

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

Guy Van den Broeck

TTI/Vanguard - Sep 13 2022

# Outline

1. The paradox of learning to reason from data

*~~deep learning~~*

2. Learning with symbolic knowledge

*logical reasoning + deep learning*

# Outline

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



*~~deep learning~~*

## 2. 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.*

## Logical 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

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.

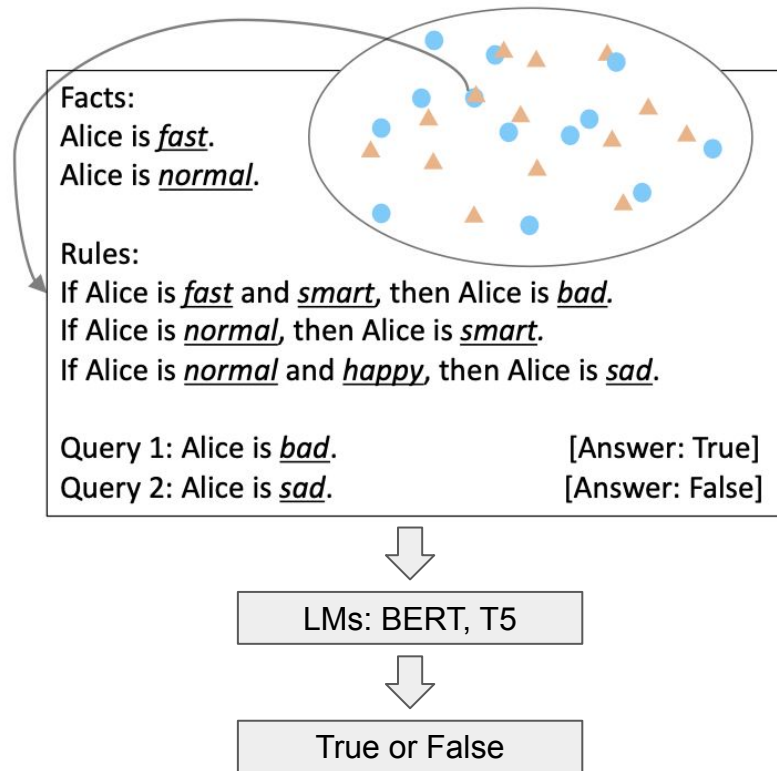
Facts: Alice sensible Alice condemned Alice thoughtless Alice polite Alice scared Alice average

Query: Alice is shy ?

# Problem Setting: SimpleLogic

The easiest of reasoning problems:

1. **Propositional logic** fragment
  - a. bounded vocabulary & number of rules
  - b. bounded reasoning depth ( $\leq 6$ )
  - 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

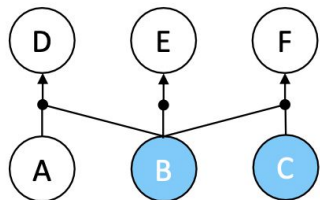


# Training a BERT model on SimpleLogic

(1) Randomly sample facts & rules.

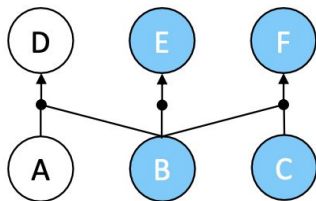
Facts: B, C

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:

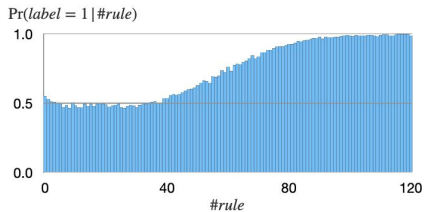
*Any rules can be freely added to the hypothesis of any proven fact.*



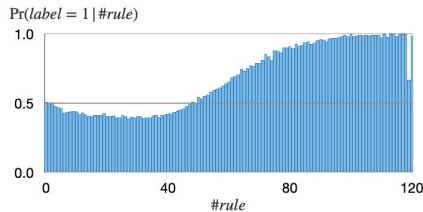
The more rules given, the more likely a predicate will be proved.



