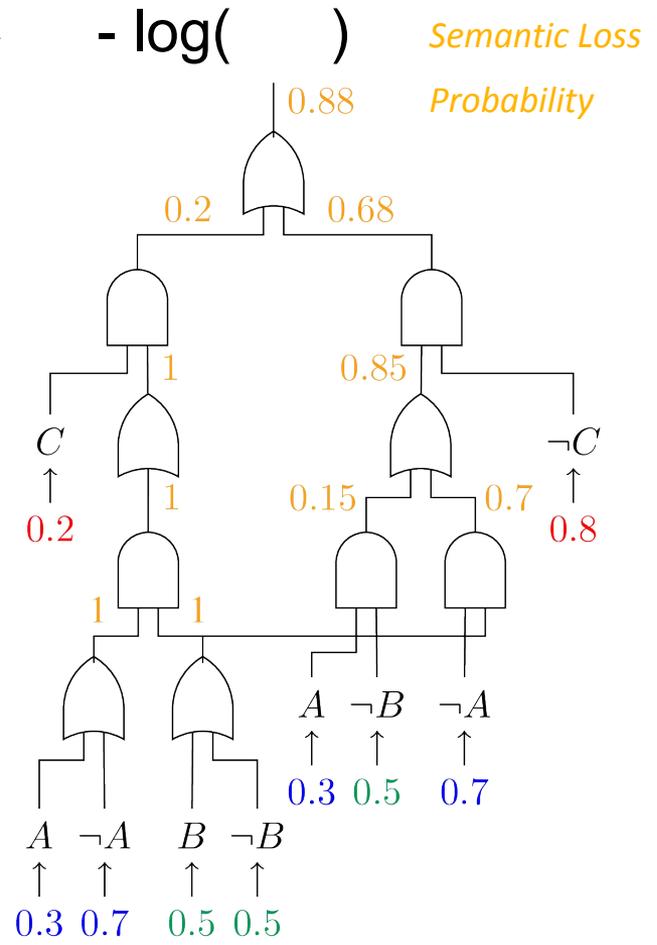
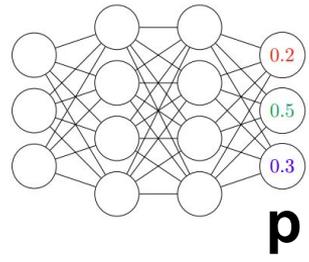
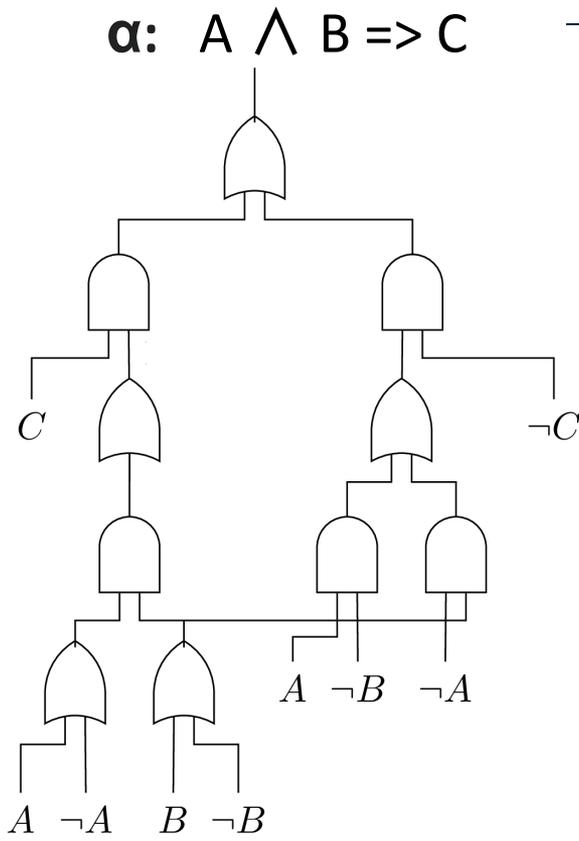
$$L(\alpha, \mathbf{p}) = L(\text{Circuit}, \mathbf{p})$$

