

Neural Symbolic Language Understanding

Ni Lao

11.12.2017

Everything presented here is publicly available.
The opinions stated here are my own, not those of Google.

Plan

- **Language & control**
 - *Neural Symbolic Machines*: Semantic Parsing on Freebase with Weak Supervision
 - Chen Liang, Jonathan Berant, Quoc Le, Kenneth Forbus, Ni Lao
- Knowledge & scalability
 - Learning to Organize Knowledge with An *N-Gram Machine*
 - Fan Yang, Jiazhong Nie, William Cohen, Ni Lao

Language & Reasoning

- Language was primarily invented for reasoning
- Communication comes later



Robert C. Berwick • Noam Chomsky

Language, Translation & Control



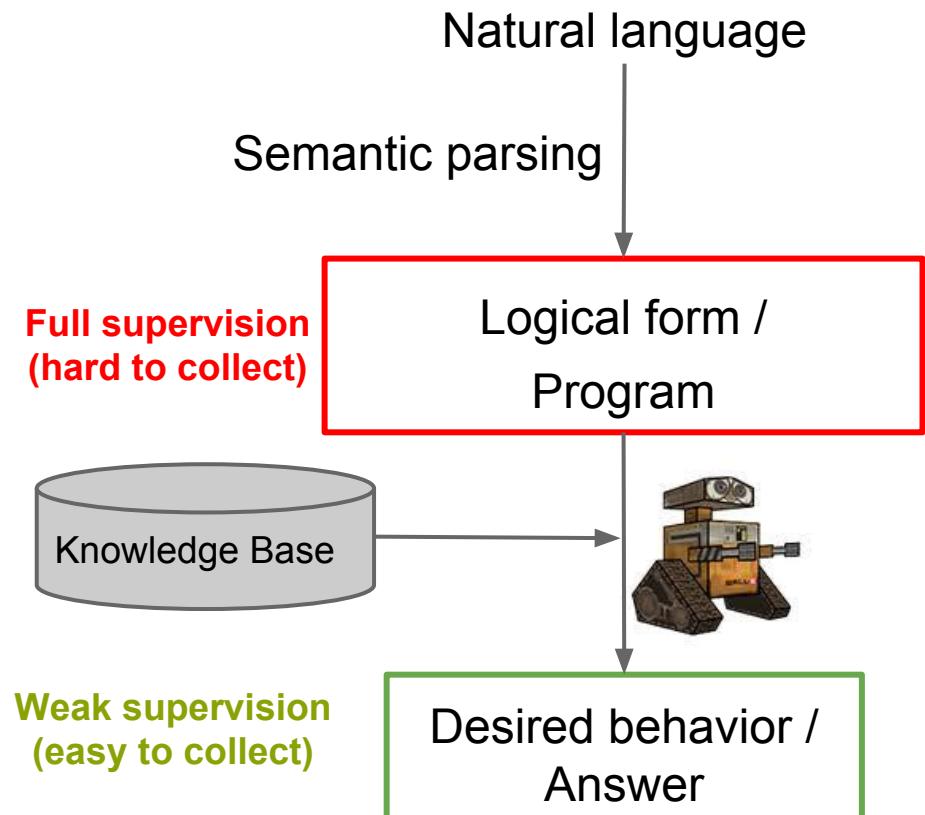
LOGIC AND
MATHEMATICS ARE
NOTHING BUT
SPECIALISED
LINGUISTIC
STRUCTURES.

Jean Piaget

- 1) Natural languages are programming languages to control human behavior
- 2) For machines and human to understand each other, they just need translation models trained with control theory

Semantic Parsing

- Question answering with structured data (KG / tables / personal data)
- Voice to action
- Personal assistant
- Entertainment

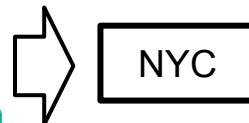
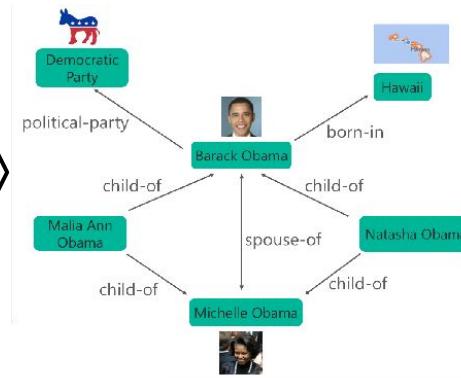


Question Answering with Knowledge Base

Largest city in US?



```
GO  
(Hop V1 CityIn)  
(Argmax V2 Population)  
RETURN
```



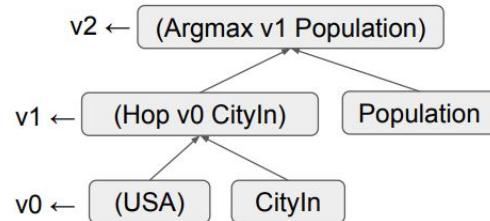
NYC



0. Paraphrase

Many ways to ask the same question, e.g.,
“What was the date that Minnesota became a state?”
“When was the state Minnesota created?”

1. Compositionality



2. Large Search Space

E.g., Freebase:
23K predicates,
82M entities,
417M triplets

3. Optimization

Reinforcement Learning is known to be hard given sparse reward

WebQuestionsSP Dataset

- 5,810 questions Google Suggest API & Amazon MTurk¹
- Remove invalid QA pairs²
- 3,098 training examples, 1,639 testing examples remaining
- Open-domain, and contains grammatical error
- Multiple entities as answer => macro-averaged F1

Grammatical error

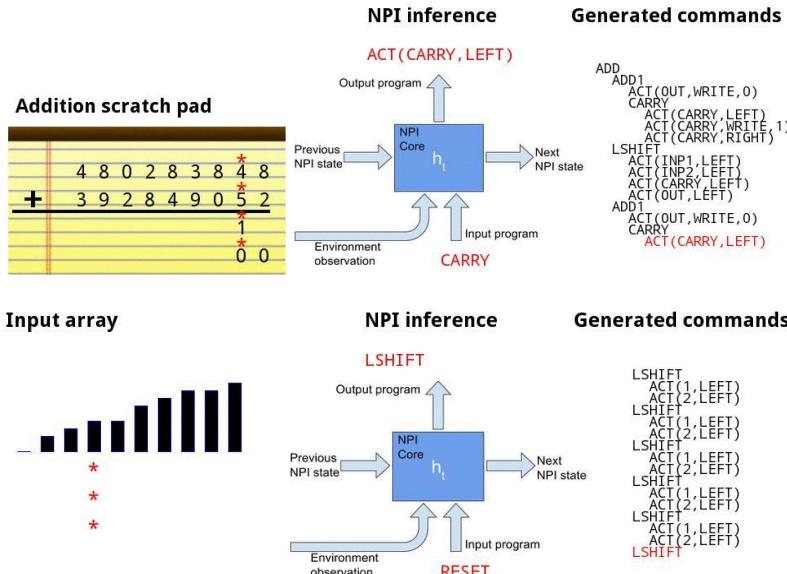
- What **do** Michelle Obama do for a living?
- What character did Natalie Portman play in Star Wars?
- What currency do you use in Costa Rica?
- What did Obama study in school?
- What killed Sammy Davis Jr?

Multiple entities

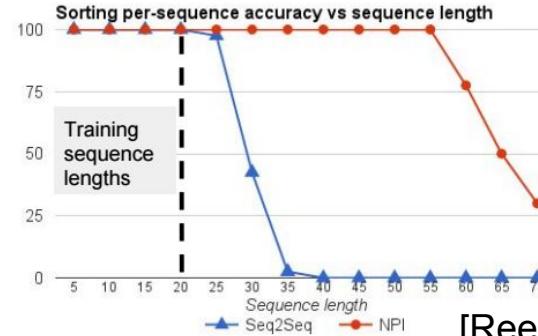
writer, lawyer
Padme Amidala
Costa Rican colon
political science
throat cancer

(Scalable) Neural Program Induction

- Impressive works to show NN can learn addition and sorting, but...



- The learned operations are not as scalable and precise.



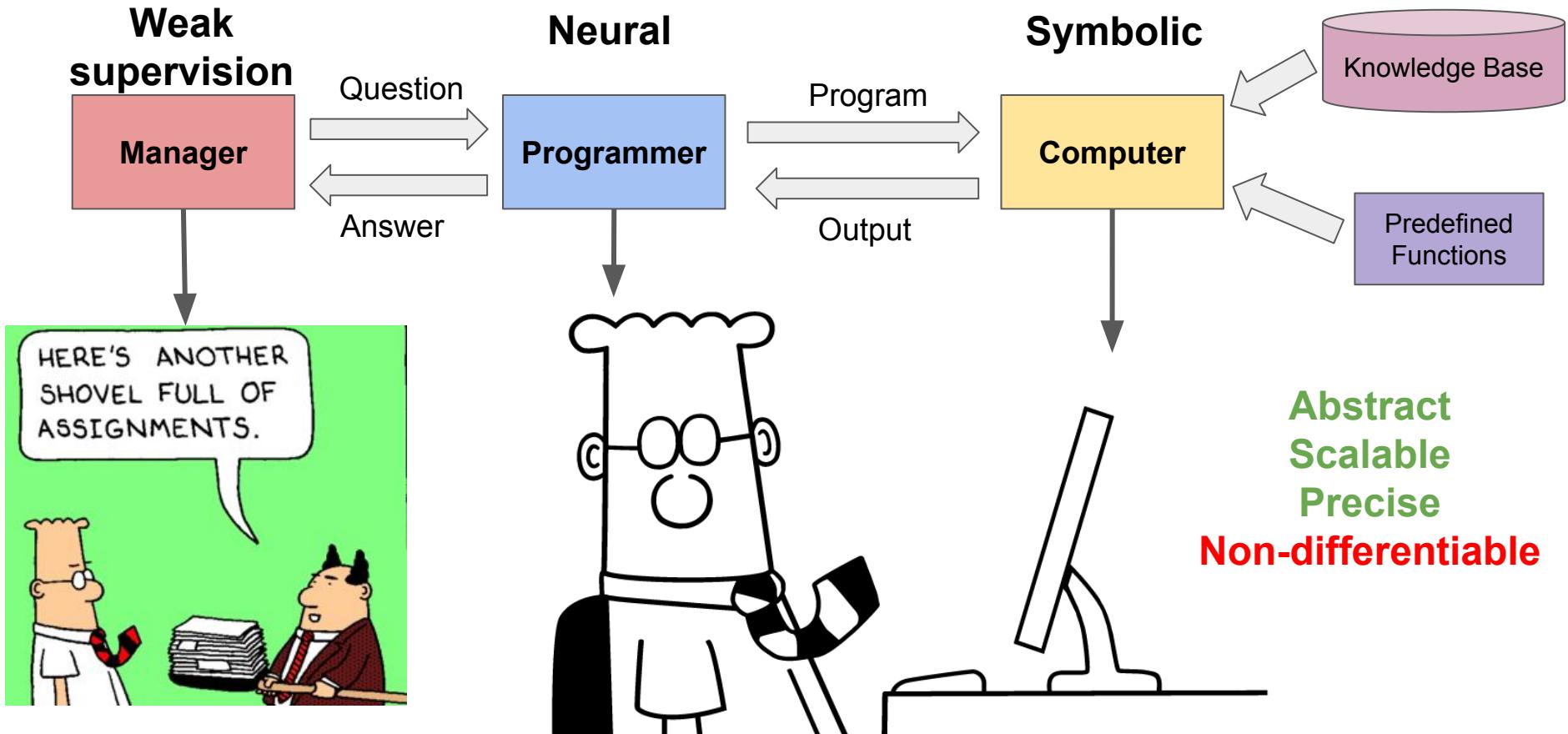
[Reed & Freitas 2015]

- Why not use existing modules that are scalable, precise and interpretable?



