

Globally Normalized Reader

Jonathan Raiman and John Miller

Task: Extractive QA (SQuAD)

Who was the first to recognize that the Analytical Engine had applications beyond pure calculation?

page KH FRS (/ˈbæbɪdʒ/; 26 December 1791 – 18 October 1871) was an English polymath.[1] A mathematician, philosopher, inventor, and mechanical engineer, Babbage originated the concept of a programmable computer.[2] He is often credited by some to be a "father of the computer",[2][3][4][5] and is credited with inventing the first mechanical computer. His work led to more complex electronic designs, though all the basic ideas of modern computers are to be found in his designs for the Analytical Engine.[2][6] His varied work in other fields has led to him being described as "pre-eminent" among the many great scientists of his century.[1] Babbage's incomplete mechanisms are on display in the Babbage Institute in London. In 1991, a functioning difference engine No. 2 was constructed from Babbage's original plans. Built to tolerances of the 19th century, the success of the finished engine demonstrated that Babbage's machine would have worked. Babbage's birthplace is disputed, but according to the Oxford Dictionary of National Biography he was most likely born at 44 Waterhouse Alley, Walworth Road, London, England.[7] A blue plaque on the wall of Larcom Street and Walworth Road commemorates Babbage's birth. Babbage's birth was given in his obituary in The Times as 26 December 1792; but then a nephew wrote to say that Babbage was born much earlier, in 1791. The parish register of St. Mary's Church, London, shows that Babbage was baptised on 6 January 1791, suggesting a birth year of 1791.[9][10][11]

Ada Lovelace (née Byron; 10 December 1815 – 27 November 1852) was an English mathematician and writer, chiefly known for her work on Charles Babbage's proposed mechanical general-purpose computer, the Analytical Engine. She was the first to recognise that the machine had applications beyond pure calculation, and created the first algorithm intended to be carried out by such a machine. As a result, she is often regarded as the first to recognise the full potential of a "computing machine" and the first computer programmer. Ada Lovelace was the only legitimate child of the poet Lord Byron, and his wife Anne Isabella Milbanke ("Annabella"), Lady Wentworth. All of Byron's other children were born out of wedlock to other women. Byron separated from his wife a month after Ada was born and left England forever four months later, eventually dying of disease in the Greek War of Independence when Ada was eight years old. Her mother remained bitter towards Lord Byron and promoted Ada's interest in mathematics and logic in an effort to prevent her from developing what she saw as the insanity seen in her father, but Ada remained interested in him despite this (and was, upon her eventual death, buried next to him at her request). Often ill, she spent most of her childhood sick. Ada married William King in 1835. King was made Earl of Lovelace in 1838, and she became Countess of Lovelace. Her educational and social exploits brought her into contact with scientists such as Andrew Crosse, Sir David Brewster, Charles Wheatstone, Michael Faraday and the author Charles Dickens, which she used to further her education. Ada described her approach as...

A computer is a device that can be instructed to carry out sequences of arithmetic or logical operations automatically. The ability of computers to follow generalized sets of operational instructions called programs, enables them to perform an extremely wide variety of tasks. Such computers are used as control systems for a very wide range of industrial and consumer devices. This includes simple devices like microwave ovens and remote controls, as well as more complex devices such as industrial robots and computer assisted manufacturing systems. Computers are also used in general purpose devices like personal computers, smartphones, and servers. The Internet is run on computers, and it connects millions of other computers. Since ancient times, simple manual devices like the abacus have been used by people in doing calculations. Early in the Industrial Revolution, mechanical devices were built to automate long tedious tasks, such as guiding patterns for looms. More sophisticated computers like machines did specialized analog calculations in the 19th century. The first digital electronic calculating machine was developed during World War II. The speed, power, and reliability of computers has increased continuously and dramatically over time. Conventionally, a modern computer consists of a central processing element, typically a central processing unit (CPU), and some form of memory. The processing element performs arithmetic and logical operations, and a sequencing and control element can change the order of operations in response to instructions or information. Peripheral devices include input devices (such as mice, joysticks, etc.), output devices (monitor screens, printers, etc.), and communication devices (network cards, modems, etc.).

~300 words

Goal

Who was the first to recognize that the Analytical Engine had applications beyond pure calculation?



millions of documents

Charles Babbage KH FRS (/ˈbæbɪdʒ/; 26 December 1791 – 18 October 1871) was an English polymath.[1] A mathematician, philosopher, inventor, and mechanical engineer, Babbage originated the concept of a programmable computer.[2] He is often credited by some to be a "father of the computer",[2][3][4][5] and is also credited with inventing the first mechanical computer. His work eventually led to more complex electronic designs, though all the basic ideas of modern computers are to be found in his designs for the analytical engine.[2][6] His varied work in other fields has earned him the description as "pre-eminent" among the many great minds of his century.[1] Babbage's incomplete mechanisms are on display in the Science Museum in London. In 1991, a functioning difference engine No. 2 was constructed from Babbage's original plans. Built to tolerances of the 19th century, the success of the finished engine demonstrated that Babbage's machine would have worked.

Ada Lovelace (née Byron; 10 December 1815 – 27 November 1852) was an English mathematician and writer, chiefly known for her work on Charles Babbage's proposed mechanical general-purpose computer, the Analytical Engine. She was the first to recognise that the machine had applications beyond pure calculation, and created the first algorithm intended to be carried out by such a machine. As a result, she is often regarded as the first to recognise the full potential of a "computing machine" and the first computer programmer. Ada Lovelace was the only legitimate child of the poet Lord Byron, and his wife Anne Isabella Milbanke ("Annabella"), Lady Wentworth. All of Byron's other children were born out of wedlock to other women. Byron separated from his wife a month after Ada was born and left England forever four months later, eventually dying of disease in the Greek War of Independence when Ada was eight years old. Her mother remained bitter towards Lord Byron and promoted

A computer is a device that can be instructed to carry out sequences of arithmetic or logical operations automatically. The ability of computers to follow generalized sets of operations, called programs, enables them to perform an extremely wide variety of tasks. Such computers are used as control systems for a very wide range of industrial and consumer devices. This includes simple devices like microwave ovens and remote controls, as well as more complex devices such as industrial robots and computer assisted medical equipment. Computers are also used in general purpose devices like personal computers and mobile devices such as smartphones. The Internet is run on computers, and it connects millions of other computers. Since ancient times, simple manual devices like the abacus have helped people in doing calculations. Early in the Industrial Revolution, mechanical devices were built to automate long tedious tasks, such as guiding patterns for looms. More sophisticated

How do we get there?