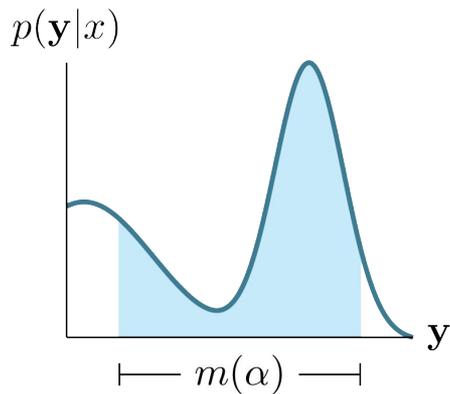


$$= -\log(\text{Circuit})$$



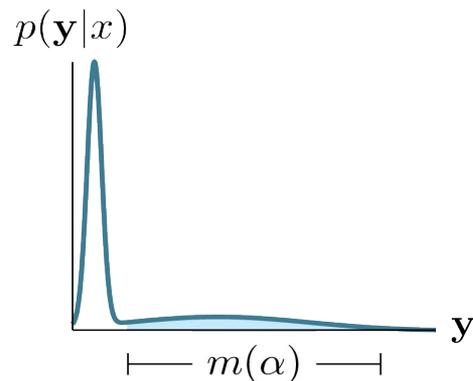
- Compilation into circuit by SAT solvers (once)
- Add circuit to neural network output in pytorch/tensorflow/...





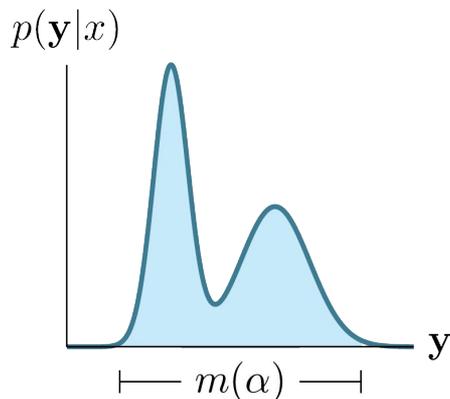
a) A network uncertain over both valid & invalid predictions

**Entropy
Regularization**



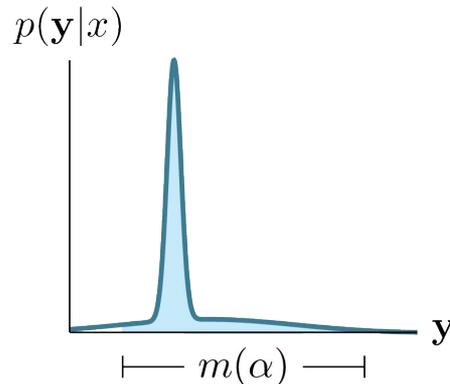
b) A network allocating most of its mass to an invalid prediction.