[Zaremba & Sutskever 2016]

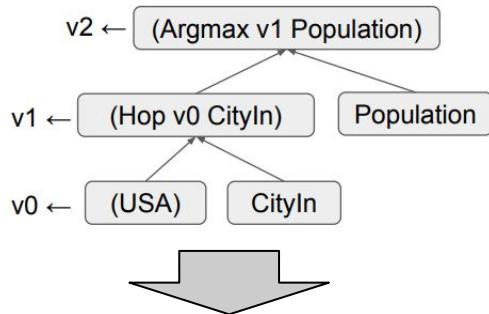
Neural Symbolic Machines



Contributions

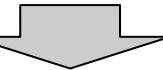
- Simple Seq2Seq model is not enough

1. Compositionality



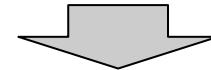
2. Large Search Space

E.g., Freebase:
23K predicates,
82M entities,
417M triplets



3. Optimization

Reinforcement
Learning is known
to be hard given
sparse reward

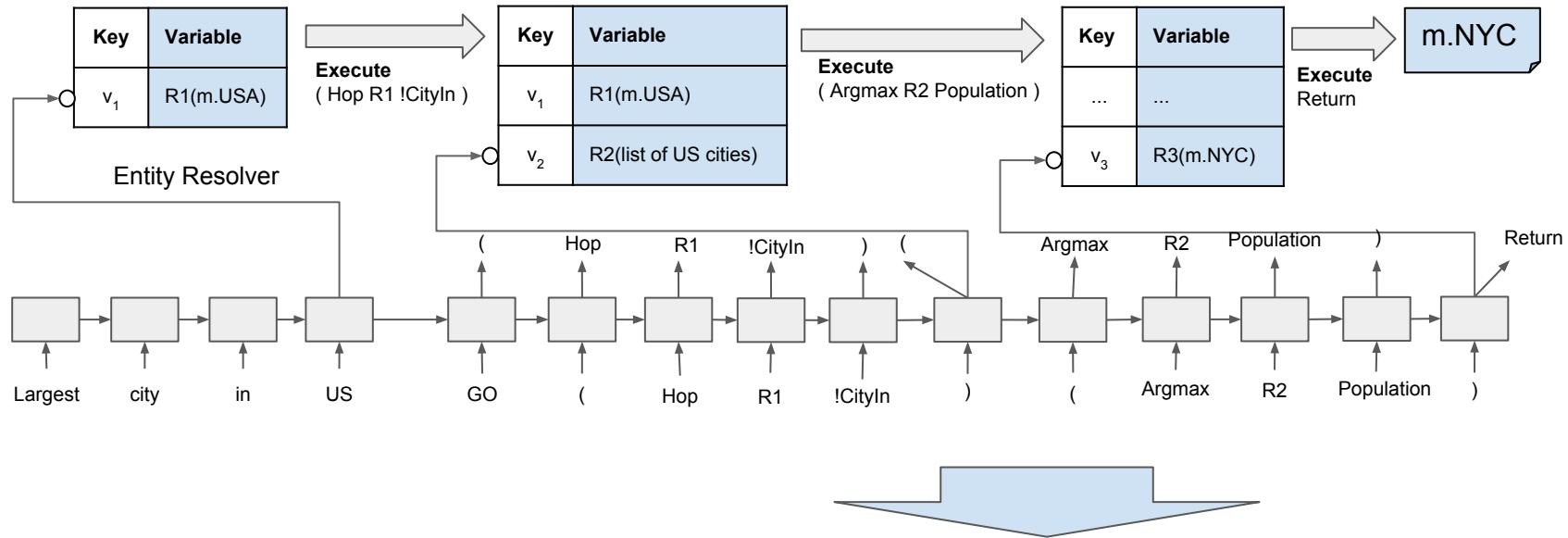


1. Key-Variable Memory

2. Beam search + Code Assistance

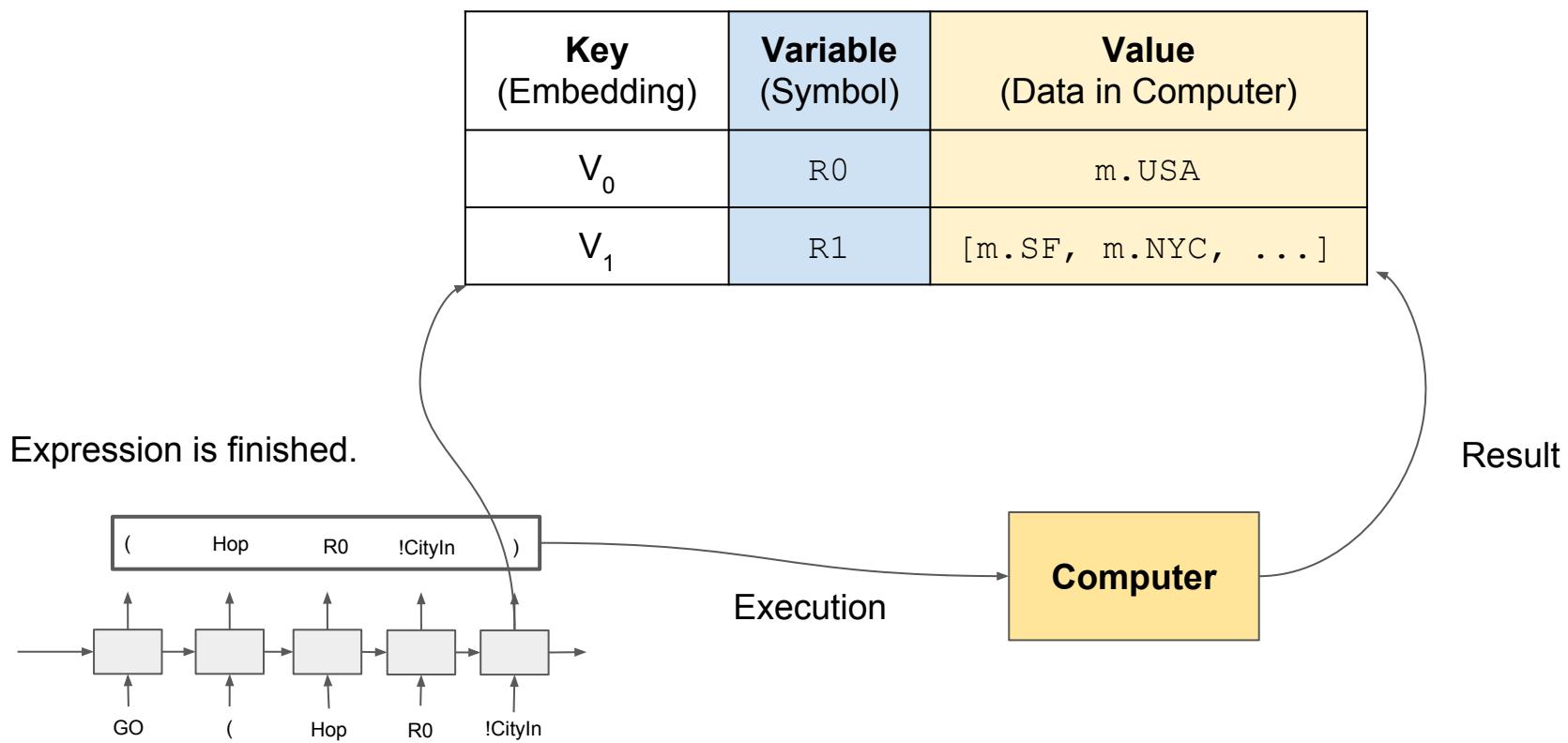
3. Augmented REINFORCE

Key-Variable Memory for Compositionality

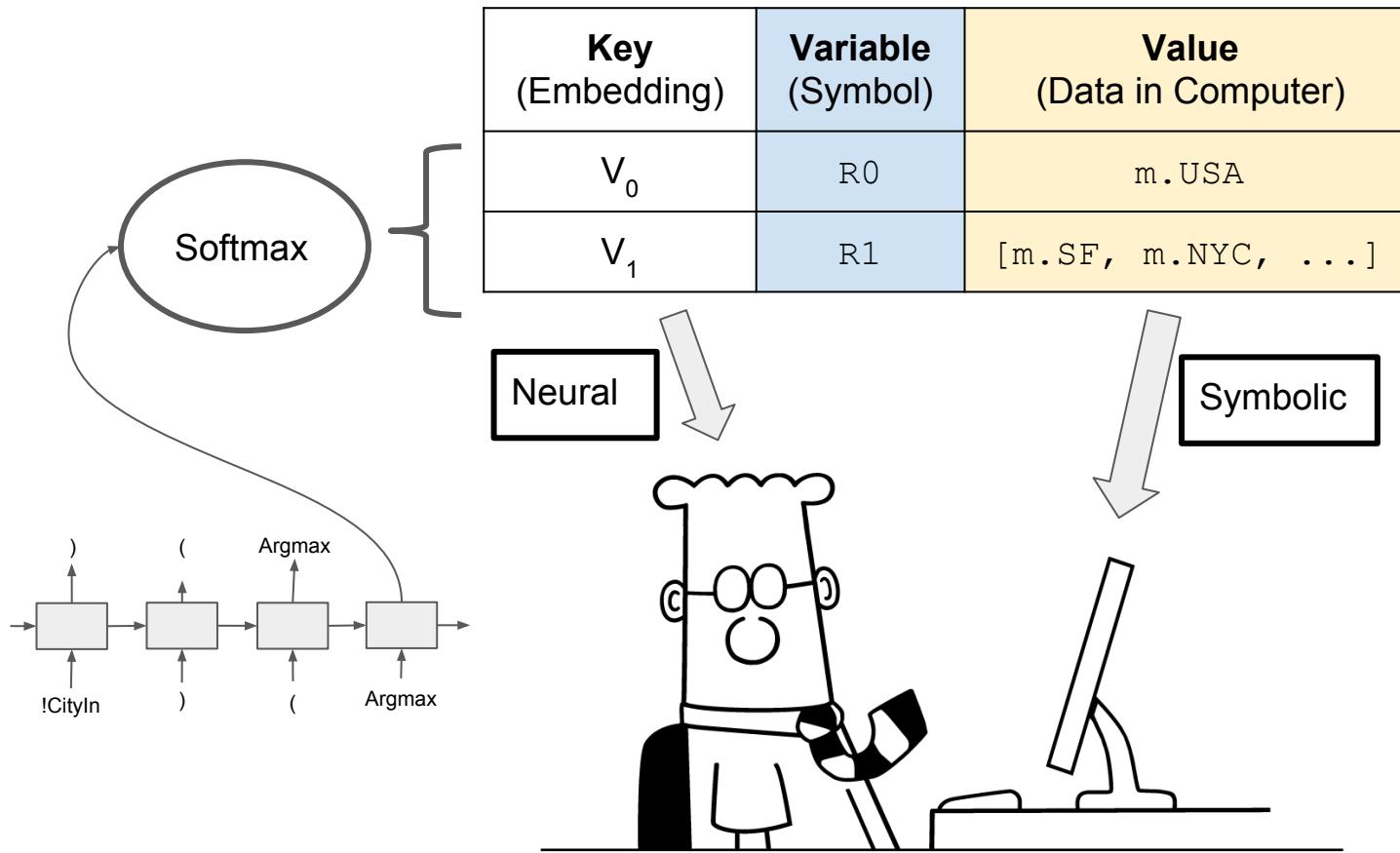


- A linearised bottom-up derivation of the recursive program

Key-Variable Memory: Save Intermediate Value



Key-Variable Memory: Reuse Intermediate Value



Code Assistance: Prune Search Space



Pen and paper

The screenshot shows a C++ IDE interface with the title "Dev-C++ 4.9.7.0 [Rhythmix] Rhythmix5.dev". The project tree on the left shows files under "Rhythmix" and "Rhythmix5". The main code editor window displays a portion of the "MidiThread.cpp" file:

```
bool MidiThread::Poll()
{
}

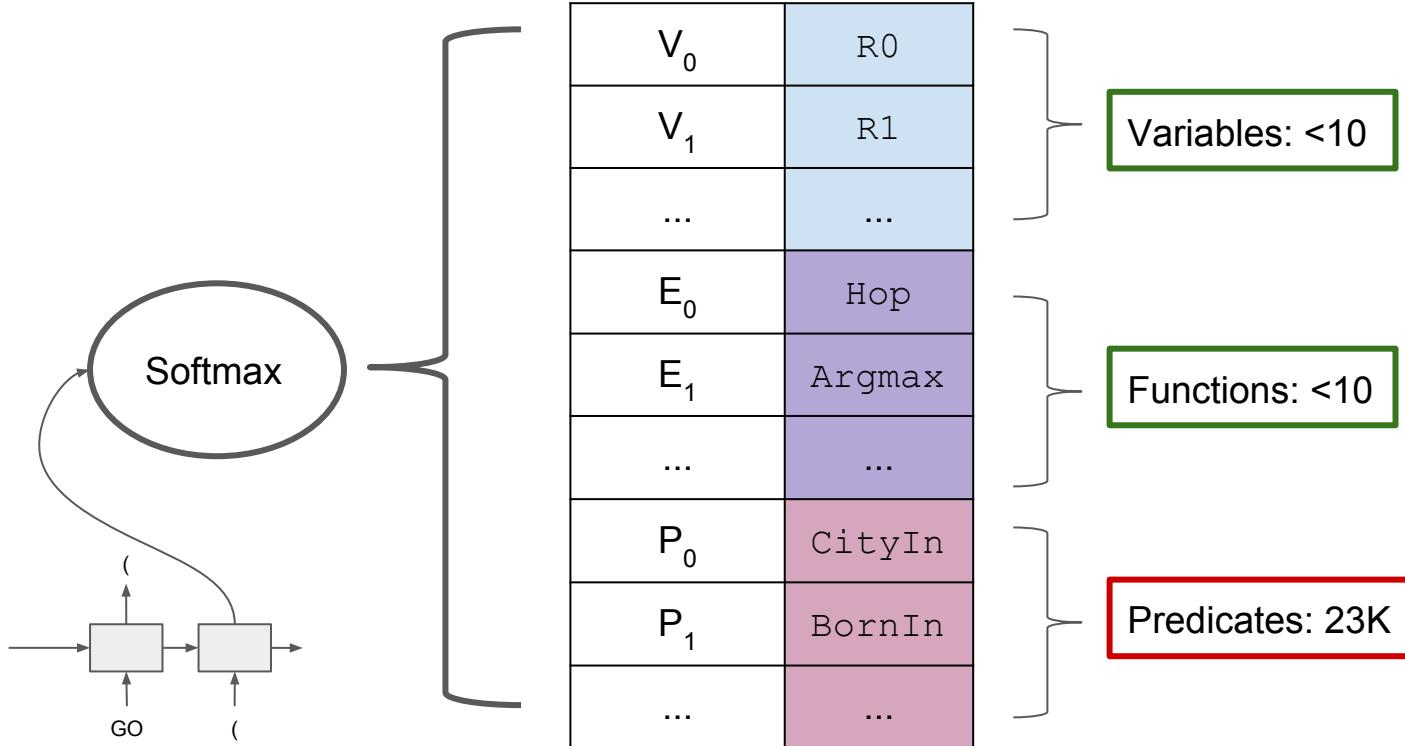
void MidiThread::AnalyseMidi()
{
    if ((midi_msg[0] == M_START) || (midi_msg[0] == M_CONT))
    {
        frame->Play(true);
    }
    else if (midi_msg[0] == X_STOP)
    {
        frame->Stop();
        BeatSeq->
    }
    else if (Cont
    Controller->
    else if (STA
    {
        BeatSeq->
    }
    else if (STAT
    {
        Destruitor *BeatSequencer()
        {
            void AddTrack(TrackData *t)
            BeatSequencer (wNameParent *parent,
            void ClearBeats()
            void ClearTracks()
