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

Guy Van den Broeck

UC Berkeley DREAM/CPAR Seminar - Nov 12 2023

# Outline

1. The paradox of learning to reason from data

   ~~*deep learning*~~

2. Architectures for learning and reasoning

   *logical reasoning + deep learning*

   a. Constrained generative AI
   b. Constrained structured prediction

# Outline

1. **The paradox of learning to reason from data**

   ~~*deep learning*~~

2. Architectures for learning and reasoning

   *logical reasoning + deep learning*

   a. Constrained generative AI
   b. Constrained structured prediction

# Can Language Models Perform Logical Reasoning?

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



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.

Q: How is Carol related to Justin ?

A: Carol is the grandmother of Justin

Reasoning Example from the CLUTRR dataset

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

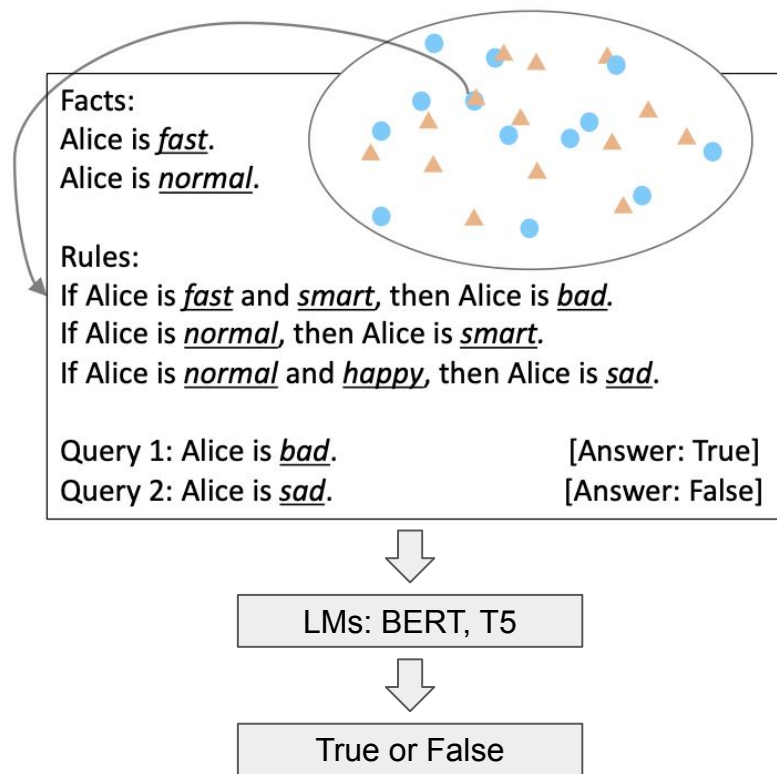Language Models:
*input → ? → Carol is the grandmother of Justin.*

Logical Reasoning:
*input → Justin in 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 (≈ 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

Facts:
Alice is _fast_.
Alice is _normal_.

Rules:
If Alice is _fast_ and _smart_, then Alice is _bad_.
If Alice is _normal_, then Alice is _smart_.
If Alice is _normal_ and _happy_, then Alice is _sad_.

Query 1: Alice is _bad_.            [Answer: True]
Query 2: Alice is _sad_.            [Answer: False]

LMs: BERT, T5

True or False

Honghua Zhang, Liunian Harold Li, Tao Meng, Kai-Wei Chang and Guy Van den Broeck. On the Paradox of Learning to Reason from Data, 2022

# SimpleLogic

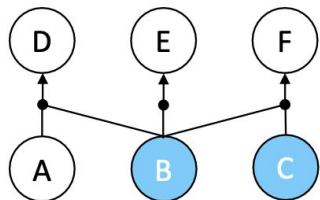Generate textual train and test examples of the form:

Rules: If witty, then diplomatic. If careless and condemned and attractive, then blushing. If dishonest and inquisitive and average, then shy. If average, then stormy. If popular, then blushing. If talented, then hurt. If popular and attractive, then thoughtless. If blushing and shy and stormy, then inquisitive. If adorable, then popular. If cooperative and wrong and stormy, then thoughtless. If popular, then sensible. If cooperative, then wrong. If shy and cooperative, then witty. If polite and shy and thoughtless, then talented. If polite, then condemned. If polite and wrong, then inquisitive. If dishonest and inquisitive, then talented. If blushing and dishonest, then careless. If inquisitive and dishonest, then troubled. If blushing and stormy, then shy. If diplomatic and talented, then careless. If wrong and beautiful, then popular. If ugly and shy and beautiful, then stormy. If shy and inquisitive and attractive, then diplomatic. If witty and beautiful and frightened, then adorable. If diplomatic and cooperative, then sensible. If thoughtless and inquisitive, then diplomatic. If careless and dishonest and troubled, then cooperative. If hurt and witty and troubled, then dishonest. If scared and diplomatic and troubled, then average. If ugly and wrong and careless, then average. If dishonest and scared, then polite. If talented, then dishonest. If condemned, then wrong. If wrong and troubled and blushing, then scared. If attractive and condemned, then frightened. If hurt and condemned and shy, then witty. If cooperative, then attractive. If careless, then polite. If adorable and wrong and careless, then diplomatic. Facts: Alice sensible Alice condemned Alice thoughtless Alice polite Alice scared Alice average
Query: Alice is shy ?