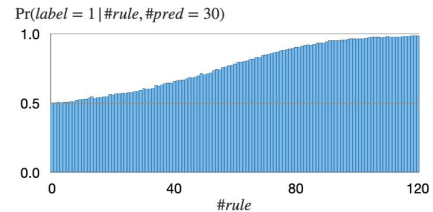
$\Pr(\text{label} = \text{True} \mid \text{Rule \#} = x)$  should increase (roughly) monotonically with  $x$



(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.

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

# Outline

1. The paradox of learning to reason from data

*~~deep learning~~*

2. **Learning with symbolic knowledge**

*logical reasoning + deep learning*

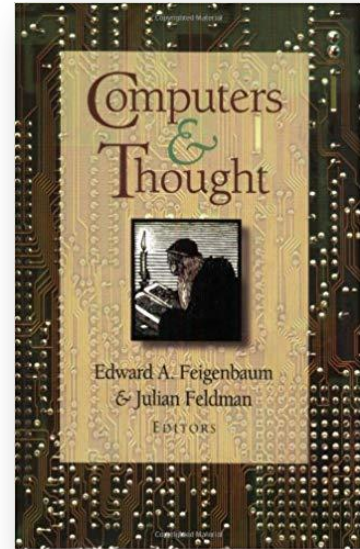
# The AI Dilemma of 2022

## **Deep learning**

*approaches the problem of designing intelligent machines by postulating a large number of very simple information processing elements, arranged in a [...] network, and certain processes for facilitating or inhibiting their activity.*

## **Knowledge representation and reasoning**

*take a much more macroscopic approach [...]. They believe that intelligent performance by a machine is an end difficult enough to achieve without “starting from scratch” , and so they build into their systems as much complexity of information processing as they are able to understand and communicate to a computer.*



Edward Feigenbaum  
and Julian Feldman

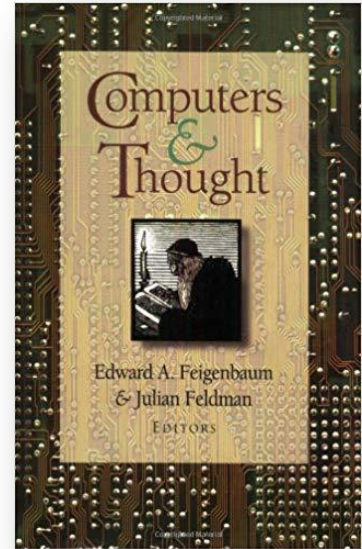
# The AI Dilemma of ~~2022~~ 1963

## **Neural cybernetics**

*approaches the problem of designing intelligent machines by postulating a large number of very simple information processing elements, arranged in a [...] network, and certain processes for facilitating or inhibiting their activity.*

## **Cognitive model builders**

*take a much more macroscopic approach [...]. They believe that intelligent performance by a machine is an end difficult enough to achieve without “starting from scratch”, and so they build into their systems as much complexity of information processing as they are able to understand and communicate to a computer.*