            void HideBeat (Rect *b)
            void OnMouseEvent (wMouseEvent &ev,
            void OnScroll (wScorEvent &event)
            void OnSize (wSizeEvent &event)
            void PaintMeasure
            void ProcessMidi (int chan, int cont_numb
            void RefreshView()
            void RemoveSelected (void)
            void SelectAll (void)
        }
    }
}
```

The code editor has a red highlight over the first if-statement. A tooltip box is open over the word "Destruitor" in the BeatSequencer class definition, listing various member functions and variables. The status bar at the bottom indicates "4:1 Modified Insert 24 Lines in file".

IDE

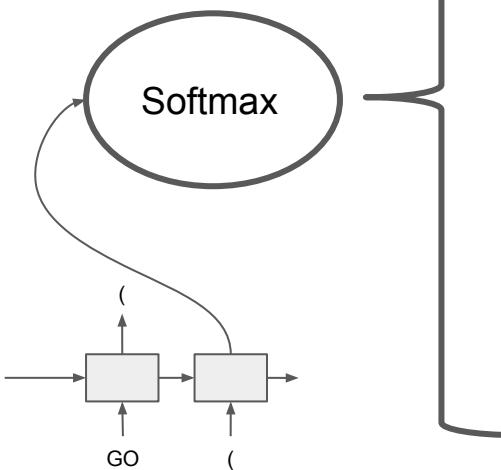
Code Assistance: Syntactic Constraint

Decoder Vocab



Code Assistance: Syntactic Constraint

Last token is '(', so
has to output a
function name next.



Decoder Vocab

The Decoder Vocab is represented as a grid of tokens. The columns are labeled V_0 , E_0 , E_1 , ..., and the rows are labeled BP , Hop , $Argmax$, ..., $Cit(A)$, Pn , In , ..., and Red 'X' marks indicate prohibited transitions or values. Brackets on the right side group specific entries:

V_0	E_0	E_1	\dots	BP	Hop	$Argmax$	\dots	$Cit(A)$	Pn	In	\dots

Variables: <10

Functions: <10

Predicates: 23K

Code Assistance: Semantic Constraint

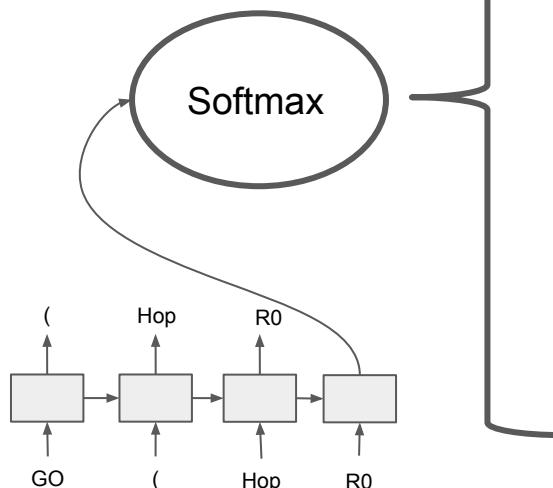
Decoder Vocab

V_0	R0
V_1	R1
...	...
E_0	Hop
E_1	Argmax
...	...
P_0	CityIn
P_1	BornIn
...	...

Variables: <10

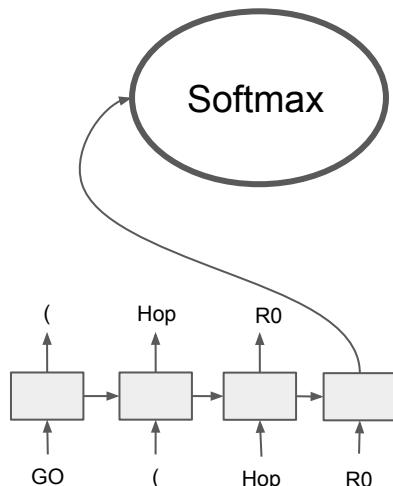
Functions: <10

Predicates: 23K

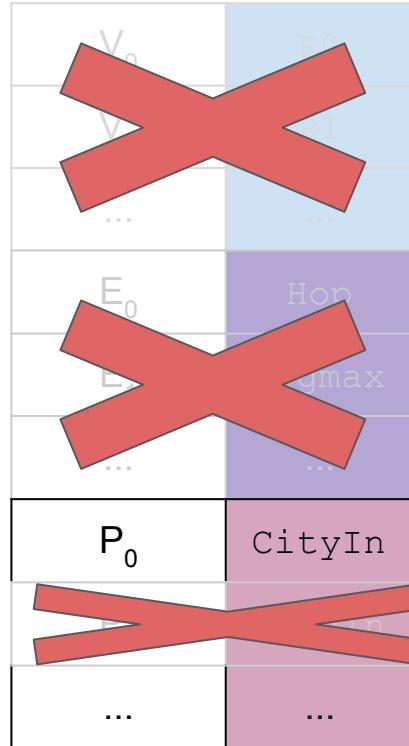


Code Assistance: Semantic Constraint

Given definition of Hop, need to output a predicate that is connected to R2 (m. USA).