Honghua Zhang, Liunian Harold Li, Tao Meng, Kai-Wei Chang and Guy Van den Broeck. On the Paradox of Learning to Reason from Data, 2022

# Training a transformer on SimpleLogic

(1) Randomly sample facts & rules.
Facts: B, C
Rules: A, B → D. B → E. B, C → F.

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

*Rule-Priority*

- - - - - - - - - - - - - - - - - - - - - -

*Label-Priority*

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

(1) Randomly assign labels to predicates.
True: B, C, E, F.
False: A, D.

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 the transformer 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 transformer parameters that compute the ground-truth reasoning function:

   <u>Theorem 1:</u> *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, the transformer has learned the ground-truth reasoning function!**

Honghua Zhang, Liunian Harold Li, Tao Meng, Kai-Wei Chang and Guy Van den Broeck. On the Paradox of Learning to Reason from Data, 2022

# 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 the transformer **has learned** to reason,
    it should not exhibit such generalization failure.

2.  If the transformer **has not learned** to reason,
    it is baffling how it achieves near-perfect in-distribution test accuracy.

Honghua Zhang, Liunian Harold Li, Tao Meng, Kai-Wei Chang and Guy Van den Broeck. On the Paradox of Learning to Reason from Data, 2022

# Why? Statistical Features

Monotonicity of entailment:

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

⬇

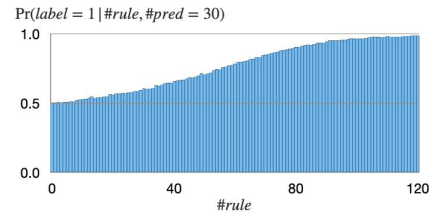The more rules given, the more likely a predicate will be proven.

⬇

Pr(label = True | Rule # = x) should increase (roughly) monotonically with x

Pr($label = 1 \,|\, \#rule$)

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

Pr($label = 1 \,|\, \#rule$)

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

Pr($label = 1 \,|\, \#rule, \#pred = 30$)

(c) Statistics for examples generated by uniform sampling;

# Model leverages statistical features to make predictions

RP_b downsamples from RP such that Pr(label = True | 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. Potentially countless statistical features

3. Such features are **inherent to the reasoning problem**, cannot make data "clean"

Honghua Zhang, Liunian Harold Li, Tao Meng, Kai-Wei Chang and Guy Van den Broeck. On the Paradox of Learning to Reason from Data, 2022

# First Conclusion

Experiments unveil the fundamental difference between

1.  learning to reason, and

2.  learning to achieve high performance on benchmarks using statistical features.

**Be careful deploying AI in applications where this difference matters.**

*FAQ: Do bigger transformers solve this problem? No, already 99% accurate…*

*FAQ: Will reasoning emerge? Perhaps on 99% of human behavior…*

# Outline

1. The paradox of learning to reason from data

   ~~*deep learning*~~

2. **Architectures for learning and reasoning**

   *logical reasoning + deep learning*

   a. Constrained generative AI
   b. Constrained structured prediction

# Outline

1. The paradox of learning to reason from data

   *~~deep learning~~*

2. **Architectures for learning and reasoning**

   *logical reasoning + deep learning*

   a. **Constrained generative AI**
   b. Constrained structured prediction

# Generative models are still hard to control



*Stable Diffusion*

Generate a sentence using "frisbee", "caught" and "dog", following the given order.

Generate a sentence using "frisbee", "caught" and "dog", following the given order.

After a perfect throw, the <u>frisbee</u> glided through the air, and the **<u>dog</u>**, with incredible agility, **<u>caught</u>** it mid-flight.
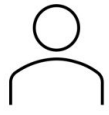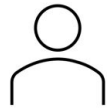
*ChatGPT*

**Generate a sentence using "frisbee", "caught" and "dog", following the given order.**

After a perfect throw, the <u>frisbee</u> glided through the air, and the <u>dog</u>, with incredible agility, <u>caught</u> it mid-flight.

*ChatGPT*

That's not correct. Generate a sentence using "frisbee", "caught" and "dog". The keywords should appear in the order as specified.