Related Work:

- Bi-Attention Flow (Seo et al., 2016)
- Dynamic Co-Attention Network (Xiong et al., 2016)
- R-Net (Wang et al., 2017)
- Rasor (Lee et al., 2016)
- Hybrid AoA Reader (2018)

Challenges:

- Bi-directional attention
- Rank all possible spans
- No available data augmentation

Globally Normalized Reader

- Factorize search into sentences, span start & end
- Globally Normalize search (Andor et al. @ ACL 2016)
- Beam Search during training w/. Early Updates

Contributions:

- **Conditional computation** (allocate computation to promising search beams)
- **Quasi-infinite training data** w/. Type Swaps
- **24.7x** speedup over bi-attention-flow
- Dev **68.4 EM, 76.21 F1** w/o bidirectional attention

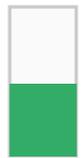
Example



Ada Lovelace was known for her work on Charles Babbage's Analytical Engine. She was the first to recognize that the machine had applications beyond calculation. As a result, she is often regarded as the first to recognise the full potential of a "computing machine" and the first computer programmer.

Who was first to recognize that the Analytical Engine had applications beyond pure calculation?

0.48



Ada Lovelace was known for her work on Charles Babbage's Analytical Engine.

0.51



She was the first to recognize that the machine had applications beyond calculation.

0.01



As a result, she is often regarded as the first to recognise the full potential of a "computing machine" and the first computer programmer.

Probability

Pick a Sentence

Who was first to recognize that the Analytical Engine had applications beyond pure calculation?

0.49



Ada Lovelace was known for her work on Charles Babbage's Analytical Engine.

0.51



She was the first to recognize that the machine had applications beyond calculation.

0.55



Ada Lovelace was known for her work on Charles Babbage's Analytical Engine.

0.09



Lovelace was known for her work on Charles Babbage's Analytical Engine.

0.36



Charles Babbage's Analytical Engine.

Start word chosen for each sentence

Who was first to recognize that the Analytical Engine had applications beyond pure calculation?

0.49



0.51



Ada Lovelace was known for her work on Charles Babbage's Analytical Engine.

She was the first to recognize that the machine had applications beyond calculation.

0.55



0.09



0.36



Ada Lovelace was known for her work on Charles Babbage's Analytical Engine.

Lovelace was known for her work on Charles Babbage's Analytical Engine.

Charles Babbage's Analytical Engine.

0.64



0.20



0.16



Ada Lovelace

Charles Babbage

Charles Babbage's Analytical Engine

Select end word among remainder

Who was first to recognize that the Analytical Engine had applications beyond pure calculation?

0.49



Ada Lovelace was known for her work on Charles Babbage's Analytical Engine.

0.51



She was the first to recognize that the machine had applications beyond calculation.

0.55



Ada Lovelace was known for her work on Charles Babbage's Analytical Engine.

0.09



Lovelace was known for her work on Charles Babbage's Analytical Engine.

0.36



Charles Babbage's Analytical Engine.

0.64



Ada Lovelace

0.20



Charles Babbage

0.16



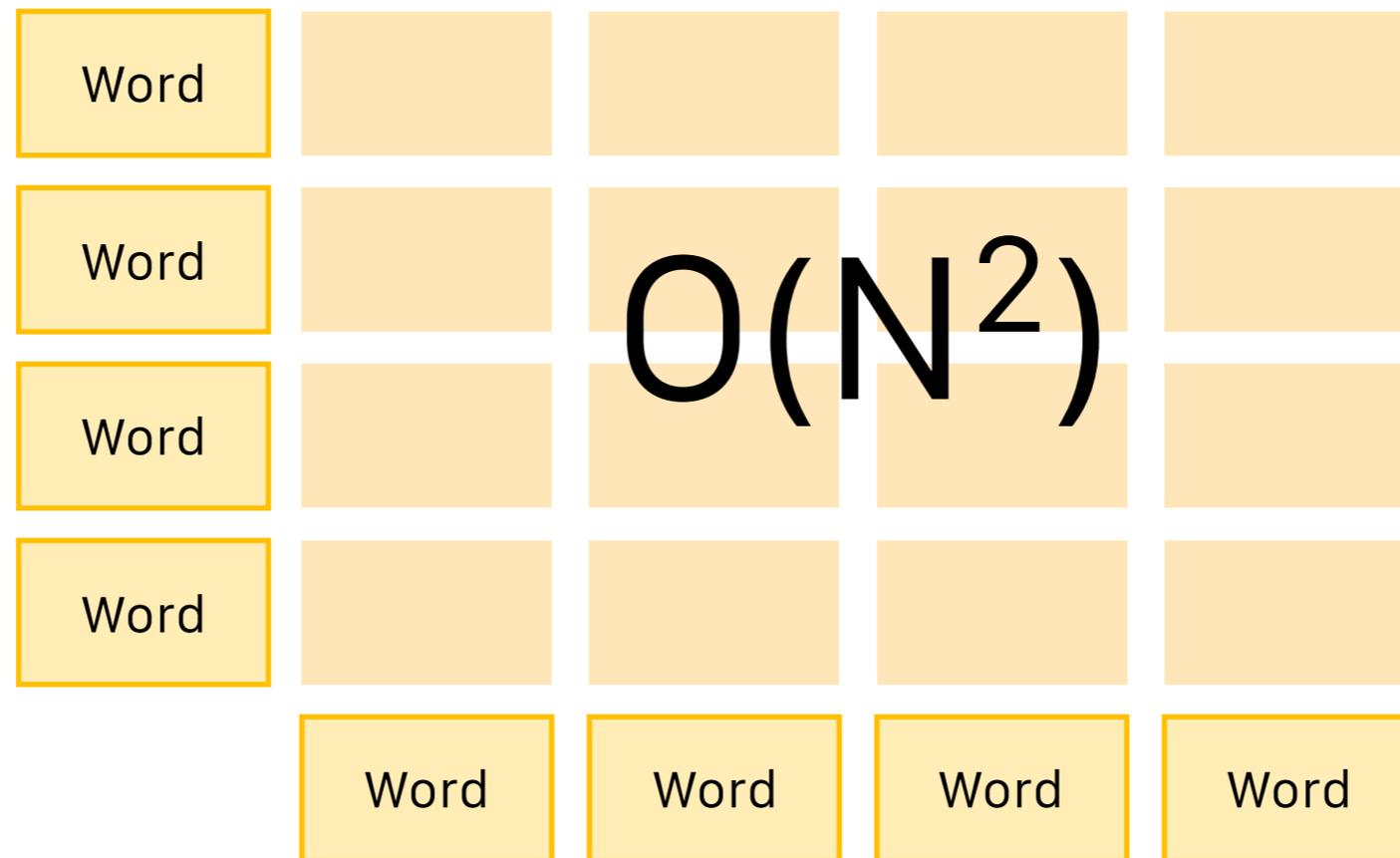
Charles Babbage's Analytical Engine

GNR's Answer: Ada Lovelace

Outline

1. Approach
 - 1) Challenges
 - 2) Architecture
 - 3) Early updates
 - 4) Conditional Computation
 - 5) Global Normalization
 - 6) Type Swaps
2. Results
 - 1) Comparison
 - 2) Data Augmentation
 - 3) Speedup