Decoder Vocab



Variables: <10

Functions: <10

Predicates: 23K

Valid Predicates:
<100

EVEN IF F IS
DISCONTINUOUS?



U CAN'T BE SERIOUS

Policy Gradient (REINFORCE)

- **MLE** optimizes log likelihood of approximate gold programs

$$J^{ML}(\theta) = \sum_q \log P(a_{0:T}^{best}(q)|q, \theta)$$

- **RL** optimizes the expected reward under a stochastic policy P

$$J^{RL}(\theta) = \sum_q \mathbb{E}_{P(a_{0:T}|q, \theta)}[R(q, a_{0:T})]$$

- The gradient is almost the same as that for MLE except for a **weight** P(R-B)

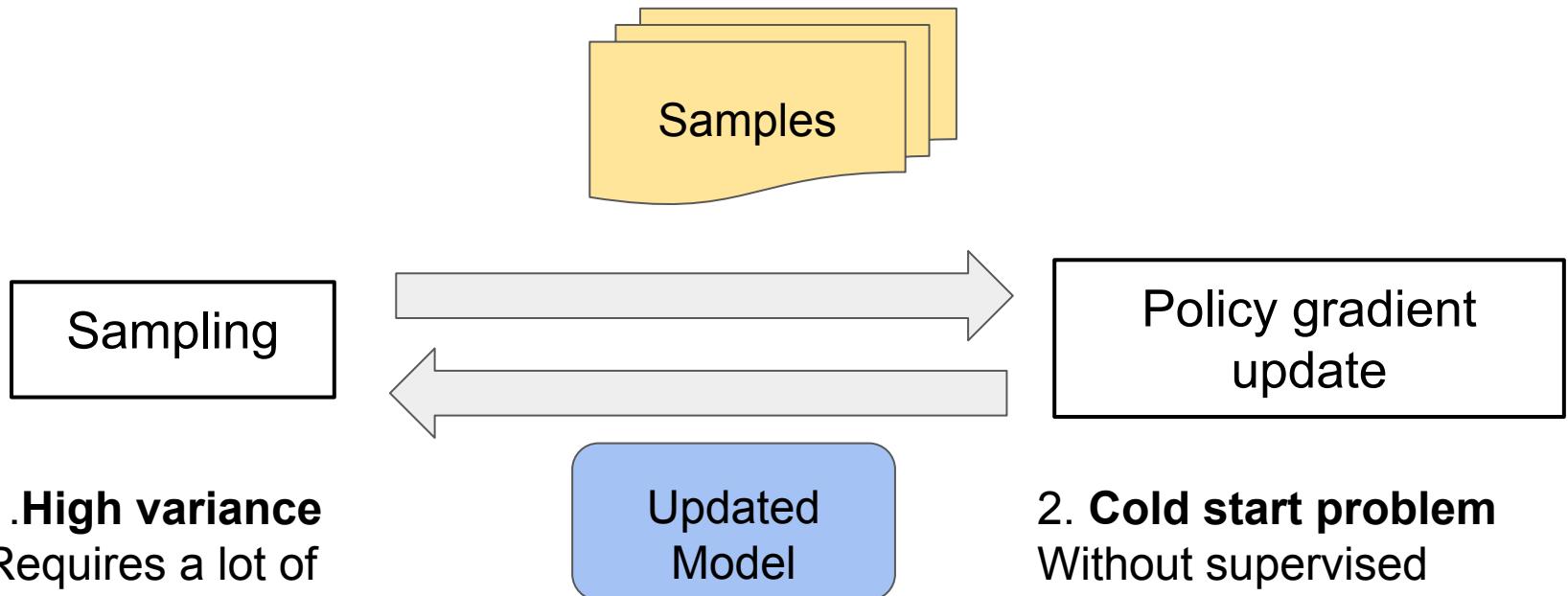
$$\nabla_{\theta} J^{RL}(\theta) = \sum_q \sum_{a_{0:T}} P(a_{0:T}|q, \theta) [R(q, a_{0:T}) - B(q)] \nabla_{\theta} \log P(a_{0:T}|q, \theta)$$

- The **baseline** does not change the solution but improves convergences, e.g.,

$$B(q) = \sum_{a_{0:T}} P(a_{0:T}|q, \theta) R(q, a_{0:T})$$

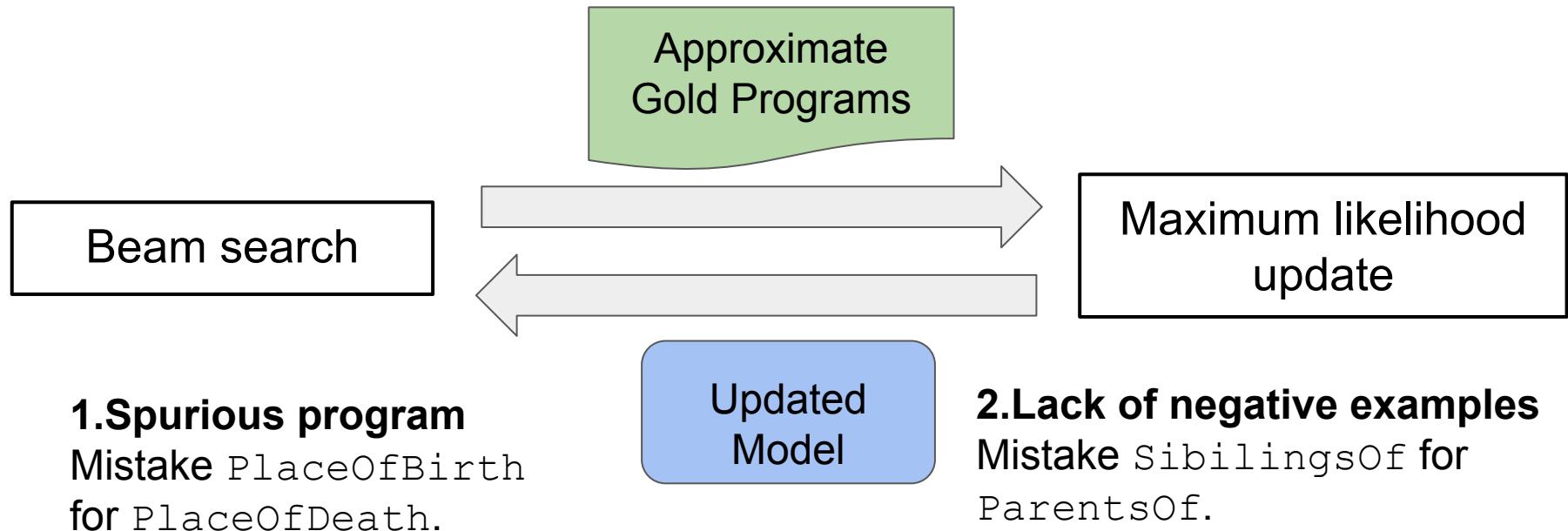
[Williams 1992]

REINFORCE Training



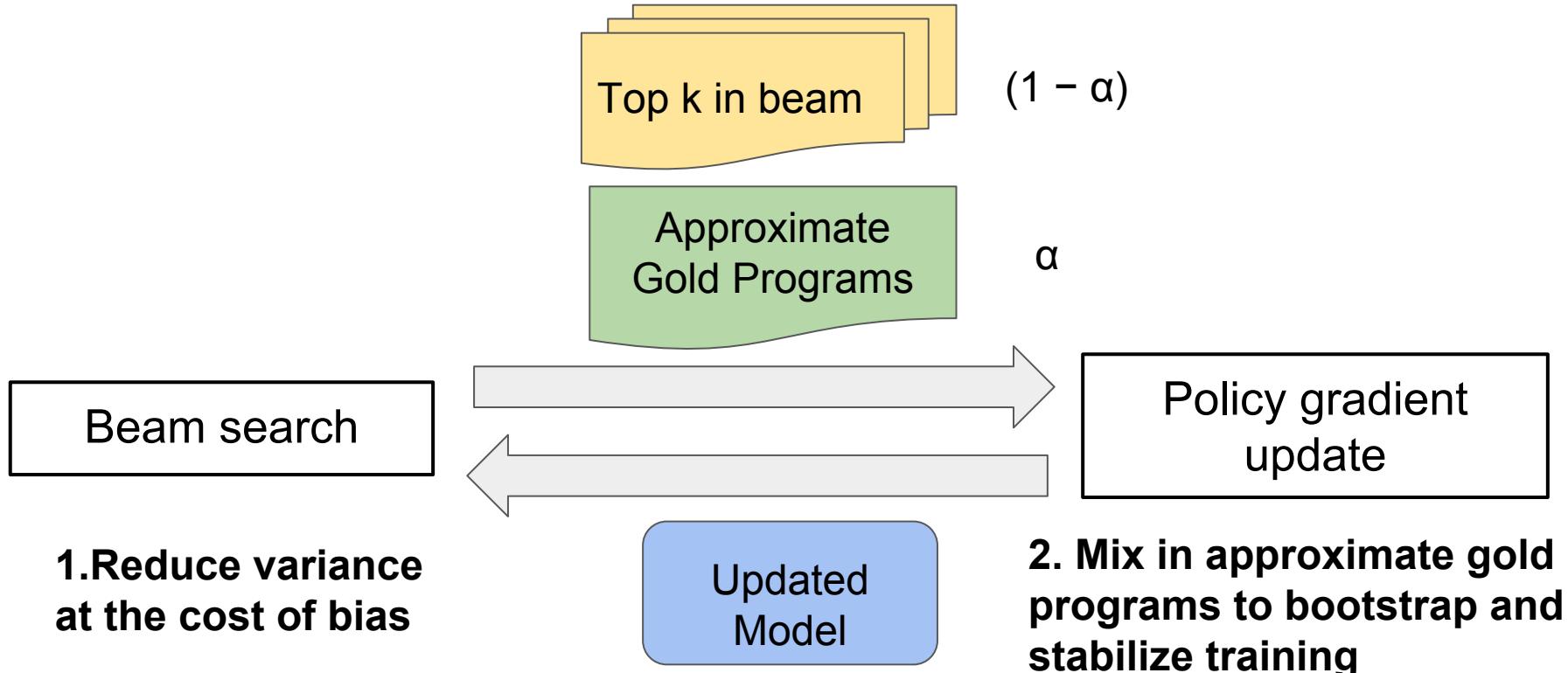
$$\nabla_{\theta} J^{RL}(\theta) = \sum_q \sum_{a_{0:T}} P(a_{0:T}|q, \theta) [R(q, a_{0:T}) - B(q)] \nabla_{\theta} \log P(a_{0:T}|q, \theta)$$