**Neuro-Symbolic
Learning**



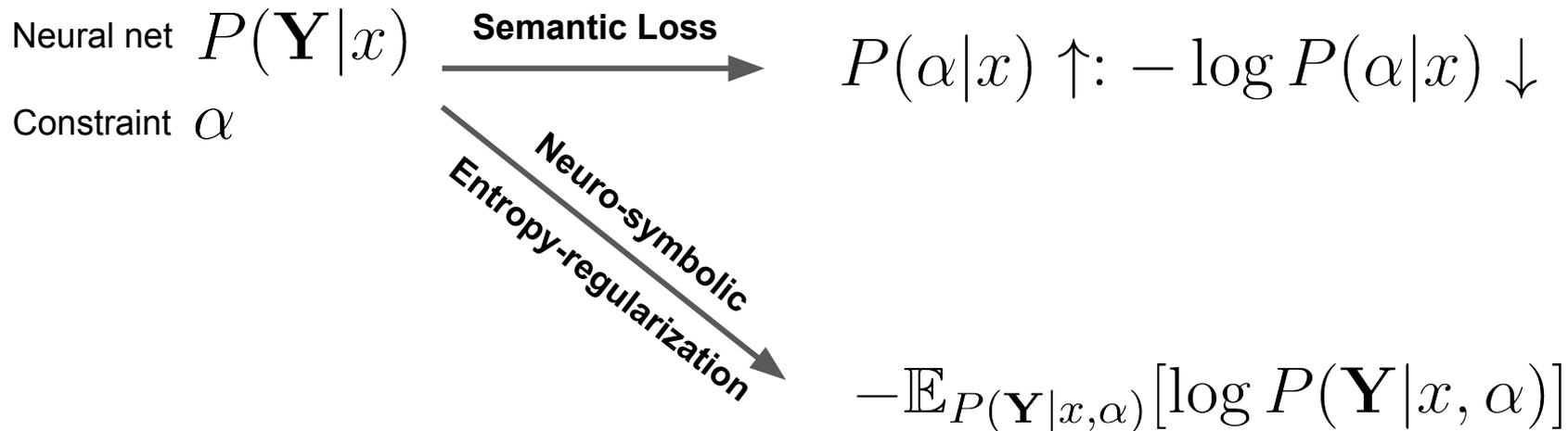
c) A network allocating most of its mass to models of constraint

**Neuro-Symbolic
Entropy Regularization**



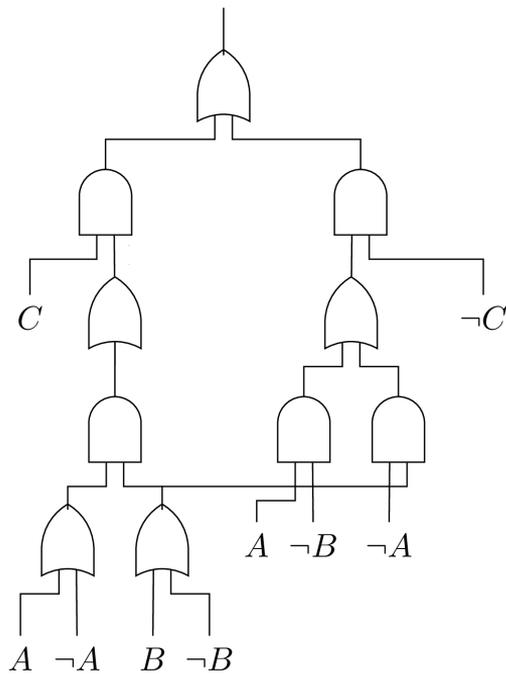
d) A network allocating most of mass to one model of formula