Here's the correct sentence: The <u>dog</u> <u>caught</u> the <u>frisbee</u> in mid-air, showing off its amazing catching skills.

*ChatGPT*

# *What do we have?*

Prefix: "The weather is"

Constraint α: text contains "winter"

Model only does $p(\text{next-token}|\text{prefix}) = $

| cold | 0.05 |
|------|------|
| warm | 0.10 |

Train some $q(.\,|\alpha)$ for a specific task distribution $\alpha \sim p_{\text{task}}$

*(amortized inference, encoder, masked model, seq2seq, prompt tuning,...)*

Train $q(\text{next-token}|\text{prefix}, \alpha)$

# *What do we need?*

Prefix: "The weather is"

Constraint α: text contains "winter"

Generate from $p(\text{next-token}|\text{prefix}, \alpha) =$

| | |
|---|---|
| cold | 0.50 |
| warm | 0.01 |

$$\propto \sum_{\text{text}} p(\text{next-token}, \text{text}, \text{prefix}, \alpha)$$

## *Marginalization!*

# CommonGen: a Challenging Benchmark

Given 3-5 concepts (keywords), our goal is to generate a sentence using all keywords, which can appear in any order and any form of inflections. e.g.,

Input: snow drive car

Reference 1: A car drives down a snow covered road.

Reference 2: Two cars drove through the snow.

$$(w_{1,1} \lor \ldots \lor w_{1,d1}) \land \ldots \land (w_{m,1} \lor \ldots \lor w_{m,dm})$$

Each clause represents the inflections for one keyword.

# Tractable Probabilistic Models

Tractable Probabilistic Models (TPMs) model joint probability distributions (just like auto-regressive LMs) and allow efficient computation of various probabilistic queries.

Probabilistic (Generating) Circuits

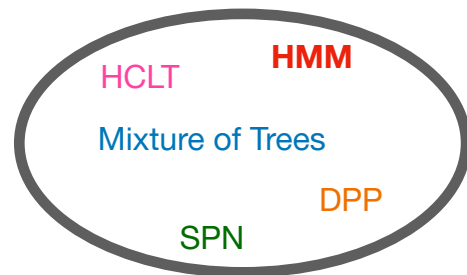HCLT

**HMM**

Mixture of Trees

DPP

SPN

e.g., efficient marginalization:

$$p_{TPM}(\text{3rd token = frisbee, 5th token = dog})$$

in particular …

$$\sum_{\text{sentence}} p_{TPM}(\textbf{sentence}, \text{next-token = "warm", prefix = "The weather is"}, \alpha)$$

➡ Efficient conditioning given lexical constraints : $p_{TPM}(\text{next-token} \mid \text{prefix}, \alpha)$

Honghua Zhang, Meihua Dang, Nanyun Peng and Guy Van den Broeck. Tractable Control for Autoregressive Language Generation, 2023.

# Step 1: Distill an HMM $p_{hmm}$ that approximates $p_{gpt}$



1. An HMM with 4096 hidden states and 50k emission tokens

2. Train the HMM on data sampled from GPT2-large (domain-adapted, either via prompting or fine-tuning), effectively minimizing $KL(p_{gpt} \parallel p_{HMM})$

3. Leverages the <u>latent variable distillation</u> technique for training probabilistic circuits at scale [ICLR 23]. (Cluster embeddings of examples to estimate latent $Z_i$)

Anji Liu, Honghua Zhang and Guy Van den Broeck. Scaling Up Probabilistic Circuits by Latent Variable Distillation, 2023.

# Computing $p_{hmm}(\alpha \mid x_{1:t+1})$

For **α** in conjunctive normal form (CNF):

$$(w_{1,1} \lor \ldots \lor w_{1,d1}) \land \ldots \land (w_{m,1} \lor \ldots \lor w_{m,dm})$$

where each $w_{ij}$ is a keyword (i.e. a **string of tokens**),
representing the constraint that $w_{ij}$ appears in the generated text.

e.g.,  α = ("swims" ∨ "like swimming") ∧ ("lake" ∨ "pool")

<u>Efficient algorithm</u>:
For m clauses and sequence length n, time-complexity for generation is $O(2^{|m|}n)$.

<u>Trick</u>: dynamic programming with clever preprocessing and local belief updates

Honghua Zhang, Meihua Dang, Nanyun Peng and Guy Van den Broeck. <u>Tractable Control for Autoregressive Language Generation</u>, 2023.