Edward Feigenbaum  
and Julian Feldman

# 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

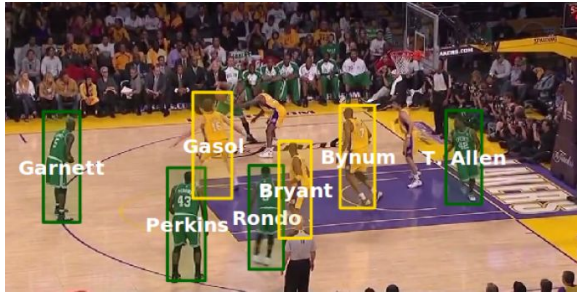


# The AI Dilemma



Integrate reasoning into modern deep learning algorithms

# Knowledge in Vision, Robotics, NLP, Activity Recognition

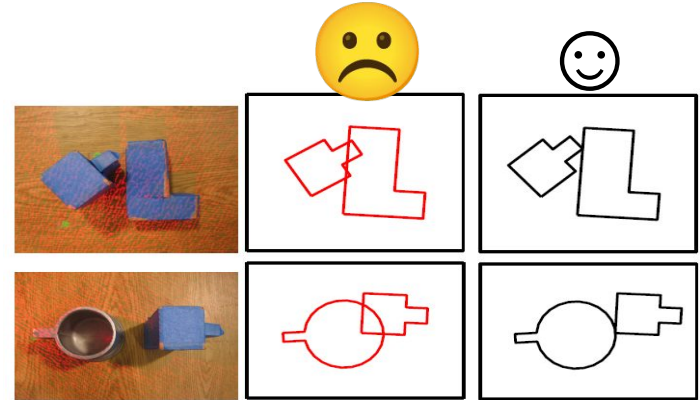


People appear at most once in a frame

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

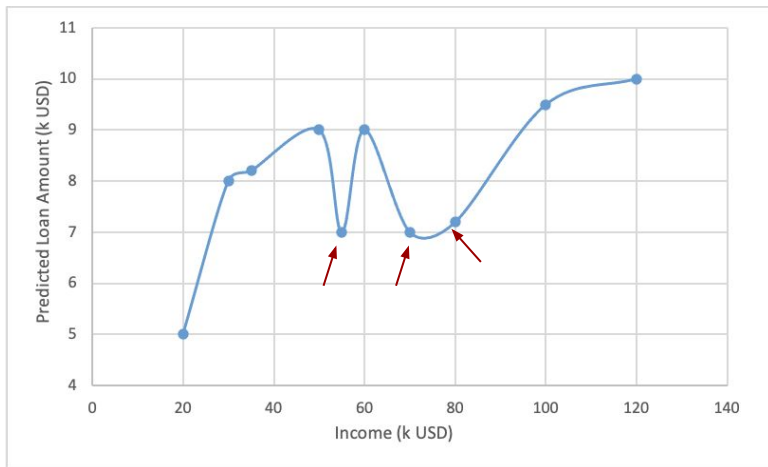


Cut the orange before squeezing the orange



Rigid objects don't overlap

# Predict Loan Amount



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

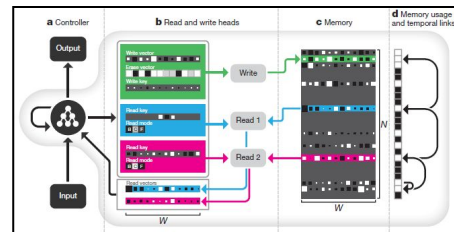
Monotonicity (Prior Knowledge):

**Increasing income should increase the approved loan amount**

# 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 article title is 'DeepMind's AI has learned to navigate the Tube using memory'. Below the title is a photograph of a person with long blonde hair looking at a London Underground tube map. The article includes social media sharing icons and a date of 'DAILY NEWS 12 October 2016'.

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 is a short summary: 'DeepMind's latest technique uses external memory to solve tasks that require logic and reasoning — a step toward more human-like AI.' The article is dated '2016 November' and is categorized as 'News & Comment'.



[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 ...



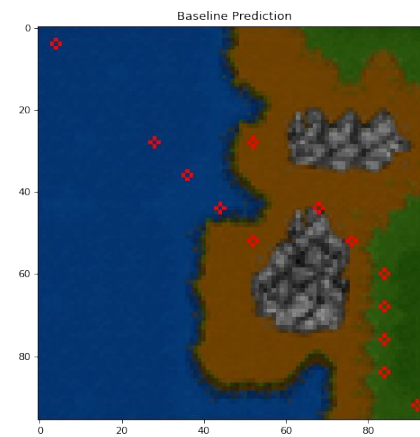
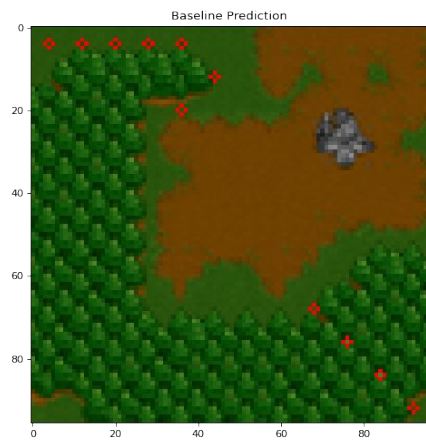
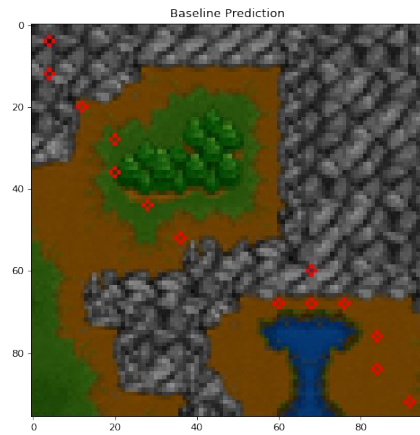
optimal planner recalculating a shortest path to the end node. To ensure that the network always moved to a valid node, the output distribution was renormalized over the set of possible triples outgoing from the current node. The performance