Two complementary neuro-symbolic losses

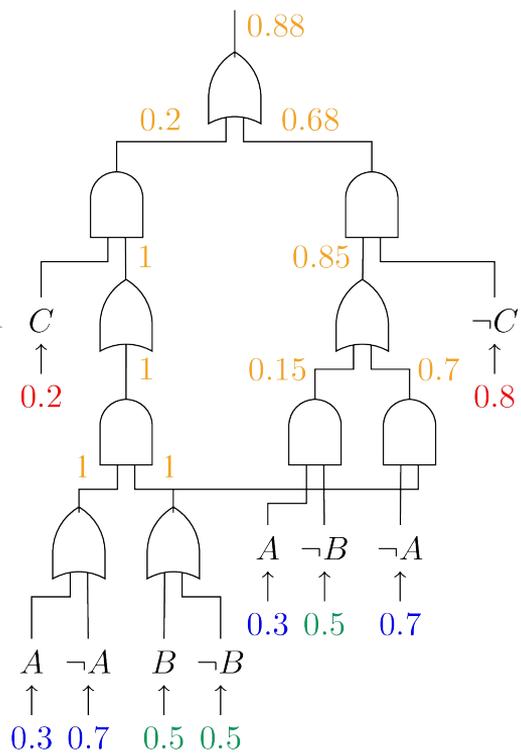


Warcraft min-cost simple-path prediction results

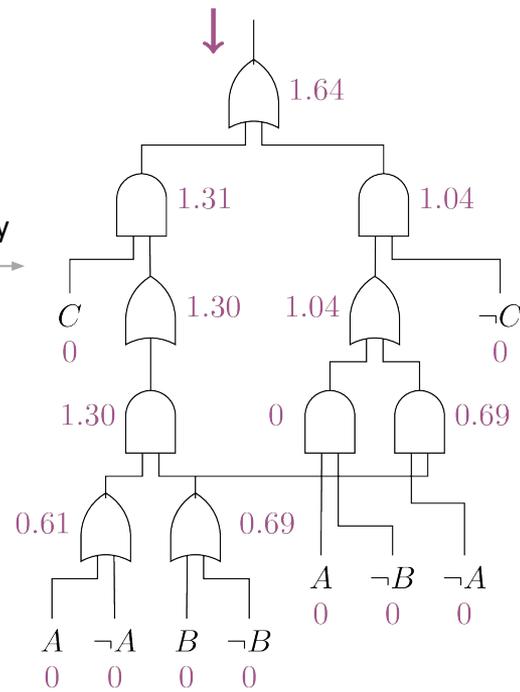
Test accuracy %	Coherent	Incoherent	Constraint
ResNet-18	44.8	97.7	56.9
Semantic loss	50.9	97.7	67.4
+ Full Entropy	51.5	97.6	67.7
+ NeSy Entropy	55.0	97.9	69.8