Iterative Maximum Likelihood Training (Hard EM)



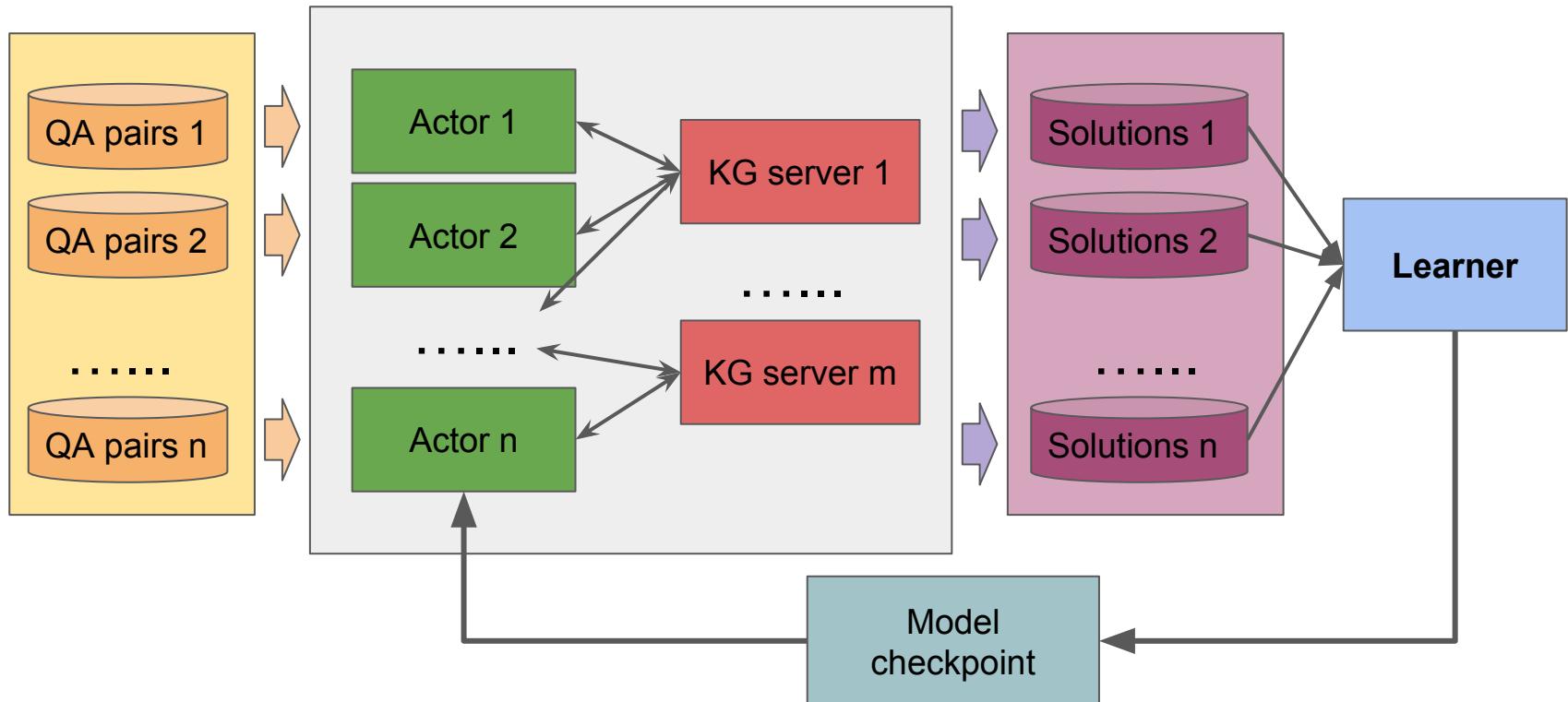
$$J^{ML}(\theta) = \sum_q \log P(a_{0:T}^{best}(q) | q, \theta)$$

Augmented REINFORCE



Distributed Architecture

- 200 actors, 100 KG servers, 1 learner



New State-of-the-Art on *WebQuestionsSP*

- First end-to-end neural network to achieve SOTA on semantic parsing with weak supervision over large knowledge base
- The performance is approaching SOTA with full supervision

Model	Avg. Prec.@1	Avg. Rec.@1	Avg. F1@1	Acc.@1
<i>STAGG</i>	67.3	73.1	66.8	58.8
<i>NSM – our model</i>	70.8	76.0	69.0	59.5
<i>STAGG (full supervision)</i>	70.9	80.3	71.7	63.9

Augmented REINFORCE

- REINFORCE get stuck at local maxima
- Iterative ML training is not directly optimizing the F1 score
- Augmented REINFORCE obtains the best performances

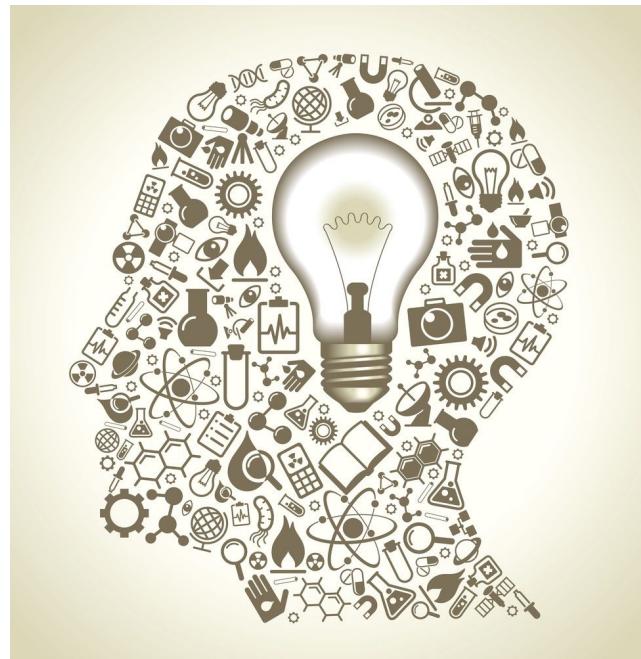
Settings	Train Avg. F1@1	Valid Avg. F1@1
<i>iterative ML only</i>	68.6	60.1
<i>REINFORCE only</i>	55.1	47.8
<i>Augmented REINFORCE</i>	83.0	67.2

Plan

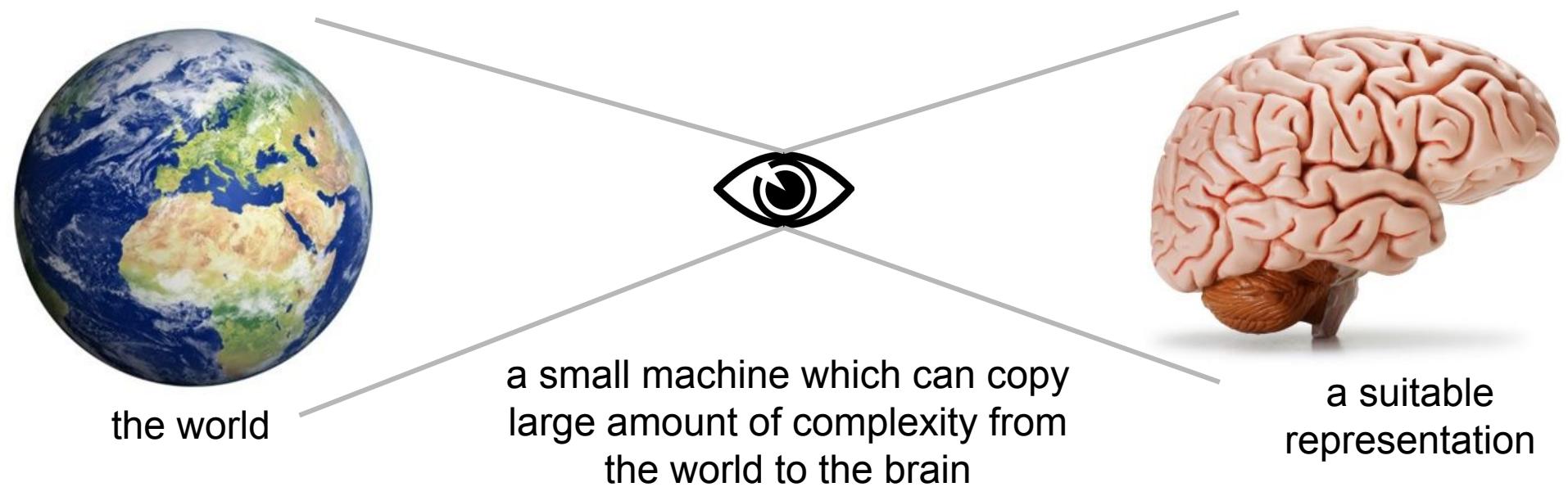
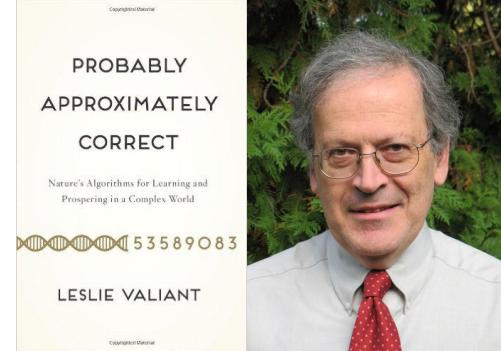
- Language & control
 - *Neural Symbolic Machines*: Semantic Parsing on Freebase with Weak Supervision
 - Chen Liang, Jonathan Berant, Quoc Le, Kenneth Forbus, Ni Lao
- **Knowledge & scalability**
 - Learning to Organize Knowledge with An *N-Gram Machine*
 - Fan Yang, Jiazhong Nie, William Cohen, Ni Lao

Where does knowledge come from?

- The human brain contains roughly 100 billion neurons each capable of making around 1,000 connections
 - Where do we get these 100 TB parameters?
 - How many lines of code do I need to write if I want to achieve AI?



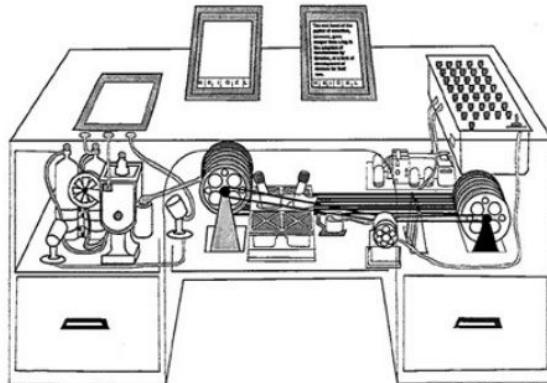
The Mind's Eye



Knowledge and Scalability

- How information should be organized for scalability?