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'

# Warcraft Shortest Path

Predicting the minimum-cost path







# 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

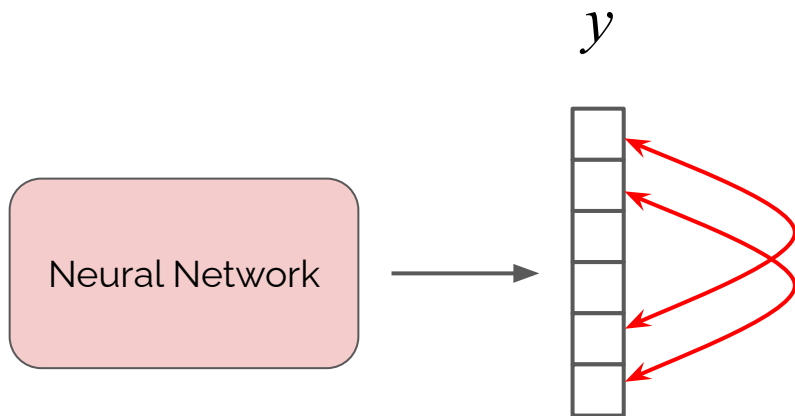
# pylon

A PyTorch Framework for Learning with Constraints

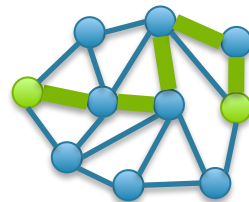
Kareem Ahmed   Tao Li   Thy Ton   Quan Guo,  
Kai-Wei Chang   Parisa Kordjamshidi   Vivek Srikumar  
Guy Van den Broeck   Sameer Singh

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

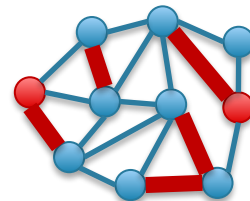
# Declarative Knowledge of the Output



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

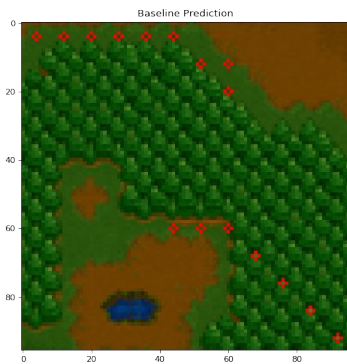
```
def check(y):  
    ...  
    return isValid
```

2 Add as loss to training

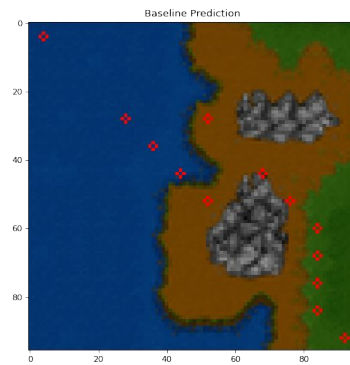
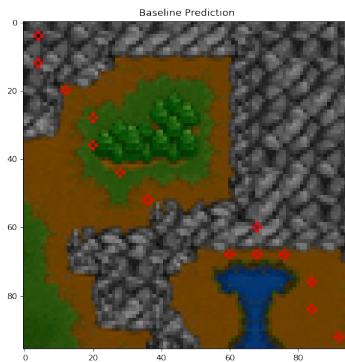
```
loss += constraint_loss(check)
```