Span Selection/Attention



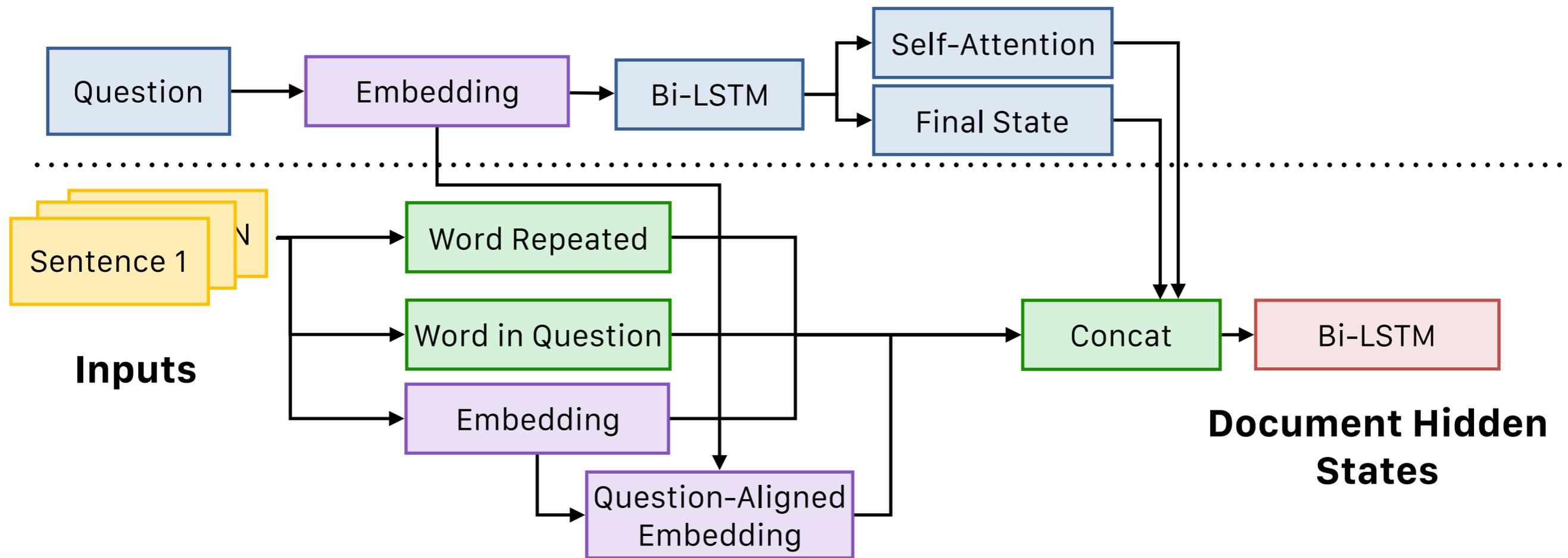
Overfitting

- 100k QA Pairs
- Dropout
- Weight Decay
- Tuning
- Pretrained word vectors
- Ensembles
- Label bias

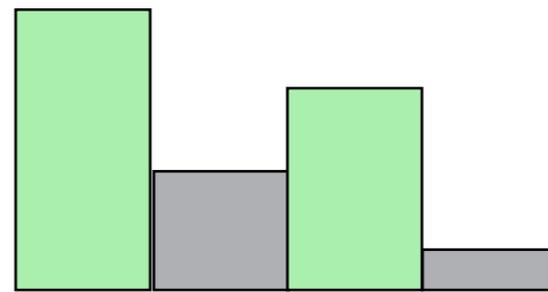
Approach

- Search to shrink candidate space & scope of attention
- Data augmentation

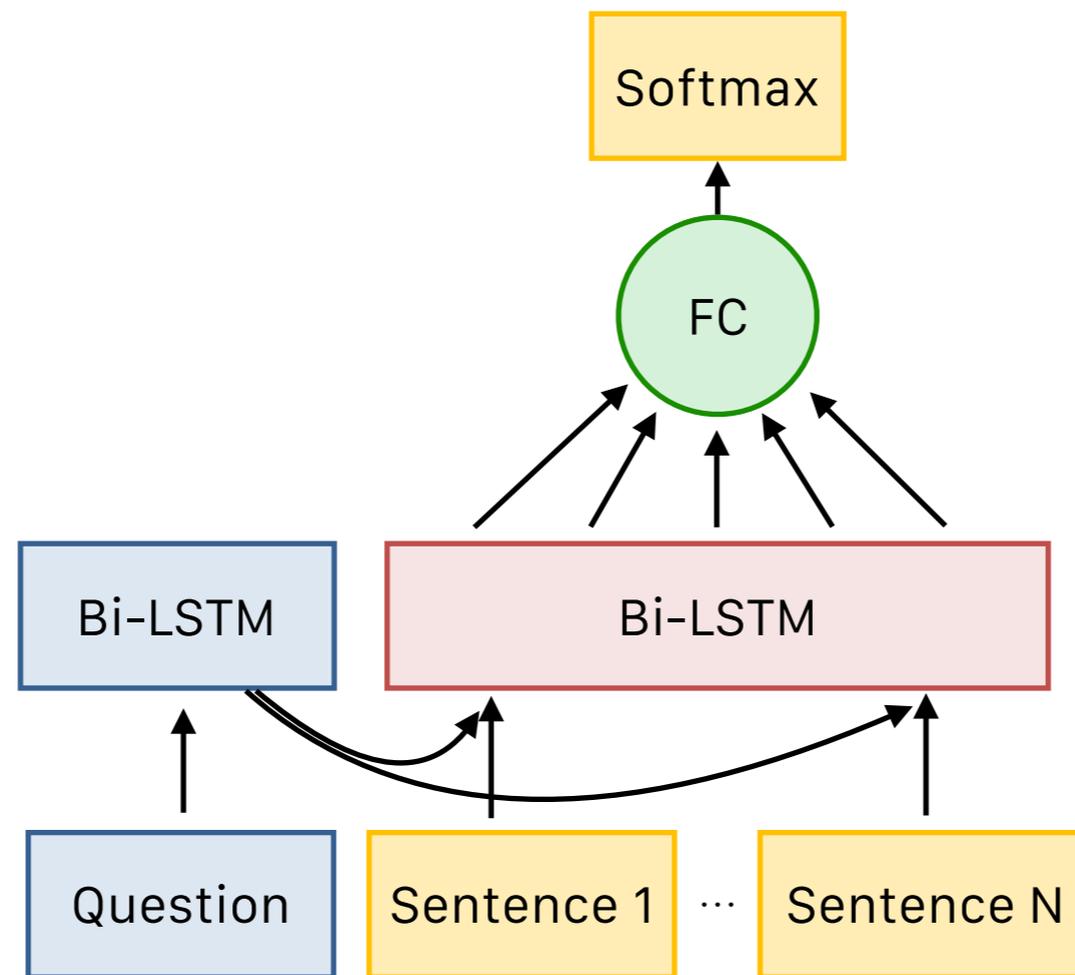
Architecture: Question-Aware Document Encoding