"AS WE MAY THINK" (1945)

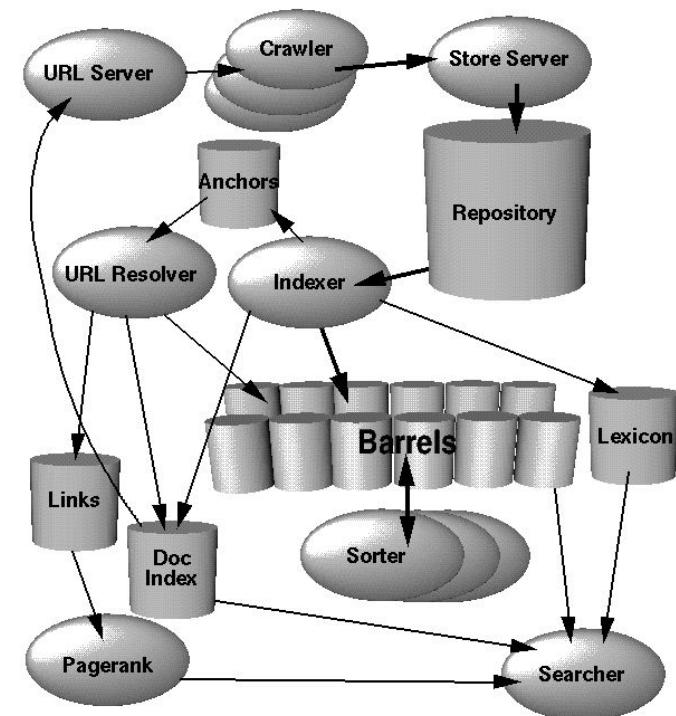


“

Consider a future device for individual use, which is a sort of mechanized private file and library. It needs a name, and to coin one at random, memex will do. A memex is a device in which an individual stores all his books, records, and communications, and which is mechanized so that it may be consulted with exceeding speed and flexibility. It is an enlarged intimate supplement to his memory.

Scalability of modern search engines

- Can respond to user's requests within a fraction of a second
- But are weak at text understanding and complex reasoning



Scalability of mammal memory



- Very **rapid adaptation** (in just one or a few trials) is necessary for survival
 - E.g., associating taste of food and sickness
- Need **fast responses** based on large amount of knowledge
 - Needs good representation of knowledge
- However, **good representation** can only be learnt gradually
 - During sleeps
 - to prevent interference with established associations

[Garcia+ 1966]

[Wickman 2012]

[Bartol+ 2015]

Complementary Learning Theory

Connections within and among neocortical areas (green) support gradual acquisition of structured knowledge through interleaved learning

Encoder

Bidirectional connections (blue) link neocortical representations to the hippocampus/MTL for storage, retrieval, and replay

Record & Replay

Episodic memory

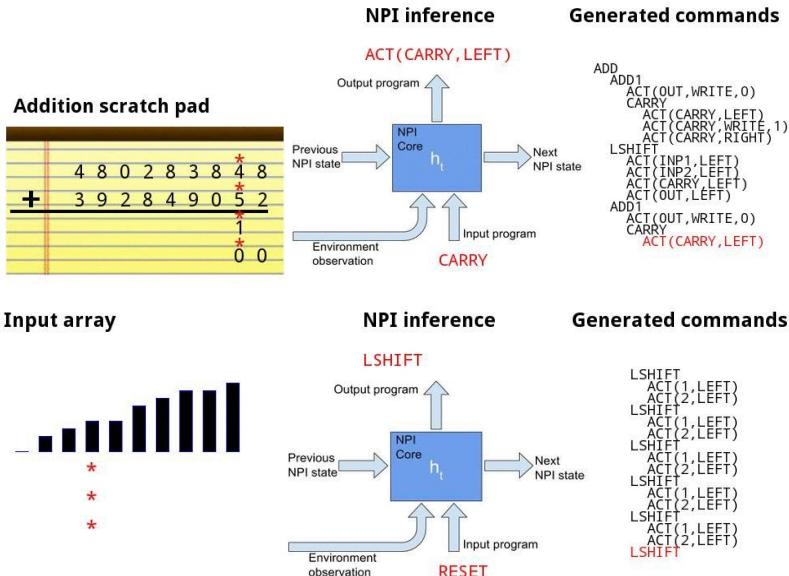
Observations

Rapid learning in connections within hippocampus (red) supports initial learning of arbitrary new information

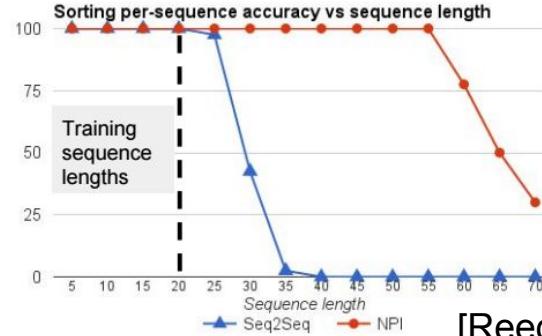
Primary sensory and motor cortices

(Scalable) Neural Nets

- Impressive works to show NN can learn addition and sorting, but...



- The learned operations are not as scalable and precise.



[Reed & Freitas 2015]

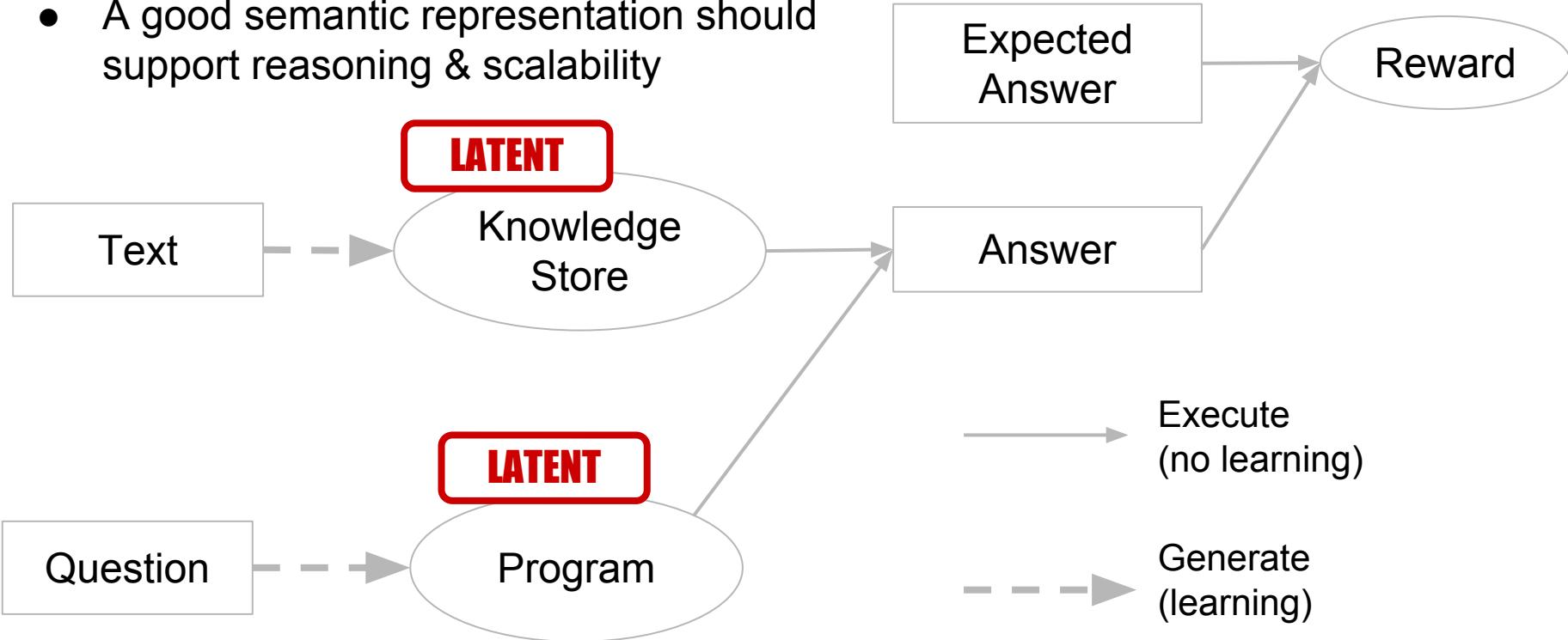
- Why not use existing modules that are scalable, precise and interpretable?



[Zaremba & Sutskever 2016]

Question answering as a simple test bed

- A good semantic representation should support reasoning & scalability



Facebook bAbI Tasks

- Simulated question answering tasks to test the ability to "**understand**"
- We introduce a special version ("life-long bAbI"), which has stories of up to 10 million sentences

Sam walks into the kitchen.	Brian is a lion.	Mary journeyed to the den.
Sam picks up an apple.	Julius is a lion.	Mary went back to the kitchen.
Sam walks into the bedroom.	Julius is white.	John journeyed to the bedroom.
Sam drops the apple.	Bernhard is green.	Mary discarded the milk.
Q: Where is the apple?	Q: What color is Brian?	Q: Where was the milk before the den?
A. Bedroom	A. White	A. Hallway

The Knowledge Storage

- Each tuple in the KG is an n-gram, which is a list of symbols
- Each tuple is also associated with a timestamp to reason over time

Table 1: Example of probabilistic knowledge storage. Each sentence may be converted to a distribution over multiple tuples, but only the one with the highest probability is shown here.

Sentences	Knowledge tuples		
	Time stamp	Symbols	Probability
Mary went to the kitchen.	1	mary to kitchen	0.9
Mary picked up the milk.	2	mary the milk	0.4
John went to the bedroom.	3	john to bedroom	0.7
Mary journeyed to the garden.	4	mary to garden	0.8

The Programs

- A **program** C is a list of **expressions** $c_1 \dots c_N$, where each c_i is either a special expression Return indicating the end of the program, or is of the form $(F, A_1 \dots A_L)$

Table 2: Functions in N-Gram Machines. The knowledge storage on which the programs can execute is Γ , and a knowledge tuple Γ_i is represented as $(i, (\gamma_1, \dots, \gamma_N))$. “FR” means *from right*.