# GeLaTo Overview



**Lexical Constraint** $\alpha$: sentence contains keyword "winter"

**Constrained Generation**: $\mathrm{Pr}(x_{t+1} \,|\, \alpha, x_{1:t} = \text{"the weather is"})$

**X** intractable

✓ efficient

Pre-trained Language Model

Tractable Probabilistic Model

Minimize KL-divergence

| $x_{t+1}$ | $\mathrm{Pr}_{LM}(x_{t+1} \,|\, x_{1:t})$ |
|---|---|
| cold | 0.05 |
| warm | 0.10 |

| $x_{t+1}$ | $\mathrm{Pr}_{TPM}(\alpha \,|\, x_{t+1}, x_{1:t})$ |
|---|---|
| cold | 0.50 |
| warm | 0.01 |

Honghua Zhang, Meihua Dang, Nanyun Peng and Guy Van den Broeck. Tractable Control for Autoregressive Language Generation, 2023.

# GeLaTo Overview

**Lexical Constraint** $\alpha$: sentence contains keyword "winter"

**Constrained Generation**: $\mathrm{Pr}(x_{t+1} \,|\, \alpha, x_{1:t} = $ "the weather is")

✗ **intractable**                    ✓ **efficient**

Pre-trained Language Model

Tractable Probabilistic Model

*Minimize KL-divergence*

| $x_{t+1}$ | $\mathrm{Pr}_{LM}(x_{t+1} \,|\, x_{1:t})$ |
|---|---|
| cold | 0.05 |
| warm | 0.10 |

| $x_{t+1}$ | $\mathrm{Pr}_{TPM}(\alpha \,|\, x_{t+1}, x_{1:t})$ |
|---|---|
| cold | 0.50 |
| warm | 0.01 |

| $x_{t+1}$ | $p(x_{t+1} \,|\, \alpha, x_{1:t})$ |
|---|---|
| cold | 0.025 |
| warm | 0.001 |

Honghua Zhang, Meihua Dang, Nanyun Peng and Guy Van den Broeck. Tractable Control for Autoregressive Language Generation, 2023.

# Step 2: Control $p_{gpt}$ via $p_{hmm}$

## Unsupervised

Language model is not
fine-tuned/prompted to satisfy constraints

By Bayes rule:
$$p_{gpt}(x_{t+1} \mid x_{1:t}, \alpha) \propto p_{gpt}(\alpha \mid x_{1:t+1}) \cdot p_{gpt}(x_{t+1} \mid x_{1:t})$$

Assume $p_{hmm}(\alpha \mid x_{1:t+1}) \approx p_{gpt}(\alpha \mid x_{1:t+1})$, we
generate from:

$$p(x_{t+1} \mid x_{1:t}, \alpha) \propto p_{hmm}(\alpha \mid x_{1:t+1}) \cdot p_{gpt}(x_{t+1} \mid x_{1:t})$$

| Method | Generation Quality | | | | | | | | Constraint Satisfaction | | | |
| | ROUGE-L | | BLEU-4 | | CIDEr | | SPICE | | Coverage | | Success Rate | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| *Unsupervised* | *dev* | *test* | *dev* | *test* | *dev* | *test* | *dev* | *test* | *dev* | *test* | *dev* | *test* |
| InsNet (Lu et al., 2022a) | - | - | 18.7 | - | - | - | - | - | **100.0** | - | **100.0** | - |
| NeuroLogic (Lu et al., 2021) | - | 41.9 | - | 24.7 | - | 14.4 | - | 27.5 | - | 96.7 | - | - |
| A*esque (Lu et al., 2022b) | - | **44.3** | - | 28.6 | - | 15.6 | - | 29.6 | - | 97.1 | - | - |
| NADO (Meng et al., 2022) | - | - | 26.2 | - | - | - | - | - | 96.1 | - | - | - |
| GeLaTo | **44.6** | 44.1 | **29.9** | **29.4** | **16.0** | 15.8 | **31.3** | **31.0** | **100.0** | **100.0** | **100.0** | **100.0** |

Honghua Zhang, Meihua Dang, Nanyun Peng and Guy Van den Broeck. Tractable Control for Autoregressive Language Generation, 2023.

# Step 2: Control $p_{gpt}$ via $p_{hmm}$

## *Supervised*

Language model is fine-tuned to perform constrained generation (e.g. seq2seq)