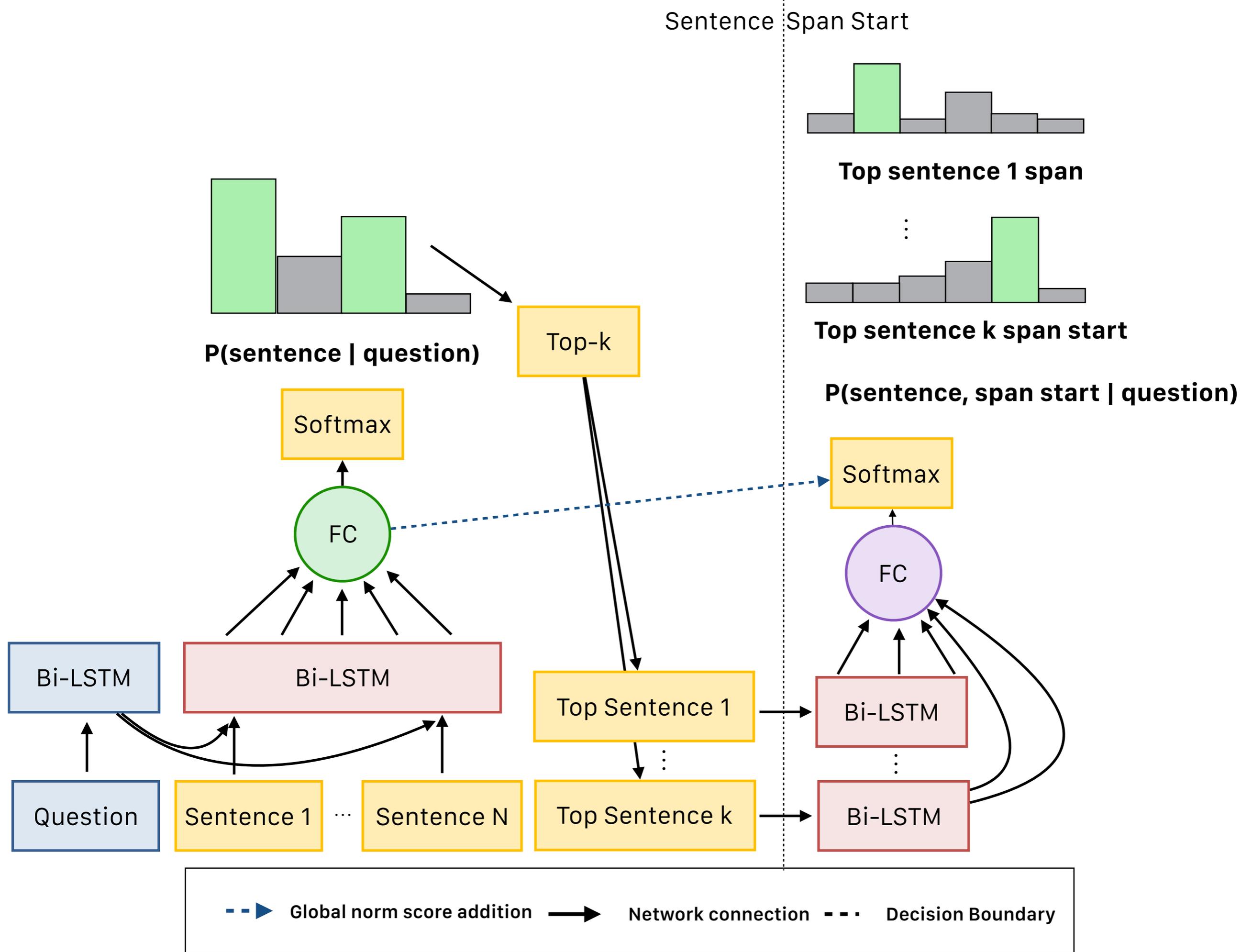
Architecture

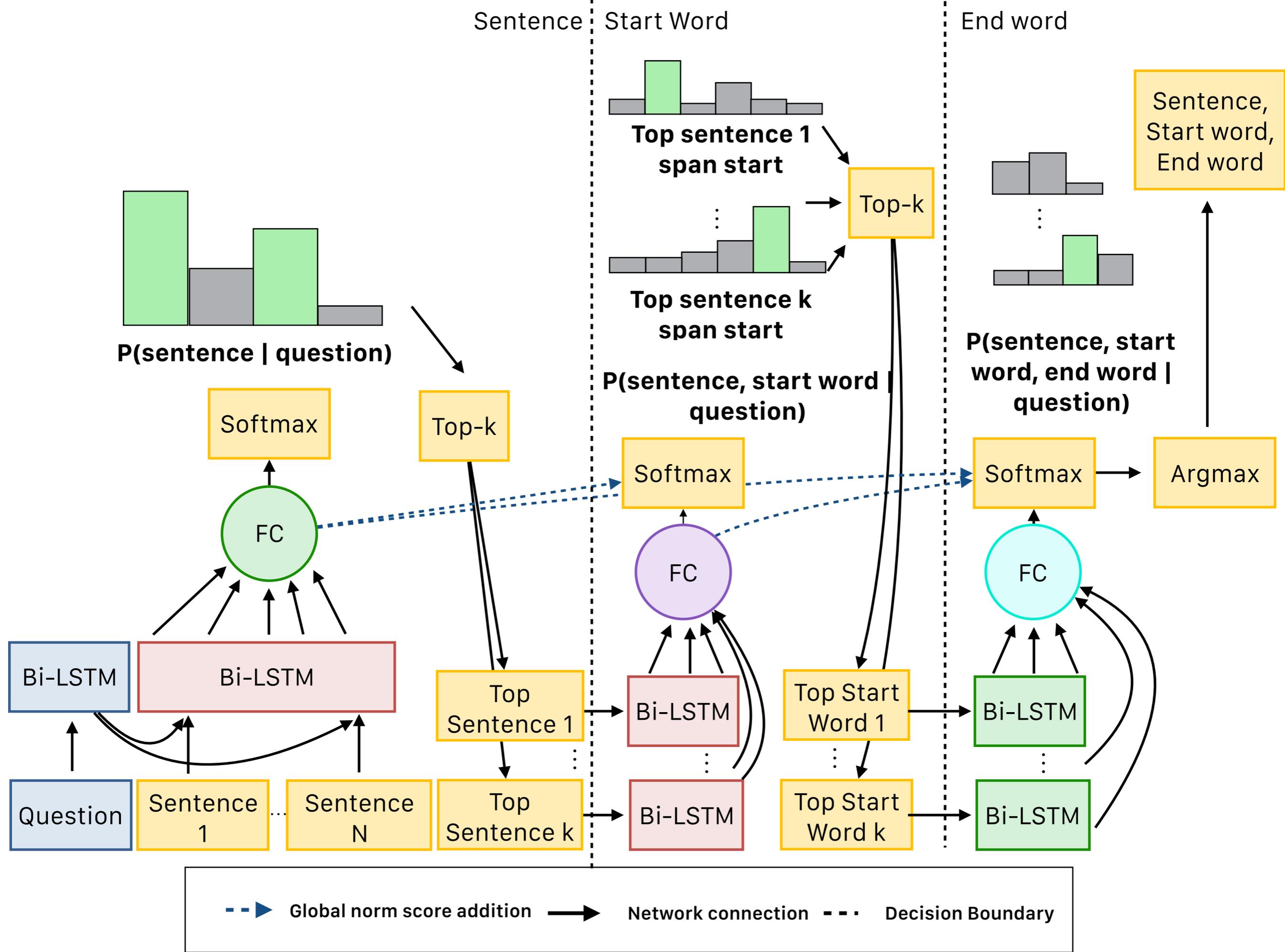


$P(\text{sentence} \mid \text{question})$



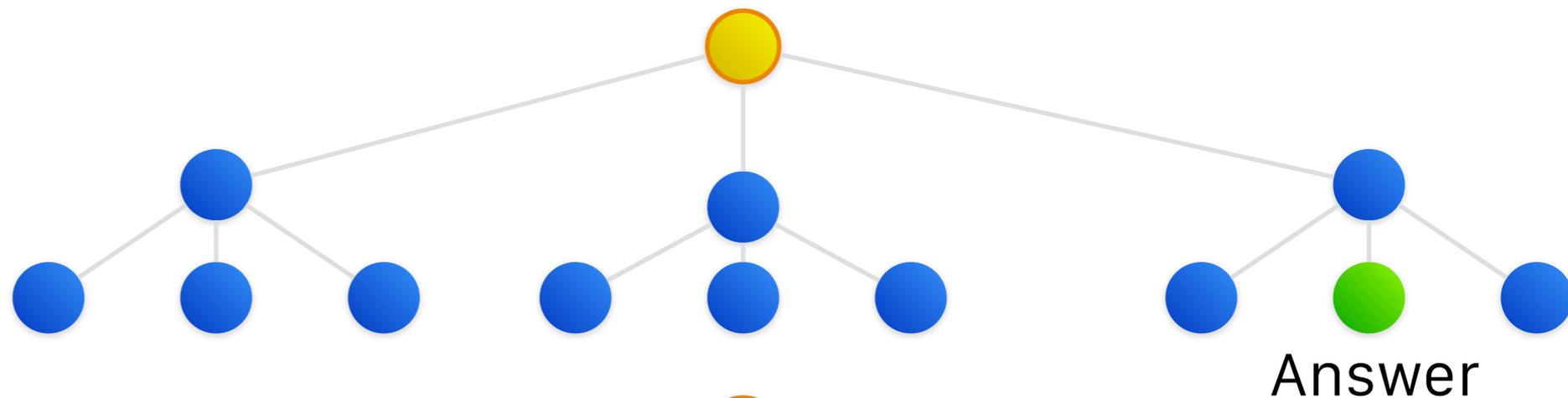
 Global norm score addition  Network connection  Decision Boundary