3 pylon derives the gradients  
(solves a combinatorial problem)

*without constraint*



*without constraint*





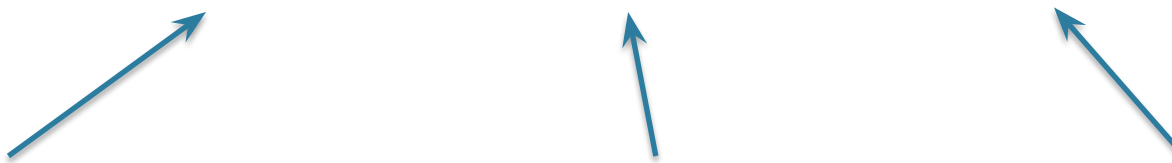
# Warcraft min-cost simple-path prediction results

Test accuracy %	Coherent	Incoherent	Constraint
ResNet-18	44.8	<b>97.7</b>	56.9

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

*Are individual edge predictions correct?*

*Is output a path?*



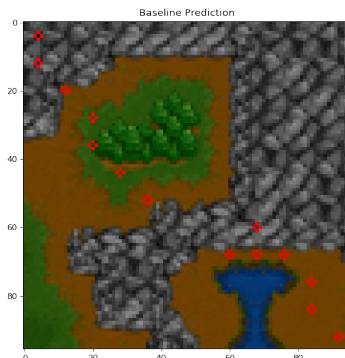
*without constraint*



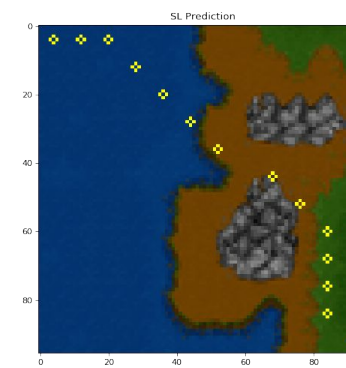
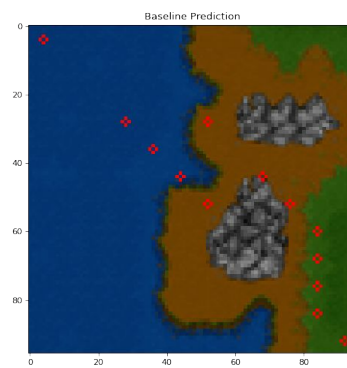
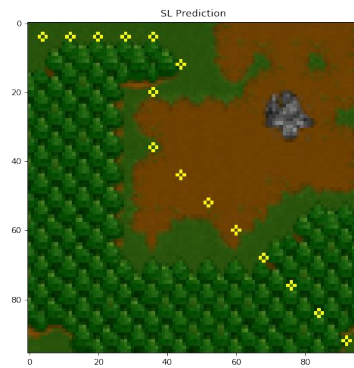
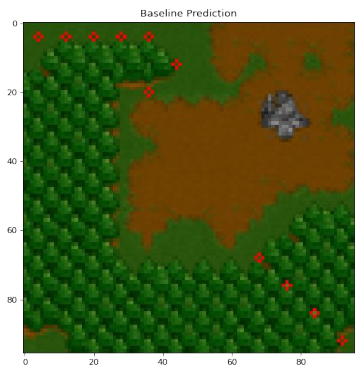
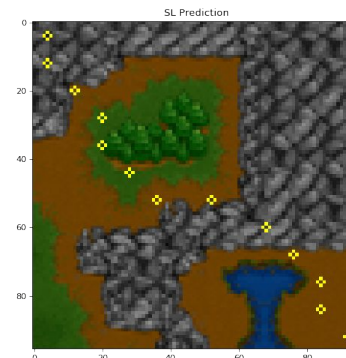
*with constraint*