Name	Inputs	Return
Hop	$v_1 \dots v_L$	$\{\gamma_{L+1} \mid \text{if } (\gamma_1 \dots \gamma_L) == (v_1, \dots, v_L), \forall \Gamma \in \Gamma\}$
HopFR	$v_1 \dots v_L$	$\{\gamma_{N-L} \mid \text{if } (\gamma_{N-L+1} \dots \gamma_N) == (v_L, \dots, v_1), \forall \Gamma \in \Gamma\}$
Argmax	$v_1 \dots v_L$	$\text{argmax}_i \{(\gamma_{L+1}, i) \mid \text{if } (\gamma_1 \dots \gamma_L) == (v_1, \dots, v_L), \forall \Gamma_i \in \Gamma\}$
ArgmaxFR	$v_1 \dots v_L$	$\text{argmax}_i \{(\gamma_{N-L}, i) \mid \text{if } (\gamma_{N-L+1} \dots \gamma_N) == (v_L, \dots, v_1), \forall \Gamma_i \in \Gamma\}$

Example KS & Program

Table 6: Task 2 Two Supporting Facts

Story	Knowledge Storage
Sandra journeyed to the hallway.	Sandra journeyed hallway
John journeyed to the bathroom.	John journeyed bathroom
Sandra grabbed the football.	Sandra got football
Daniel travelled to the bedroom.	Daniel journeyed bedroom
John got the milk.	John got milk
John dropped the milk.	John got milk
Question	Program
Where is the milk?	ArgmaxFR milk got Argmax V1 journeyed

Seq2Seq components

- A **knowledge encoder** that converts sentences to knowledge tuples and defines a distribution $P(\Gamma_i | s_i, s_{i-1}; \theta_{\text{enc}})$ s_{i-1} for co-references
 - The probability of a KS $\Gamma = \{\Gamma_1, \dots, \Gamma_n\}$ is the product of its tuples' probabilities:
$$P(\Gamma | \mathbf{s}; \theta_{\text{enc}}) = \prod_{\Gamma_i \in \Gamma} P(\Gamma_i | s_i, s_{i-1}; \theta_{\text{enc}})$$
- A **knowledge decoder** that converts tuples back to sentences and defines a distribution $P(s_i | \Gamma_i, s_{i-1}; \theta_{\text{dec}})$ -- it will enable unsupervised training
- A **programmer** that converts questions to programs and defines a distribution $P(C | q, \Gamma; \theta_{\text{prog}})$ Γ for code assist

Inference

Given an example (s, q, a) from our training set, we would like to maximize the **expected reward**

$$O^{QA}(\theta_{\text{enc}}, \theta_{\text{prog}}) = \sum_{\Gamma} \sum_C P(\Gamma|s; \theta_{\text{enc}}) P(C|q, \Gamma; \theta_{\text{prog}}) R(\Gamma, C, a),$$

For gradient estimations we apply **beam search** instead of **MCMC**, which has huge variances

It leads to a **hard search problem**, which we solve by having

1) a stabilized **auto-encoding** objective to bias the encoder to more interesting hypotheses;

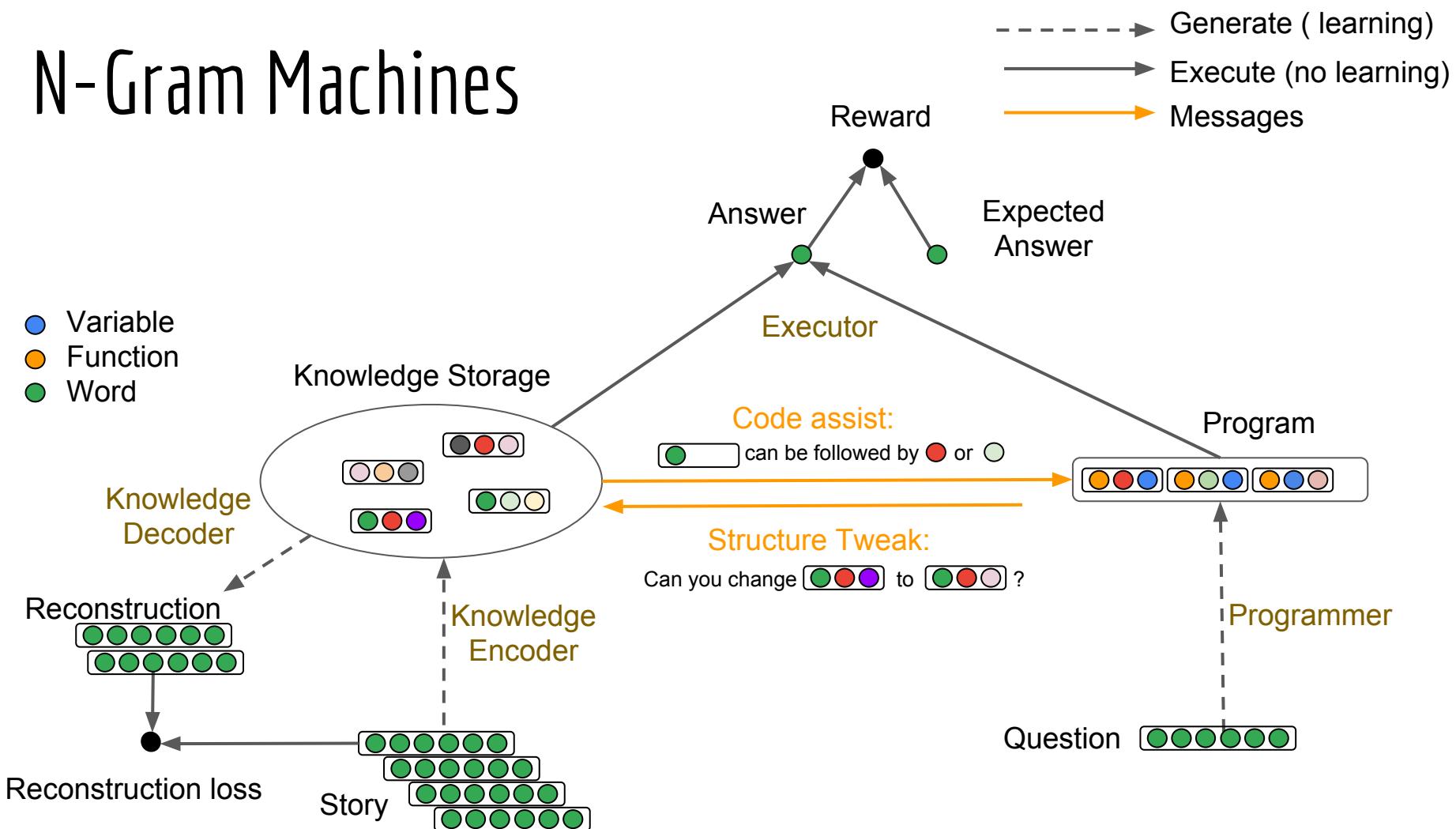
$$O^{\text{AE}}(\theta_{\text{enc}}, \theta_{\text{dec}}) = \mathbb{E}_{p(z|x; \theta_{\text{enc}})} [\log p(x|z; \theta_{\text{dec}})] + \sum_{z \in \mathbf{Z}^N(x)} \log p(x|z; \theta_{\text{dec}}),$$

$\mathbf{Z}^N(x)$: all tuples of length N which only consist of words from x

2) a **structural tweak** procedure which retrospectively corrects the inconsistency among multiple hypotheses so that reward can be achieved

- While code assist uses the **knowledge storage** to inform the **programmer**, structure tweak adjusts the **knowledge encoder** to cooperate with an uninformed **programmer**.

N-Gram Machines



Optimization

- For training stability and tweaking, we augment the training objective with **experience replays**

$$\nabla_{\theta_{\text{dec}}} O'(\theta) = \sum_{s_i \in \mathbf{s}} \sum_{\Gamma_i} [\beta(\Gamma_i) + P(\Gamma_i | s_i, s_{i-1}; \theta_{\text{enc}})] \nabla_{\theta_{\text{dec}}} \log P(s_i | \Gamma, s_{i-1}; \theta_{\text{dec}}),$$

$\beta(\Gamma_i)$ is 1 if Γ_i only contains tokens from s_i and 0 otherwise

$$\nabla_{\theta_{\text{enc}}} O'(\theta) = \sum_{s_i \in \mathbf{s}} \sum_{\Gamma_i} [P(\Gamma_i | s_i, s_{i-1}; \theta_{\text{enc}}) \log P(s_i | \Gamma_i, s_{i-1}; \theta_{\text{dec}}) + \mathcal{R}(\mathcal{G}'(\Gamma_i)) + \mathcal{R}(\mathcal{G}(\Gamma_i))] \nabla_{\theta_{\text{enc}}} \log P(\Gamma_i | s_i, s_{i-1}; \theta_{\text{enc}}),$$

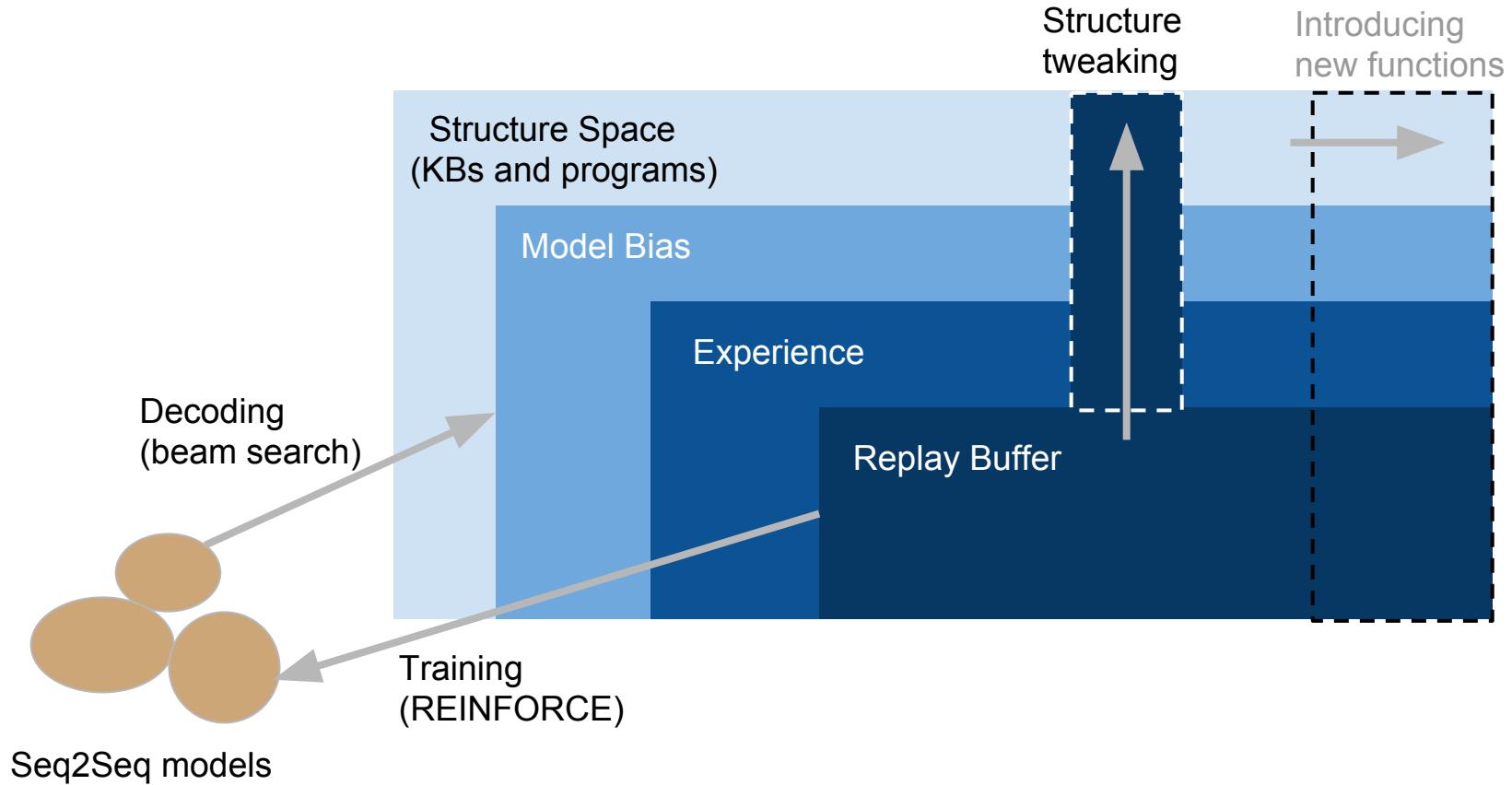
where $\mathcal{R}(\mathcal{G}) = \sum_{\mathbf{\Gamma} \in \mathcal{G}} \sum_C P(\mathbf{\Gamma} | \mathbf{s}; \theta_{\text{enc}}) P(C | q, \mathbf{\Gamma}; \theta_{\text{prog}}) R(\mathbf{\Gamma}, C, a)$ is the total expected reward for a set of valid knowledge stores \mathcal{G} , $\mathcal{G}(\Gamma_i)$ is the set of knowledge stores which contains the tuple Γ_i , and $\mathcal{G}'(\Gamma_i)$ is the set of knowledge stores which contains the tuple Γ_i through tweaking.