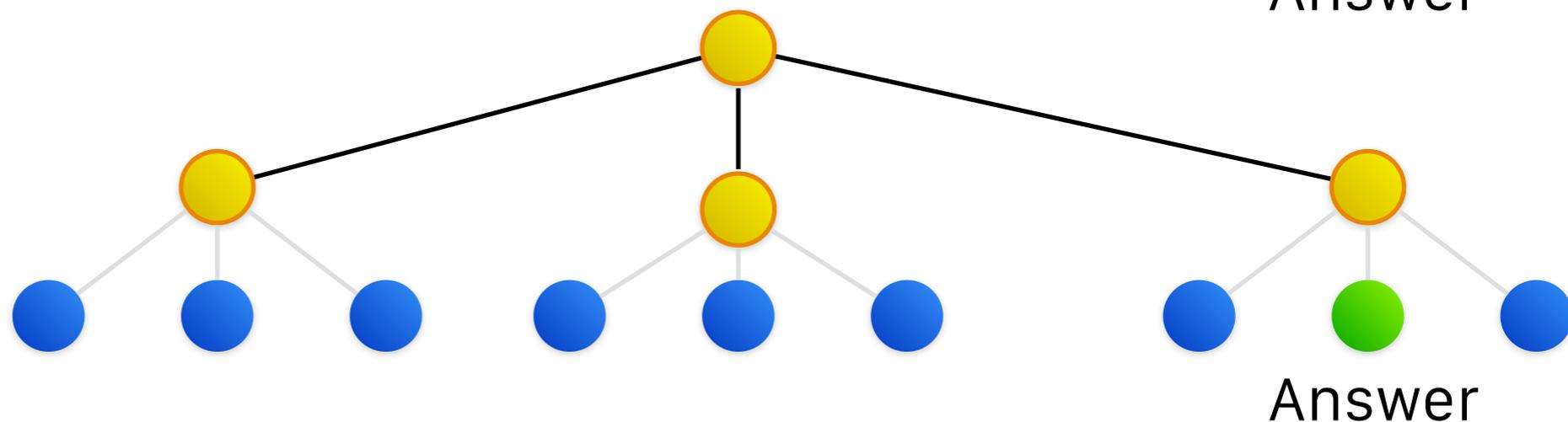


Early Updates

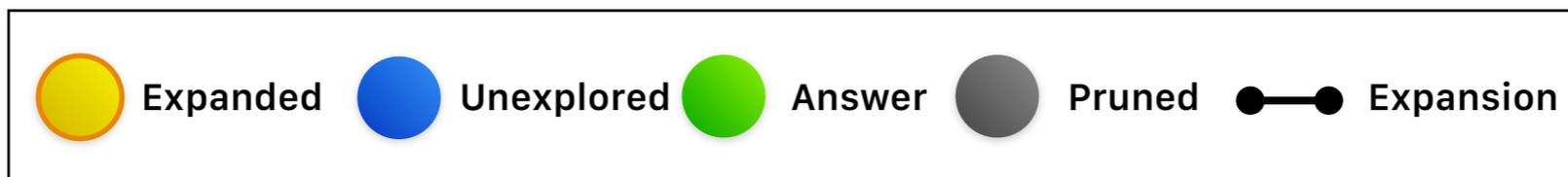
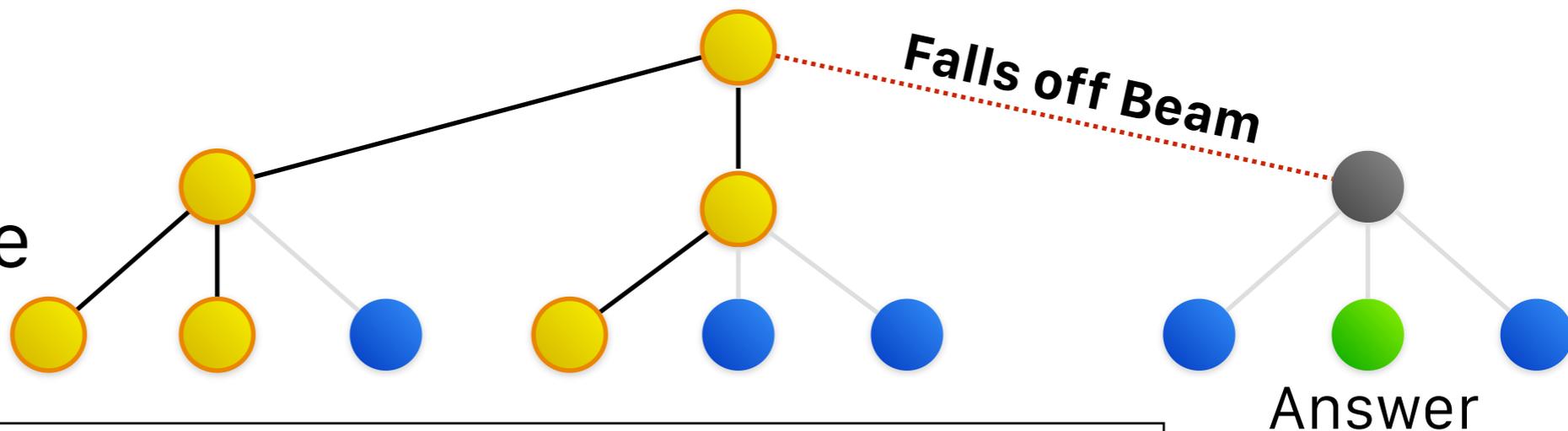
1) Begin search