Empirically $p_{HMM}(\alpha \mid x_{1:t+1}) \approx p_{gpt}(\alpha \mid x_{1:t+1})$ does not hold well enough;

we view $p_{HMM}(x_{t+1} \mid x_{1:t}, \alpha)$ and $p_{gpt}(x_{t+1} \mid x_{1:t})$ as classifiers trained for the same task with different biases; thus we generate from their *weighted geometric mean*:

$$p(x_{t+1} \mid x_{1:t}, \alpha) \propto p_{hmm}(x_{t+1} \mid x_{1:t}, \alpha)^w \cdot p_{gpt}(x_{t+1} \mid x_{1:t})^{1-w}$$

| Method | *Generation Quality* | | | | | | | | *Constraint Satisfaction* | | | |
| | ROUGE-L | | BLEU-4 | | CIDEr | | SPICE | | Coverage | | Success Rate | |
| *Supervised* | *dev* | *test* | *dev* | *test* | *dev* | *test* | *dev* | *test* | *dev* | *test* | *dev* | *test* |
| NeuroLogic (Lu et al., 2021) | - | 42.8 | - | 26.7 | - | 14.7 | - | 30.5 | - | 97.7 | - | 93.9[†] |
| A*esque (Lu et al., 2022b) | - | 43.6 | - | 28.2 | - | 15.2 | - | 30.8 | - | 97.8 | - | 97.9[†] |
| NADO (Meng et al., 2022) | 44.4[†] | - | 30.8 | - | 16.1[†] | - | **32.0**[†] | - | 97.1 | - | 88.8[†] | - |
| GeLaTo | **46.0** | **45.6** | **34.1** | **32.9** | **16.7** | **16.8** | 31.3 | **31.9** | **100.0** | **100.0** | **100.0** | **100.0** |

Honghua Zhang, Meihua Dang, Nanyun Peng and Guy Van den Broeck. Tractable Control for Autoregressive Language Generation, 2023.

# Advantages of our framework:

1. Constraint **α** is <u>guaranteed to be satisfied</u>: for any next-token $x_{t+1}$ that would make **α** unsatisfiable, $p(x_{t+1} \mid x_{1:t}, \alpha) = 0$ for both settings.

2. Training $p_{hmm}$ <u>does not depend on **α**</u>, which is only imposed at inference (generation) time. Once $p_{hmm}$ is trained, we can impose whatever **α**.

3. We can impose <u>additional tractable constraints</u>:
   - The keywords are generated following a particular order.
   - (Some) keywords must appear at a particular position.
   - (Some) keywords must **not** appear in the generated sentence.

Conclusion: you can control an intractable generative model using a tractable generative model for (symbolic) reasoning.

# Inpainting/constrained generation is still challenging



Diffusion models are good at fine-grained details, but not so good at global consistency of generated images.

# Inpainting/constrained generation is still challenging



**Tiramisu**

# Constrained posterior in diffusion models

Unconstrained denoising step: $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) := \sum_{\tilde{\mathbf{x}}_0} q(\mathbf{x}_{t-1}|\tilde{\mathbf{x}}_0, \mathbf{x}_t) \cdot p_\theta(\tilde{\mathbf{x}}_0|\mathbf{x}_t)$



⇩ Constraint c on the generated image (e.g., inpainting)

Constrained denoising step: $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, c) := \sum_{\tilde{\mathbf{x}}_0} q(\mathbf{x}_{t-1}|\tilde{\mathbf{x}}_0, \mathbf{x}_t) \cdot p_\theta(\tilde{\mathbf{x}}_0|\mathbf{x}_t, c)$

Computing or sampling from the constrained posterior $p_\theta(\tilde{\mathbf{x}}_0|\mathbf{x}_t, c)$ is **intractable** for diffusion models.

Denoising $p(\tilde{\boldsymbol{x}}_0|\boldsymbol{x}_t, \boldsymbol{x}_0^{\mathrm{k}}) \propto p_{\mathrm{DM}}(\tilde{\boldsymbol{x}}_0|\boldsymbol{x}_t, \boldsymbol{x}_0^{\mathrm{k}})^{\alpha} \cdot p_{\mathrm{TPM}}(\tilde{\boldsymbol{x}}_0|\boldsymbol{x}_t, \boldsymbol{x}_0^{\mathrm{k}})^{1-\alpha}$



$p_{\mathrm{DM}}(\tilde{\mathbf{x}}_0|\mathbf{x}_t, c)$

From the diffusion model:
Good at generating vivid details

$p_{\mathrm{TPM}}(\tilde{\mathbf{x}}_0|\mathbf{x}_t, c)$

From the probabilistic circuit:
Exact samples – better global coherence

# Controlling the denoiser with a probabilistic circuit



CoPaint

$p_{\mathrm{DM}}(\tilde{\mathbf{x}}_0|\mathbf{x}_t, \mathbf{x}_0^{\mathrm{k}})$

Tiramisu

$p_{\mathrm{DM}}(\tilde{\mathbf{x}}_0|\mathbf{x}_t, \mathbf{x}_0^{\mathrm{k}})$

$p(\tilde{\mathbf{x}}_0|\mathbf{x}_t, \mathbf{x}_0^{\mathrm{k}})$

$p_{\mathrm{TPM}}(\tilde{\mathbf{x}}_0|\mathbf{x}_t, \mathbf{x}_0^{\mathrm{k}})$

TPM not used

Denoising process

$t = 249$ $\quad$ $t = 217$ $\quad$ $t = 201$ $\quad$ $t = 100$ $\quad$ $t = 0$ $\qquad$ $t = 249$ $\quad$ $t = 217$ $\quad$ $t = 201$ $\quad$ $t = 100$ $\quad$ $t = 0$