*without constraint*

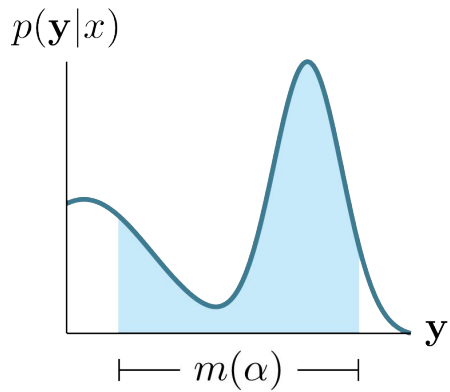


*with constraint*

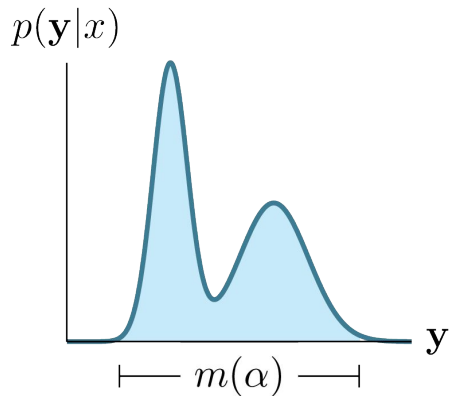


# Warcraft min-cost simple-path prediction results

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

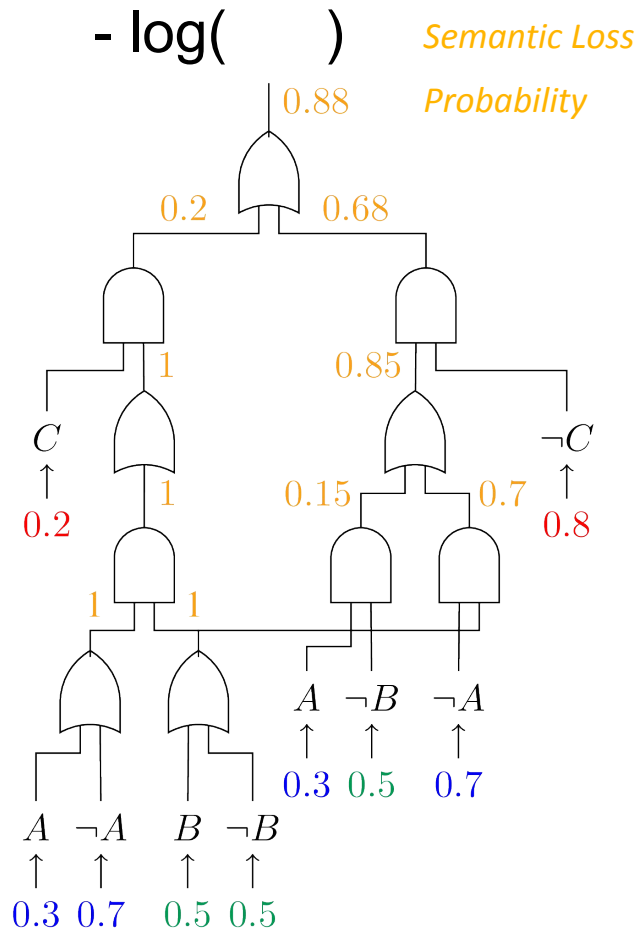
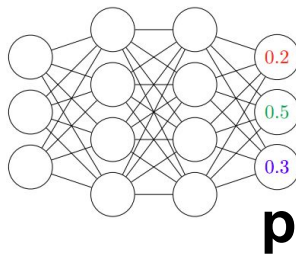
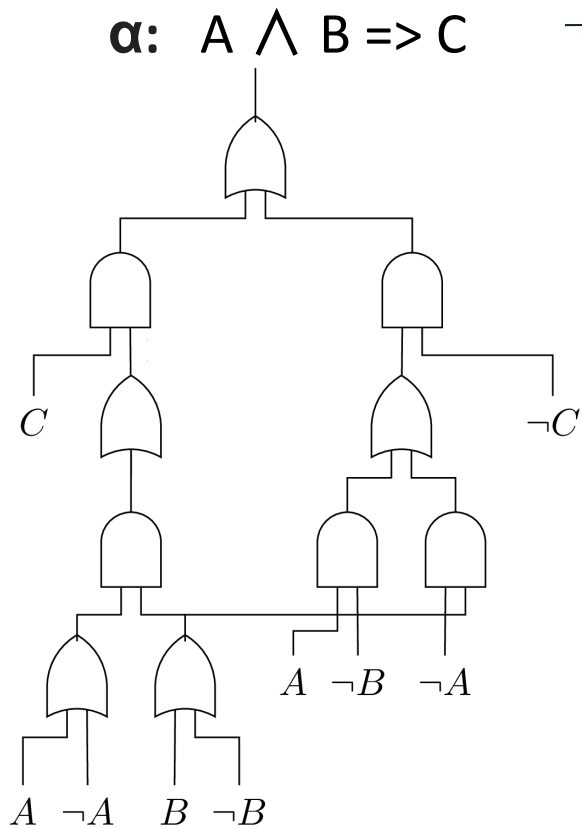