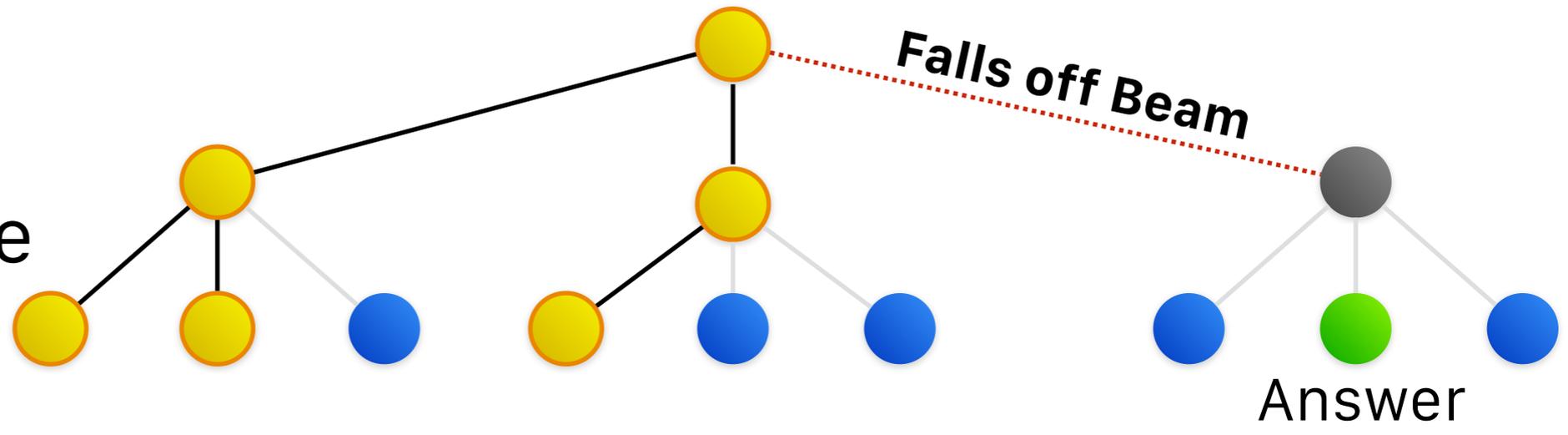
2) Expand search nodes



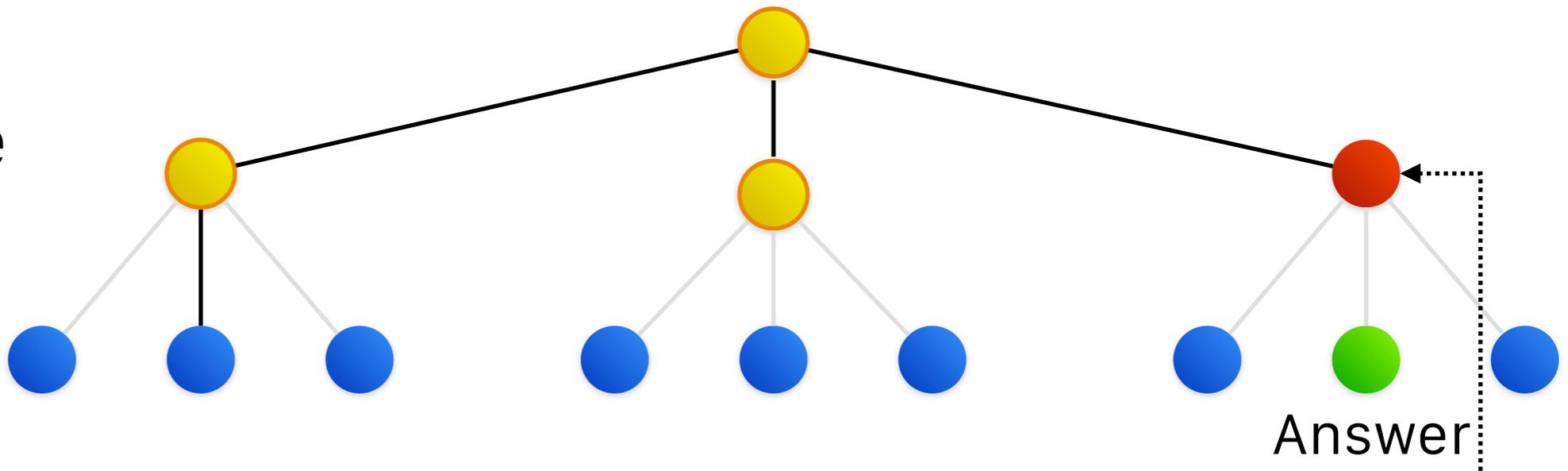
3) Expand & Prune



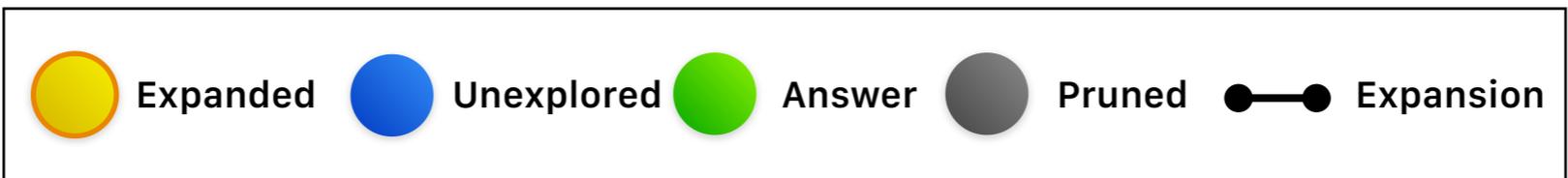
3) Expand & Prune



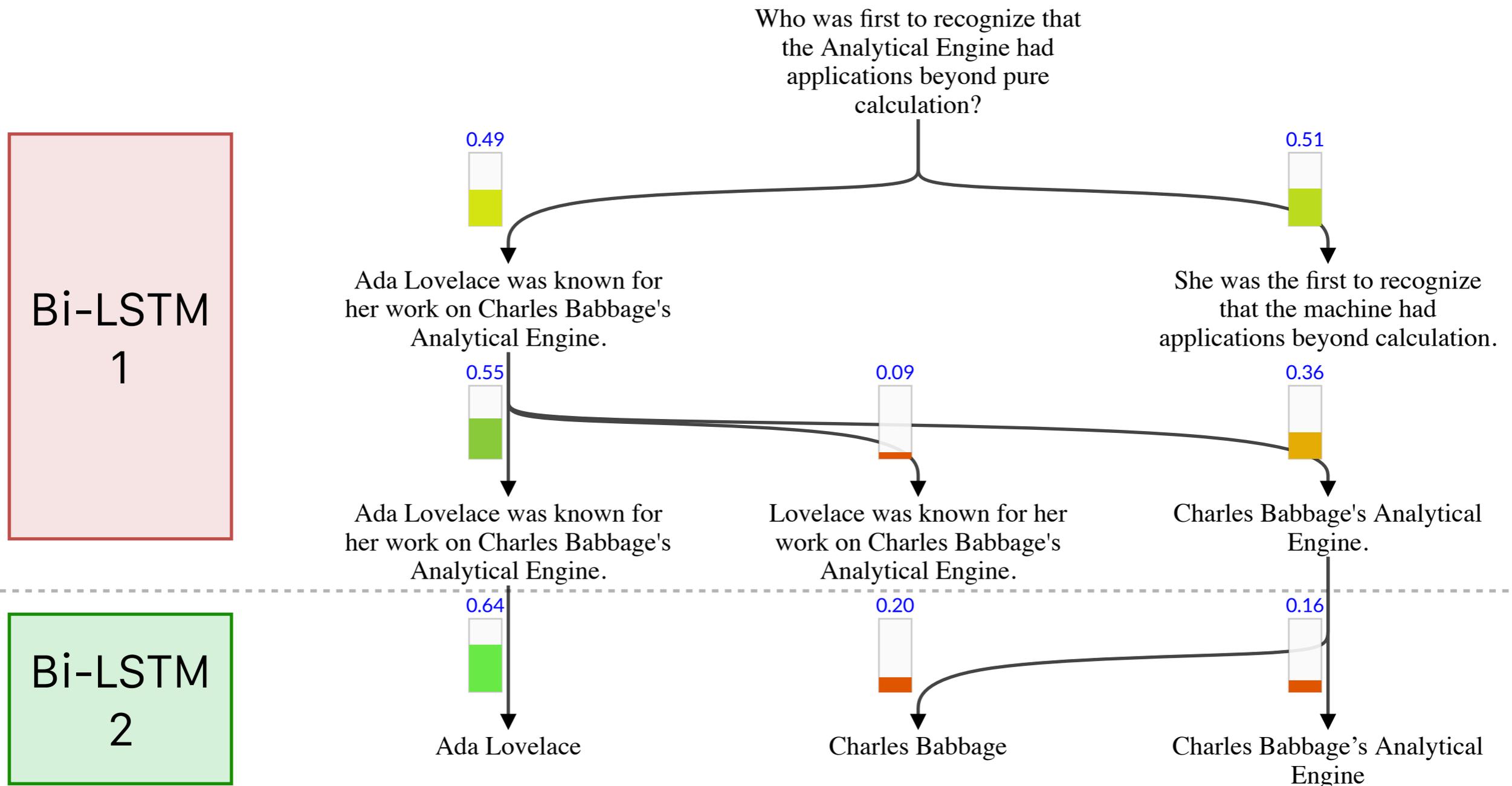
4) Early Update



Maximize Probability to avoid falling off



Conditional Computation



**Sentence prediction accuracy
88-89%: can focus computation on
subset!**

Local Normalization

$$\begin{aligned}\mathbb{P}(a|d, q) &= \mathbb{P}_{\text{sent}}(i|d, q) \cdot \mathbb{P}_{\text{sw}}(j|i, d, q) \cdot \mathbb{P}_{\text{ew}}(k|j, i, d, q) \\ &= \frac{\exp(\phi_{\text{sent}}(i, d, q))}{Z_{\text{sent}}(d, q)} \cdot \frac{\exp(\phi_{\text{sw}}(j, i, d, q))}{Z_{\text{sw}}(i, d, q)} \cdot \frac{\exp(\phi_{\text{ew}}(k, j, i, d, q))}{Z_{\text{ew}}(j, i, d, q)}\end{aligned}$$