# High-resolution image benchmarks

| Tasks | | Algorithms | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Dataset | Mask | Tiramisu (ours) | CoPaint | RePaint | DDNM | DDRM | DPS | Resampling |
| CelebA-HQ | Left | 0.189 | **0.185** | 0.195 | 0.254 | 0.275 | 0.201 | 0.257 |
| | Top | 0.187 | **0.182** | 0.187 | 0.248 | 0.267 | 0.187 | 0.251 |
| | Expand1 | **0.454** | 0.468 | 0.504 | 0.597 | 0.682 | 0.466 | 0.613 |
| | Expand2 | 0.442 | 0.455 | 0.480 | 0.585 | 0.686 | **0.434** | 0.601 |
| | V-strip | **0.487** | 0.502 | 0.517 | 0.625 | 0.724 | 0.535 | 0.647 |
| | H-strip | **0.484** | 0.488 | 0.517 | 0.626 | 0.731 | 0.492 | 0.639 |
| ImageNet | Left | **0.286** | 0.289 | 0.296 | 0.410 | 0.369 | 0.327 | 0.369 |
| | Top | **0.308** | 0.312 | 0.336 | 0.427 | 0.373 | 0.343 | 0.368 |
| | Expand1 | **0.616** | 0.623 | 0.691 | 0.786 | 0.726 | 0.621 | 0.711 |
| | Expand2 | **0.597** | 0.607 | 0.692 | 0.799 | 0.724 | 0.618 | 0.721 |
| | V-strip | 0.646 | 0.654 | 0.741 | 0.851 | 0.761 | **0.637** | 0.759 |
| | H-strip | 0.657 | 0.660 | 0.744 | 0.851 | 0.753 | **0.647** | 0.774 |
| LSUN-Bedroom | Left | **0.285** | 0.287 | 0.314 | 0.345 | 0.366 | 0.314 | 0.367 |
| | Top | **0.310** | 0.323 | 0.347 | 0.376 | 0.368 | 0.355 | 0.372 |
| | Expand1 | **0.615** | 0.637 | 0.676 | 0.716 | 0.695 | 0.641 | 0.699 |
| | Expand2 | **0.635** | 0.641 | 0.666 | 0.720 | 0.691 | 0.638 | 0.690 |
| | V-strip | **0.672** | 0.676 | 0.711 | 0.760 | 0.721 | 0.674 | 0.725 |
| | H-strip | 0.679 | 0.686 | 0.722 | 0.766 | 0.726 | **0.674** | 0.724 |
| Average | | **0.474** | 0.481 | 0.518 | 0.596 | 0.591 | 0.489 | 0.571 |

# Qualitative results on high-resolution image datasets

# Outline

1.  The paradox of learning to reason from data

    *~~deep learning~~*

2.  **Architectures for learning and reasoning**

    *logical reasoning + deep learning*

    a.  Constrained generative AI
    b.  **Constrained structured prediction**

# Warcraft Shortest Path



*// for a $12 \times 12$ grid, $2^{144}$ states but only $10^{10}$ valid ones!*
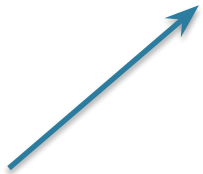
[Differentiation of Blackbox Combinatorial Solvers, Marin Vlastelica, Anselm Paulus, Vít Musil, Georg Martius, Michal Rolínek, 2019]

Baseline Prediction

Baseline Prediction

Baseline Prediction

Baseline Prediction

| ARCHITECTURE | EXACT MATCH | HAMMING SCORE | CONSISTENCY |
|---|---|---|---|
| RESNET-18+FIL | 55.0 | **97.7** | 56.9 |

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

*Are individual
edge predictions
correct?*

*Is output
a path?*

Kareem Ahmed, Eric Wang, Kai-Wei Chang and Guy Van den Broeck. Neuro-Symbolic Entropy Regularization, 2021.