Probability



Entropy

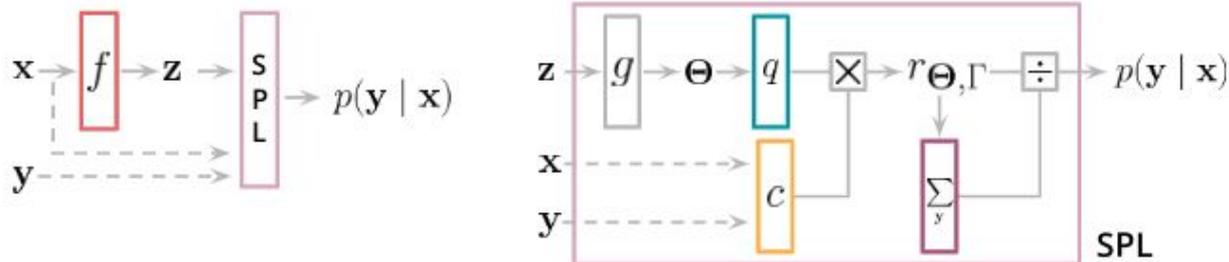


Joint entity-relation extraction in natural language processing

#		3	5	10	15	25	50	75
ACE05	Baseline	4.92 ± 1.12	7.24 ± 1.75	13.66 ± 0.18	15.07 ± 1.79	21.65 ± 3.41	28.96 ± 0.98	33.02 ± 1.17
	Self-training	7.72 ± 1.21	12.83 ± 2.97	16.22 ± 3.08	17.55 ± 1.41	27.00 ± 3.66	32.90 ± 1.71	37.15 ± 1.42
	Product t-norm	8.89 ± 5.09	14.52 ± 2.13	19.22 ± 5.81	21.80 ± 7.67	30.15 ± 1.01	34.12 ± 2.75	37.35 ± 2.53
	Semantic Loss	12.00 ± 3.81	14.92 ± 3.14	22.23 ± 3.64	27.35 ± 3.10	30.78 ± 0.68	36.76 ± 1.40	38.49 ± 1.74
	+ Full Entropy	14.80 ± 3.70	15.78 ± 1.90	23.34 ± 4.07	28.09 ± 1.46	31.13 ± 2.26	36.05 ± 1.00	39.39 ± 1.21
	+ NeSy Entropy	14.72 ± 1.57	18.38 ± 2.50	26.41 ± 0.49	31.17 ± 1.68	35.85 ± 0.75	37.62 ± 2.17	41.28 ± 0.46
SciERC	Baseline	2.71 ± 1.10	2.94 ± 1.00	3.49 ± 1.80	3.56 ± 1.10	8.83 ± 1.00	12.32 ± 3.00	12.49 ± 2.60
	Self-training	3.56 ± 1.40	3.04 ± 0.90	4.14 ± 2.60	3.73 ± 1.10	9.44 ± 3.80	14.82 ± 1.20	13.79 ± 3.90
	Product t-norm	6.50 ± 2.00	8.86 ± 1.20	10.92 ± 1.60	13.38 ± 0.70	13.83 ± 2.90	19.20 ± 1.70	19.54 ± 1.70
	Semantic Loss	6.47 ± 1.02	9.31 ± 0.76	11.50 ± 1.53	12.97 ± 2.86	14.07 ± 2.33	20.47 ± 2.50	23.72 ± 0.38
	+ Full Entropy	6.26 ± 1.21	8.49 ± 0.85	11.12 ± 1.22	14.10 ± 2.79	17.25 ± 2.75	22.42 ± 0.43	24.37 ± 1.62
	+ NeSy Entropy	6.19 ± 2.40	8.11 ± 3.66	13.17 ± 1.08	15.47 ± 2.19	17.45 ± 1.52	22.14 ± 1.46	25.11 ± 1.03

Semantic Probabilistic Layers

- How to give a 100% guarantee that Boolean constraints will be satisfied?
- Bake the constraint into the neural network as a special layer



- Secret sauce is again tractable circuits – computation graphs for reasoning

Hierarchical Multi-Label Classification

“if the image is classified as a dog, it must also be classified as an animal”

“if the image is classified as an animal, it must be classified as either cat or dog”

DATASET	EXACT MATCH	
	HMCNN	MLP+SPL
CELLCYCLE	3.05 ± 0.11	3.79 ± 0.18
DERISI	1.39 ± 0.47	2.28 ± 0.23
EISEN	5.40 ± 0.15	6.18 ± 0.33
EXPR	4.20 ± 0.21	5.54 ± 0.36
GASCH1	3.48 ± 0.96	4.65 ± 0.30
GASCH2	3.11 ± 0.08	3.95 ± 0.28
SEQ	5.24 ± 0.27	7.98 ± 0.28
SPO	1.97 ± 0.06	1.92 ± 0.11
DIATOMS	48.21 ± 0.57	58.71 ± 0.68
ENRON	5.97 ± 0.56	8.18 ± 0.68
IMCLEF07A	79.75 ± 0.38	86.08 ± 0.45
IMCLEF07D	76.47 ± 0.35	81.06 ± 0.68