a = answer

d = document

q = question

i = sentence

j = start word

k = end word

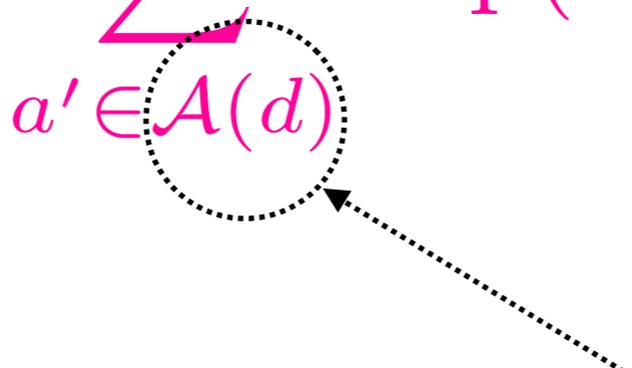
$\phi(\cdot)$ = Score function

Z = Partition function

Global Normalization

$$\text{score}(a, d, q) = \phi_{\text{sent}}(d_i) + \phi_{\text{sw}}(d_{i,j}) + \phi_{\text{ew}}(d_{i,j:k})$$

$$\mathbb{P}(a \mid d, q) = \frac{\exp(\text{score}(a, d, q))}{Z}$$

$$Z = \sum_{a' \in \mathcal{A}(d)} \exp(\text{score}(a', d, q))$$


a = answer

d = document

q = question

i = sentence

j = start word

k = end word

$\phi(\cdot)$ = Score function

Z = Partition function

Set grows exponentially.