# Declarative Knowledge of the Output

$y$



Neural Network

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

**1** Specify knowledge as a predicate

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

```
PyTorch Code

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

Kareem Ahmed, Tao Li, Thy Ton, Quan Guo, Kai-Wei Chang, Parisa Kordjamshidi, Vivek Srikumar, Guy Van den Broeck and Sameer Singh. PYLON: A PyTorch Framework for Learning with Constraints

# py**lon**

```
PyTorch Code

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

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

**②** Add as loss to training

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

Kareem Ahmed, Tao Li, Thy Ton, Quan Guo, Kai-Wei Chang, Parisa Kordjamshidi, Vivek Srikumar, Guy Van den Broeck and Sameer Singh. PYLON: A PyTorch Framework for Learning with Constraints

# py**lon**

① Specify knowledge as a predicate

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

② Add as loss to training

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

③ pylon derives the gradients
(solves a combinatorial problem)

```
PyTorch Code

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

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

Kareem Ahmed, Tao Li, Thy Ton, Quan Guo, Kai-Wei Chang, Parisa Kordjamshidi, Vivek Srikumar, Guy Van den Broeck and Sameer Singh. PYLON: A PyTorch Framework for Learning with Constraints

*without constraint*    *with constraint*    *without constraint*    *with constraint*

Baseline Prediction    SL Prediction    Baseline Prediction    SL Prediction

Baseline Prediction    SL Prediction    Baseline Prediction    SL Prediction

$p(\mathbf{y}|x)$

$\longmapsto\ m(\alpha)\ \longmapsto$ **y**

a) A network uncertain over both valid & invalid predictions

$p(\mathbf{y}|x)$

$\longmapsto\ m(\alpha)\ \longmapsto$ **y**

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

**Semantic Loss**

$$\mathrm{L}^{\mathrm{s}}(\alpha, \mathbf{p}) \propto -\log \sum_{\mathbf{x} \models \alpha} \prod_{i:\mathbf{x} \models X_i} \mathbf{p}_i \prod_{i:\mathbf{x} \models \neg X_i} (1 - \mathbf{p}_i)$$

Probability of satisfying constraint α after sampling from neural net output layer **p**

In general: #P-hard ☹

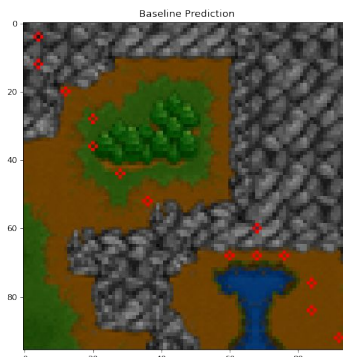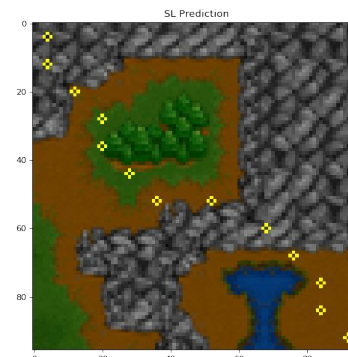Do this probabilistic-logical reasoning during learning in a computation graph

α: A ∧ B => C

Semantic Loss

- log( )

Probability

| without constraint | with constraint | without constraint | with constraint |

| ARCHITECTURE | EXACT MATCH | HAMMING SCORE | CONSISTENCY |
|---|---|---|---|