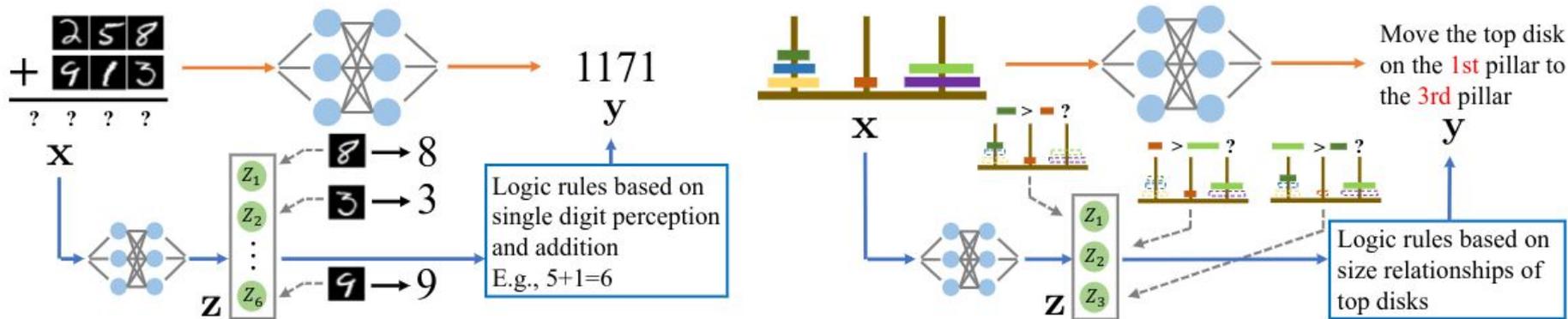
Neuro-Symbolic Learning Settings

Learn

1. **neural network** given **symbols and constraints and data**
2. **neural network and constraints** given **symbols and data**
3. **neural network and constraints and symbols** given **data**

Everyone is working on 1. Ongoing work on 2.

Neuro-Symbolic Joint Training



Learn invariant features using neural networks. Learn logic to tie it all together.

Ask Yitao Liang, Anji Liu

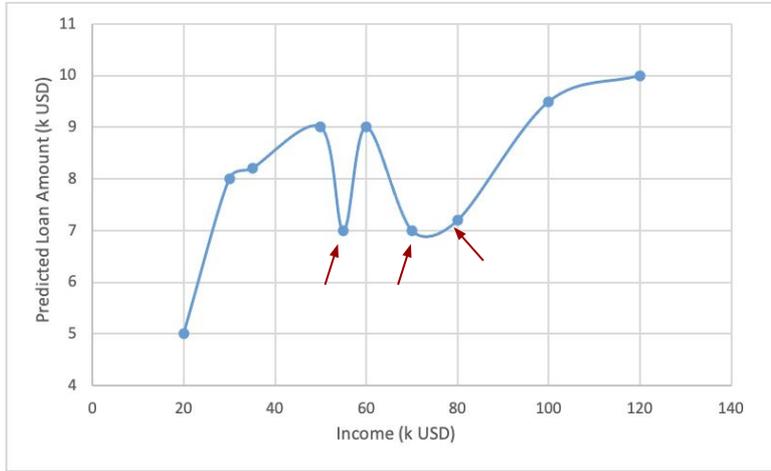
Neuro-Symbolic Joint Training

Model	Multi-digit addition [test seq length + train/test img]						Tower of Hanoi		
	5 w/ test	10 w/ test	20 w/ test	5 w/ train	10 w/ train	20 w/ train	Task #1	Task #2	Task #3
DeepProbLog [†]	88.30	77.46	timeout	94.92	89.74	timeout	89.28	97.96	89.33
LSTM	81.40	56.97	39.05	88.92	77.40	63.23	78.26	98.32	74.36
DNC	81.49	59.64	33.83	81.88	59.96	37.85	76.20	97.87	73.87
NToC(ours)	89.82	77.97	63.55	89.97	86.07	71.96	85.16	97.94	85.49

Learn invariant features using neural networks. Learn logic to tie it all together.