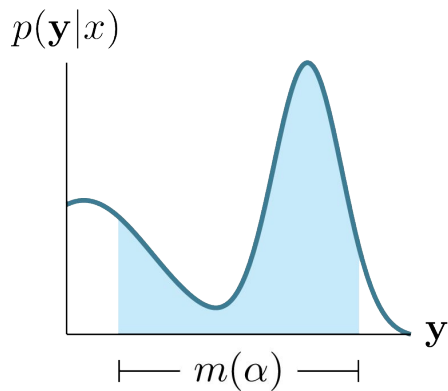
a) A network uncertain over both valid & invalid predictions



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

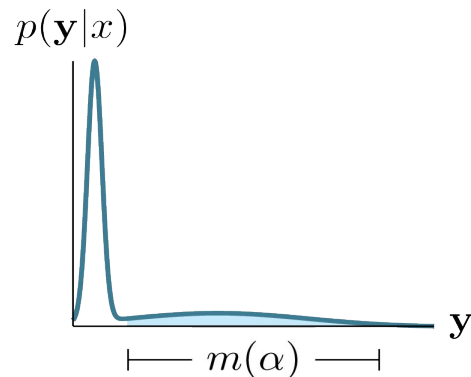
↓  
**Neuro-Symbolic Learning**





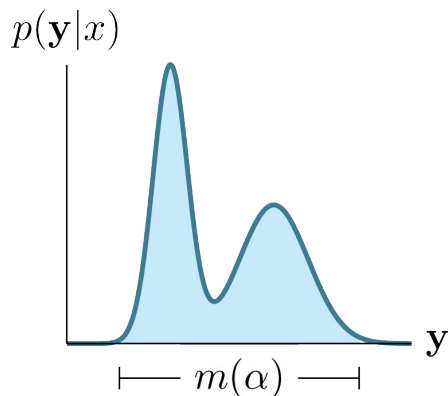
a) A network uncertain over both valid & invalid predictions

**Entropy  
Regularization**



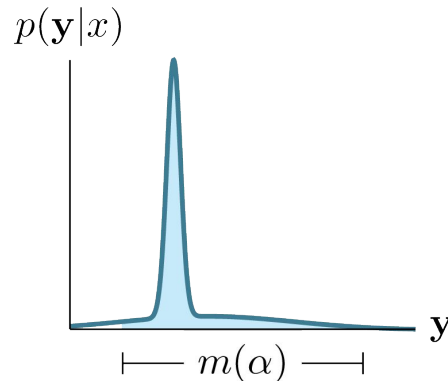
b) A network allocating most of its mass to an invalid prediction.

**Neuro-Symbolic  
Learning**



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

**Neuro-Symbolic  
Entropy Regularization**



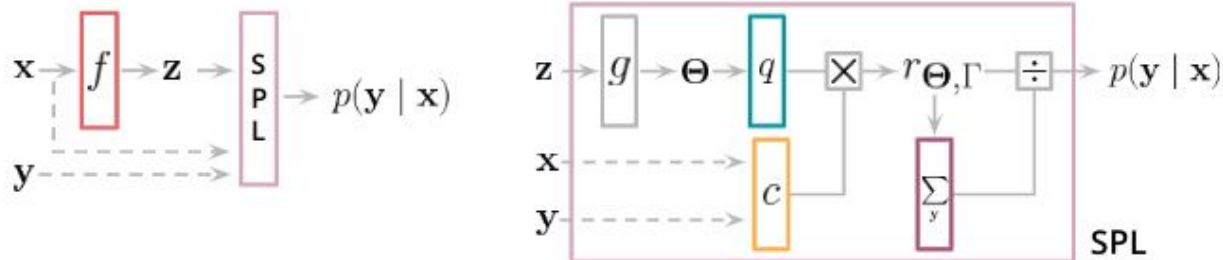
d) A network allocating most of mass to one model of formula