| RESNET-18+FIL | 55.0 | **97.7** | 56.9 |
| RESNET-18+$\mathcal{L}_{SL}$ | 59.4 | **97.7** | 61.2 |

# 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



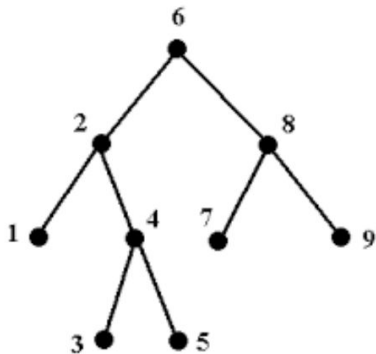- Secret sauce is again tractable circuits – computation graphs for reasoning

Kareem Ahmed, Stefano Teso, Kai-Wei Chang, Guy Van den Broeck and Antonio Vergari. Semantic Probabilistic Layers for Neuro-Symbolic Learning, 2022.

| GROUND TRUTH | RESNET-18 | SEMANTIC LOSS | SPL (ours) |

| ARCHITECTURE | EXACT MATCH | HAMMING SCORE | CONSISTENCY |
| --- | --- | --- | --- |
| RESNET-18+FIL | 55.0 | **97.7** | 56.9 |
| RESNET-18+$\mathcal{L}_{SL}$ | 59.4 | **97.7** | 61.2 |
| RESNET-18+SPL | 75.1 | 97.6 | **100.0** |
| OVERPARAM. SDD | **78.2** | 96.3 | **100.0** |

Kareem Ahmed, Stefano Teso, Kai-Wei Chang, Guy Van den Broeck and Antonio Vergari. Semantic Probabilistic Layers for Neuro-Symbolic Learning, 2022.
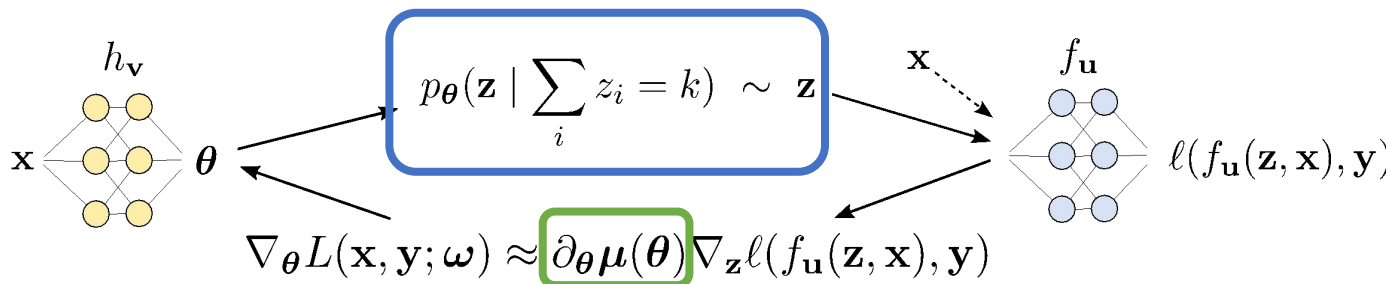
# Hierarchical Multi-Label Classification



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

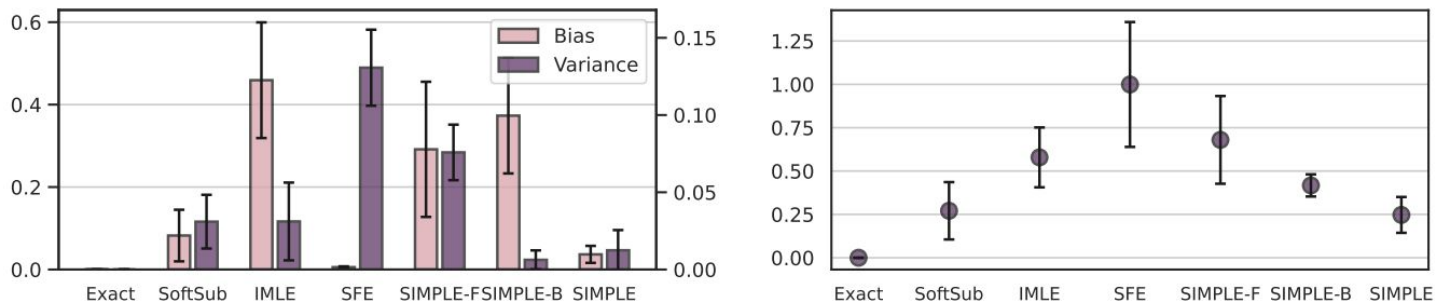"if the image is classified as an animal, it must be classified as either cat or dog"

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

# **SIMPLE**: Gradient Estimator for *k*-Subset Sampling



$$p_{\boldsymbol{\theta}}(\mathbf{z} \mid \sum_i z_i = k) \ \sim \ \mathbf{z}$$

$$\ell(f_{\mathbf{u}}(\mathbf{z}, \mathbf{x}), \mathbf{y})$$

$$\nabla_{\boldsymbol{\theta}} L(\mathbf{x}, \mathbf{y}; \boldsymbol{\omega}) \approx \boxed{\partial_{\boldsymbol{\theta}} \boldsymbol{\mu}(\boldsymbol{\theta})} \nabla_{\mathbf{z}} \ell(f_{\mathbf{u}}(\mathbf{z}, \mathbf{x}), \mathbf{y})$$

We achieve ***lower bias and variance*** by exact, discrete samples and exact derivative of conditional marginals.



and SotA Learning to Explain (L2X) and sparse discrete VAE results.

Kareem Ahmed, Zhe Zeng, Mathias Niepert, Guy Van den Broeck. SIMPLE: A Gradient Estimator for k-Subset Sampling, ICLR 2023

# *Secret Sauce: Probabilistic Circuits*

**Tutorial (3h)**



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

**Overview Paper (80p)**



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

# Outline

1. The paradox of learning to reason from data

    *deep learning*

2. Architectures for learning and reasoning

    *logical (and probabilistic) reasoning + deep learning*

    a. Constrained generative AI
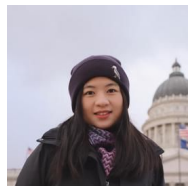    b. Constrained structured prediction

# Thanks

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



Honghua    Kareem    Zhe    Meihua    Anji

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