$$\nabla_{\theta_{\text{prog}}} O'(\theta) = \sum_{\mathbf{\Gamma}} \sum_C [\alpha I [C \in \mathcal{C}^*(\mathbf{s}, q)] + P(C | q, \mathbf{\Gamma}; \theta_{\text{prog}})] \cdot P(\mathbf{\Gamma} | \mathbf{s}; \theta_{\text{enc}}) R(\mathbf{\Gamma}, C, a) \nabla_{\theta_{\text{prog}}} \log P(C | q, \mathbf{\Gamma}; \theta_{\text{prog}}),$$

where $\mathcal{C}^*(\mathbf{s}, q)$ is the experience replay buffer for (\mathbf{s}, q) . $\alpha = 0.1$ is a constant. During training, the program with the highest weighted reward (i.e. $P(\mathbf{\Gamma} | \mathbf{s}; \theta_{\text{enc}}) R(\mathbf{\Gamma}, C, a)$) is added to the replay buffer.

- optimize by **coordinate ascent** – updating three components in alternation with **REINFORCE**

Learning to search



Results

- The bAbI dataset contains twenty tasks in total. We consider the subset of them that are extractive question answering tasks

Table 3: Test accuracy on bAbI tasks with auto-encoding (AE) and structure tweak (ST)

	Task 1	Task 2	Task 11	Task 15	Task 16
MemN2N	1.000	0.830	0.840	1.000	0.440
QA	0.007	0.027	0.000	0.000	0.098
QA + AE	0.709	0.551	1.000	0.246	1.000
QA + AE + ST	1.000	0.853	1.000	1.000	1.000

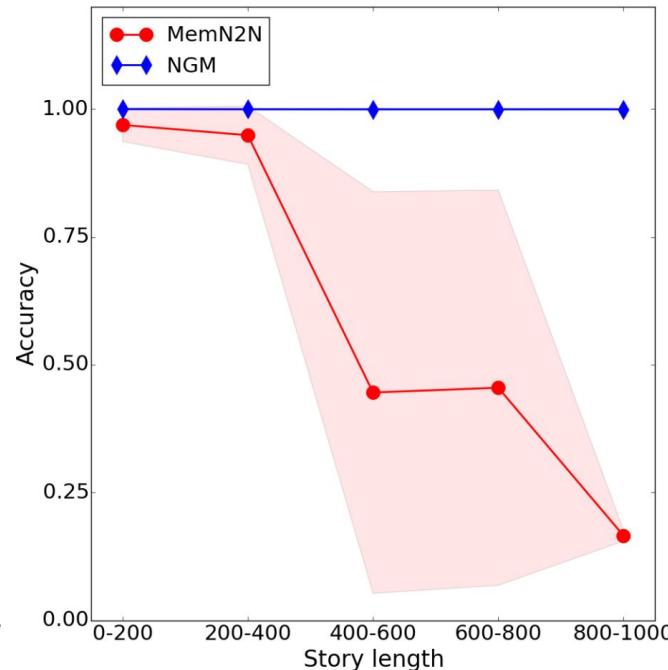
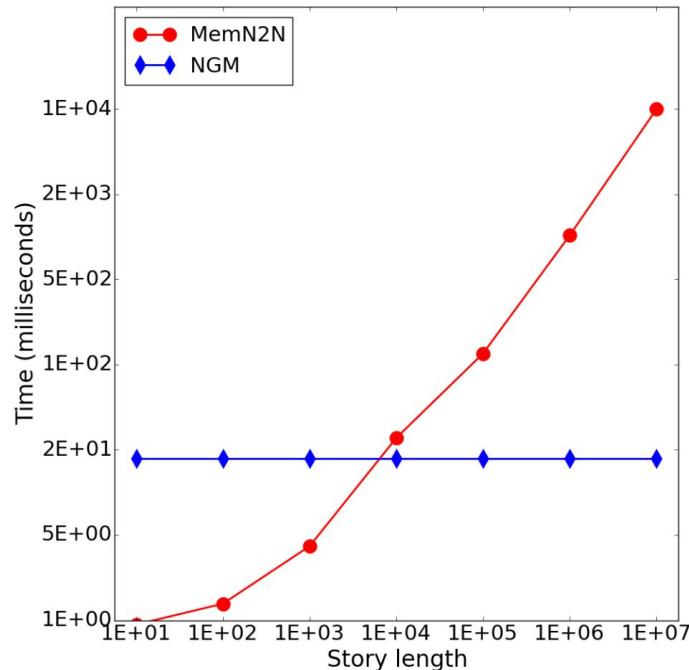
Results

- AE is a strong bias towards good representations
- ST helps to achieve consistency, e.g.,
 - "he" vs "john" (coreference)
 - "cat" vs "cats" (singular vs. plural)
 - "go to" vs "journey to" (synonyms)

QA	QA + AE	QA + AE + ST
went went went	daniel went office	daniel went office
mary mary mary	mary <u>back</u> garden	mary <u>went</u> garden
john john john	john <u>back</u> kitchen	john <u>went</u> kitchen
mary mary mary	mary <u>grabbed</u> football	mary <u>got</u> football
there there there	sandra got apple	sandra got apple
cats cats cats	<u>cats</u> afraid wolves	<u>cat</u> afraid wolves
mice mice mice	<u>mice</u> afraid wolves	<u>mouse</u> afraid wolves
is is cat	gertrude is cat	gertrude is cat

Result

- Scalability



Thanks!

Generated Programs

- **Question:** “what college did russell wilson go to?”
- **Generated program:**

```
(hop v1 /people/person/education)
(hop v2 /education/education/institution)
(filter v3 v0 /common/topic/notable_types )
<EOP>
```

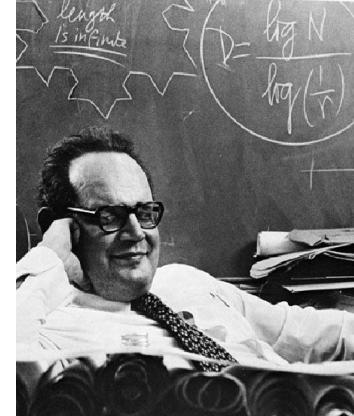
In which

v0 = “College/University” (m.01y2hn1)
v1 = “Russell Wilson” (m.05c10yf)

- **Distribution of the length of generated programs**

#Expressions	0	1	2	3
<i>Percentage</i>	0.4%	62.9%	29.8%	6.9%
<i>F1</i>	0.0	73.5	59.9	70.3

Mandelbrot Set

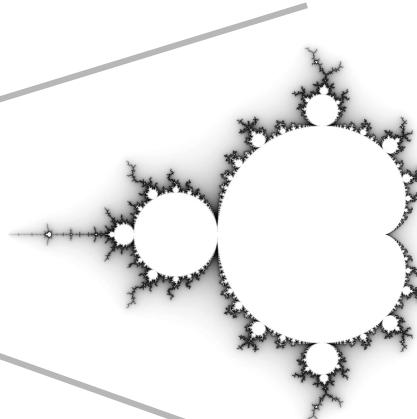


the nature
of complex
numbers



$$z_0 = 0$$

$$z_{n+1} = z_n^2 + c$$



$$c \in M \iff \lim_{n \rightarrow \infty} |z_{n+1}| \leq 2$$

Structure Tweak

- No reward
 - Caused by semantic (i.e. run-time) errors
 - Program: Hop var *cats*
 - Knowledge tuple: *Cat* afraid wolves
 - Propose knowledge tuples that allow the program with high probability to obtain positive reward
- Low expected reward
 - program with high reward in $P(\text{prog} | q, \text{kg}) \neq$ program with high probability in $P(\text{prog} | q)$
 - “Hop var *journeyed*” vs “Hop var *went*”
 - Propose knowledge tuples that allow program with high probability to obtain high rewards

Example KS & Program

Table 8: Task 15 Basic Deduction

Story	Knowledge Storage
Sheep are afraid of cats.	Sheep afraid cats
Cats are afraid of wolves.	Cat afraid wolves
Jessica is a sheep.	Jessica is sheep
Mice are afraid of sheep.	Mouse afraid sheep
Wolves are afraid of mice.	Wolf afraid mice
Emily is a sheep.	Emily is sheep
Winona is a wolf.	Winona is wolf
Gertrude is a mouse.	Gertrude is mouse
Question	Program
What is Emily afraid of?	Hop Emily is Hop V1 afraid

Example KS & Program

Table 9: Task 16 Basic Induction

Story	Knowledge Storage
Berhard is a rhino.	Bernhard a rhino
Lily is a swan.	Lily a swan
Julius is a swan.	Julius a swan
Lily is white.	Lily is white
Greg is a rhino.	Greg a rhino
Julius is white.	Julius is white
Brian is a lion.	Brian a lion
Bernhard is gray.	Bernhard is gray
Brian is yellow.	Brian is yellow
Question	Program
What color is Greg?	Hop Greg a HopFR V1 a Hop V2 is