# 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	<b>14.80 ± 3.70</b>	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	<b>18.38 ± 2.50</b>	<b>26.41 ± 0.49</b>	<b>31.17 ± 1.68</b>	<b>35.85 ± 0.75</b>	<b>37.62 ± 2.17</b>	<b>41.28 ± 0.46</b>
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	<b>6.50 ± 2.00</b>	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	<b>9.31 ± 0.76</b>	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	<b>22.42 ± 0.43</b>	24.37 ± 1.62
	+ NeSy Entropy	6.19 ± 2.40	8.11 ± 3.66	<b>13.17 ± 1.08</b>	<b>15.47 ± 2.19</b>	<b>17.45 ± 1.52</b>	22.14 ± 1.46	<b>25.11 ± 1.03</b>

# 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



# Warcraft Shortest Path



GROUND TRUTH



RESNET-18



SEMANTIC LOSS



SPL (ours)

# 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	<b>3.79 ± 0.18</b>
DERISI	1.39 ± 0.47	<b>2.28 ± 0.23</b>
EISEN	5.40 ± 0.15	<b>6.18 ± 0.33</b>
EXPR	4.20 ± 0.21	<b>5.54 ± 0.36</b>
GASCH1	3.48 ± 0.96	<b>4.65 ± 0.30</b>
GASCH2	3.11 ± 0.08	<b>3.95 ± 0.28</b>
SEQ	5.24 ± 0.27	<b>7.98 ± 0.28</b>
SPO	<b>1.97 ± 0.06</b>	<b>1.92 ± 0.11</b>
DIATOMS	48.21 ± 0.57	<b>58.71 ± 0.68</b>
ENRON	5.97 ± 0.56	<b>8.18 ± 0.68</b>
IMCLEF07A	79.75 ± 0.38	<b>86.08 ± 0.45</b>
IMCLEF07D	76.47 ± 0.35	<b>81.06 ± 0.68</b>

# Outline

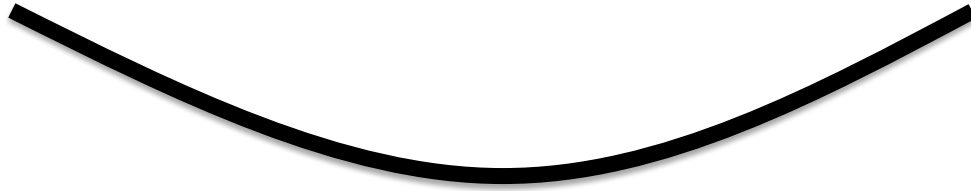
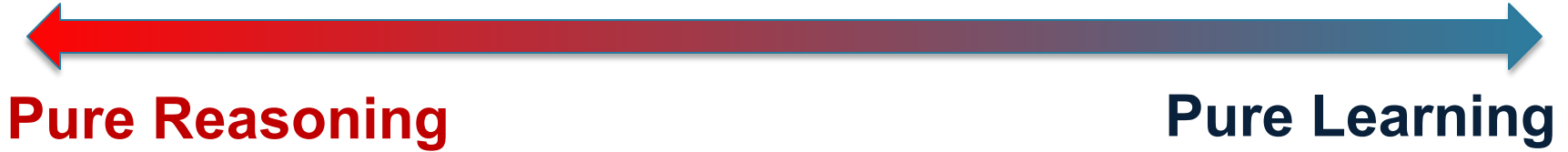
1. The paradox of learning to reason from data

*~~deep learning~~*

2. Learning with symbolic knowledge

*logical (and probabilistic) reasoning + deep learning*

# The AI Dilemma



Integrate reasoning into modern deep learning algorithms

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

# Thanks

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

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