Ask Yitao Liang, Anji Liu

Predict Loan Amount

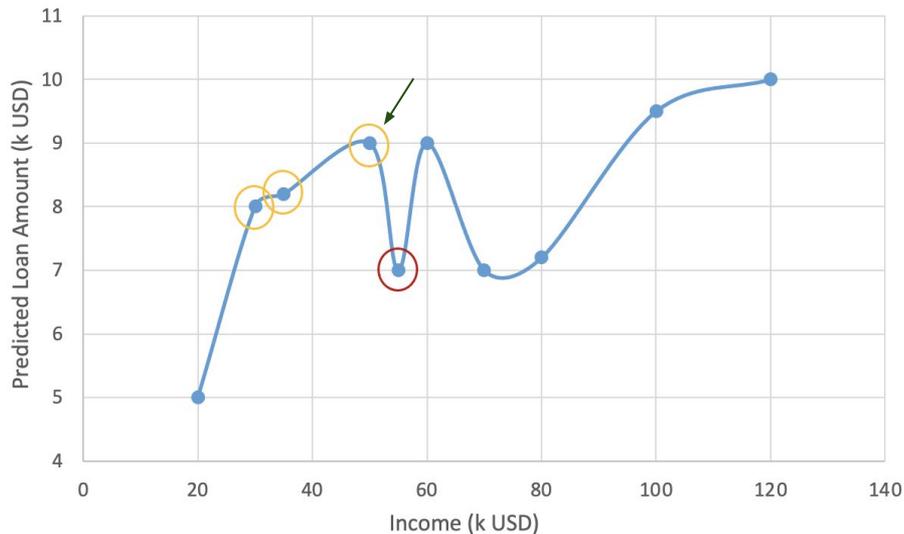


Neural Network Model: **Increasing income can decrease the approved loan amount**

Monotonicity (Prior Knowledge):

Increasing income should increase the approved loan amount

Counterexamples

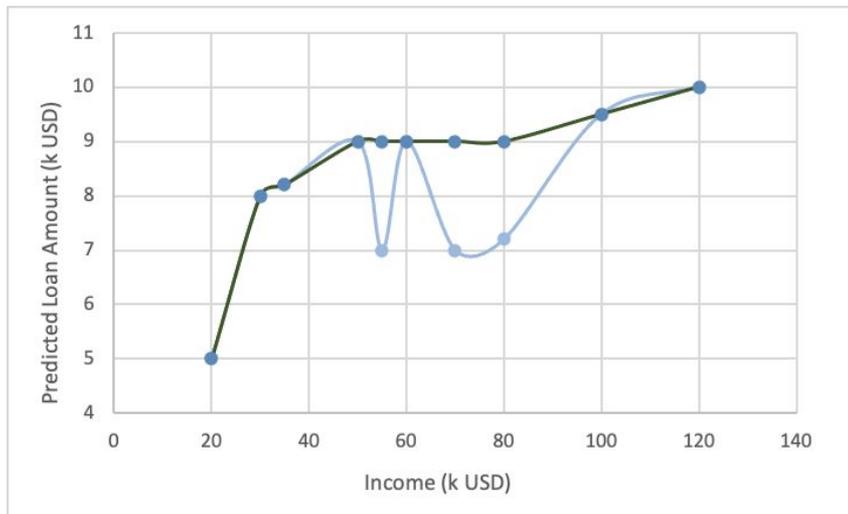


$$\exists x, y \ x \leq y \implies f(x) > f(y)$$

Computed using SMT(LRA)
logical reasoning solver

Maximal counterexamples
(largest violation) using OMT

Counterexample-Guided Predictions



Monotonic Envelope:

- Replace each prediction by its maximal counterexample
- Envelope construction is online (during prediction)
- Guarantees monotonic predictions for any ReLU neural net

- Works for high-dimensional input
- Works for multiple monotonic features

Monotonic Envelope: Performance

Dataset	Feature	NN _b	Envelope
Auto-MPG	Weight	9.33±3.22	9.19±3.41
	Displ.	9.33±3.22	9.63±2.61
	W,D	9.33±3.22	9.63±2.61
	W,D,HP	9.33±3.22	9.63±2.61
Boston	Rooms	14.37±2.4	14.19±2.28
	Crime	14.37±2.4	14.02±2.17

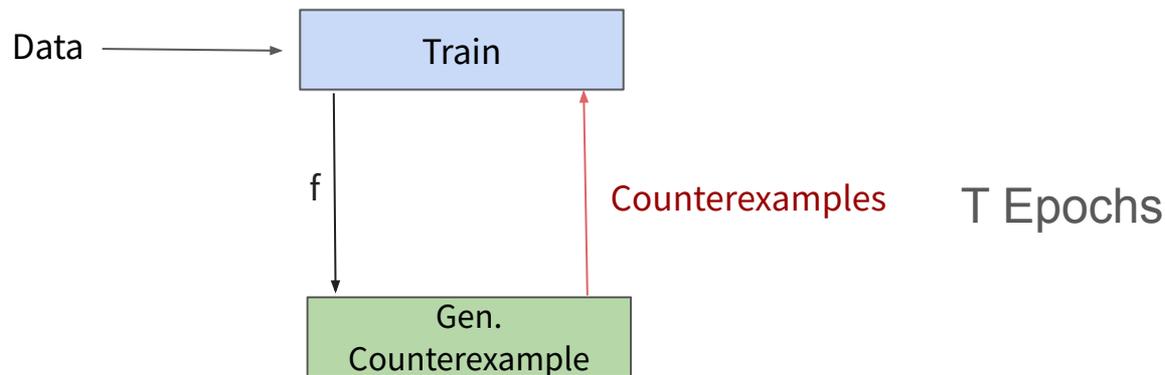
Dataset	Feature	NN _b	Envelope
Heart	Trestbps	0.85±0.04	0.85±0.04
	Chol.	0.85±0.04	0.85±0.05
	T,C	0.85±0.04	0.85±0.05
Adult	Cap. Gain	0.84	0.84
	Hours	0.84	0.84

Guaranteed monotonicity at little to no cost

Counterexample-Guided Learning

How to use monotonicity to improve model quality?

“Monotonicity as inductive bias”



Counterexample-Guided Learning: Performance

Dataset	Feature	NN _b	CGL
Auto-MPG	Weight	9.33±3.22	9.04±2.76
	Displ.	9.33±3.22	9.08±2.87
	W,D	9.33±3.22	8.86±2.67
	W,D,HP	9.33±3.22	8.63±2.21
Boston	Rooms	14.37±2.4	12.24±2.87
	Crime	14.37±2.4	11.66±2.89

Dataset	Feature	NN _b	CGL
Heart	Trestbps	0.85±0.04	0.86±0.02
	Chol.	0.85±0.04	0.85±0.05
	T,C	0.85±0.04	0.86±0.06
Adult	Cap. Gain	0.84	0.84
	Hours	0.84	0.84

Monotonicity is a *great* inductive bias for learning

COMET:

Counterexample-Guided Monotonicity Enforced Training

Table 4: Monotonicity is an effective inductive bias. COMET outperforms Min-Max networks on all datasets. COMET outperforms DLN in regression datasets and achieves similar results in classification datasets.

Dataset	Features	Min-Max	DLN	COMET
Auto-MPG	Weight	9.91 ± 1.20	16.77 ± 2.57	8.92 ± 2.93
	Displ.	11.78 ± 2.20	16.67 ± 2.25	9.11 ± 2.25
	W,D	11.60 ± 0.54	16.56 ± 2.27	8.89 ± 2.29
	W,D,HP	10.14 ± 1.54	13.34 ± 2.42	8.81 ± 1.81
Boston	Rooms	30.88 ± 13.78	15.93 ± 1.40	11.54 ± 2.55
	Crime	25.89 ± 2.47	12.06 ± 1.44	11.07 ± 2.99

Dataset	Features	Min-Max	DLN	COMET
Heart	Trestbps	0.75 ± 0.04	0.85 ± 0.02	0.86 ± 0.03
	Chol.	0.75 ± 0.04	0.85 ± 0.04	0.87 ± 0.03
	T,C	0.75 ± 0.04	0.86 ± 0.02	0.86 ± 0.03
Adult	Cap. Gain	0.77	0.84	0.84
	Hours	0.73	0.85	0.84

COMET = Provable Guarantees + SotA Results

The AI Dilemma



- Knowledge is (hidden) everywhere in ML
- A little bit of reasoning goes a long way!

Deep learning with structured output constraints
Learning monotonic neural networks

Outline

1. The paradox of learning to reason from data

~~deep learning~~

2. Tractable deep generative models

probabilistic reasoning + deep learning

3. Learning with symbolic knowledge

logical reasoning + deep learning

The AI Dilemma



Integrate reasoning into modern deep learning algorithms

Thanks

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

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