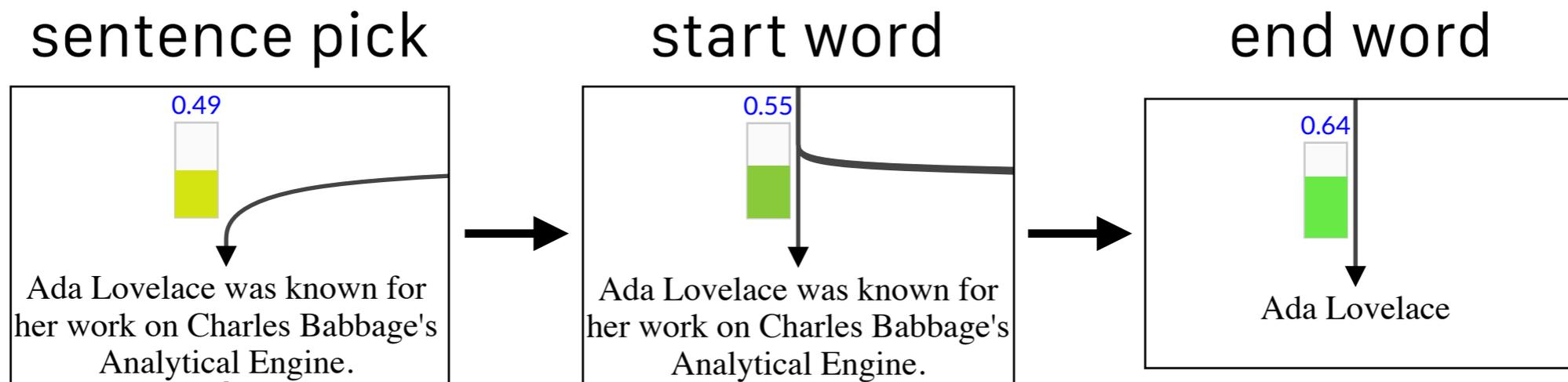
Approximate using beam search

$\mathcal{A}(d)$ = Set of all possible answer spans

$d_{i,j:k}$ = span from word j to k , in sentence i

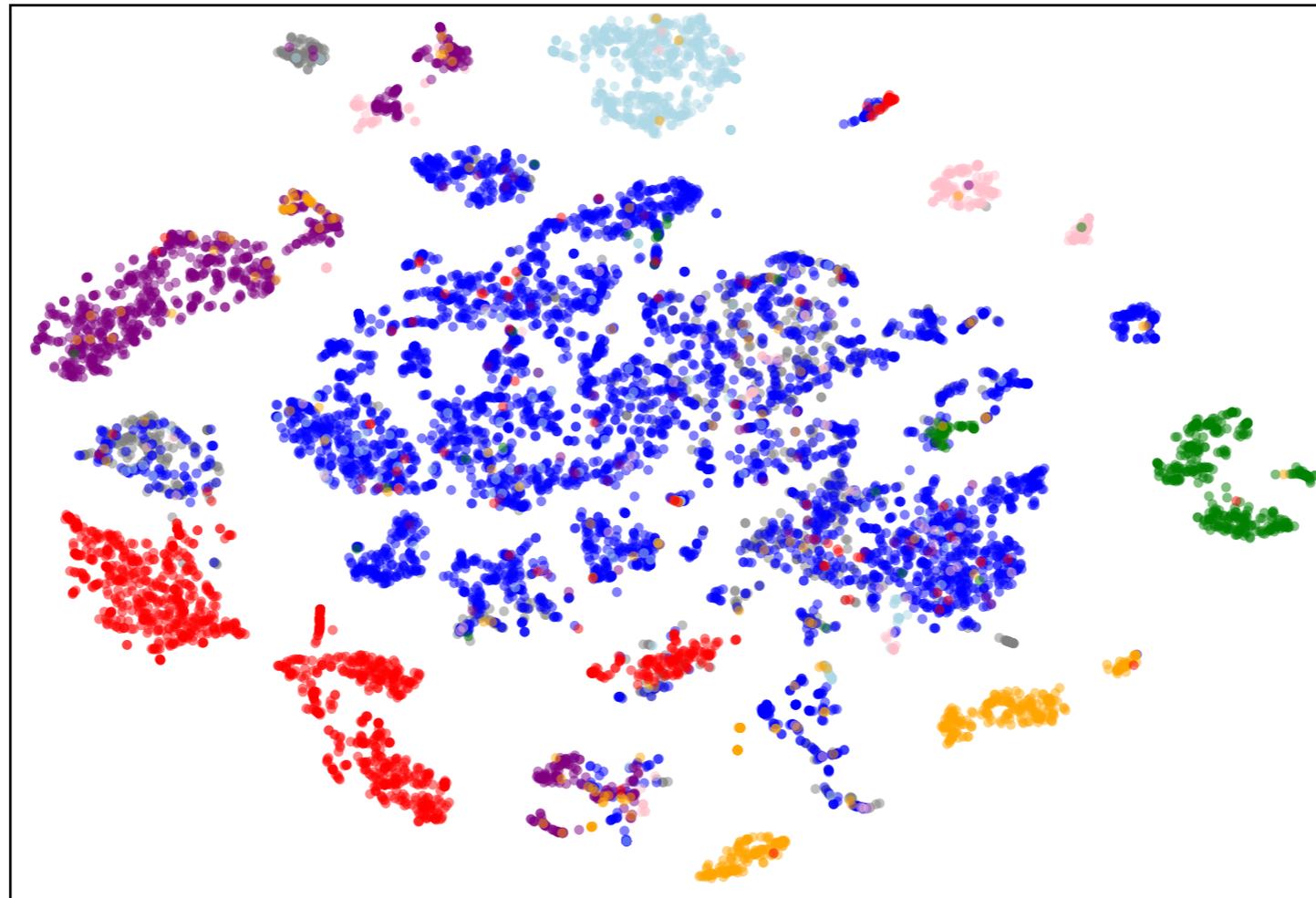
Global Normalization

Answer probability grows as search advances
(we are not multiplying probabilities!)



Note: globally normalized models remain undecided until the end word, while local models usually have spiked distributions

Type Swaps



T-SNE Question hidden state naturally clusters according to question type.

How can we exploit this?

Type Swaps

- Common SQuAD pitfall: pick wrong answer with right type (human, organization, etc.)
- Solution: increase typed-based QA pairs

Who said in **December 2012** that the fight would change from military to law enforcement?

... Basic objectives of the **Bush Administration** "war on terror", such as targeting al Qaeda and building international counterterrorism alliances, remain in place. In **December 2012**, **Jeh Johnson**, the General Counsel of the **Department of Defense** stated that the military fight will be replaced by a law enforcement operation when speaking at **Oxford University** ...

Answer: **Jeh Johnson**

Type Swaps

- Common SQuAD pitfall: pick wrong answer with right type (human, organization, etc.)
- Solution: increase typed-based QA pairs

Who said in **April 25, 2011** that the fight would change from military to law enforcement?

... Basic objectives of the **Cabinet of Japan** "war on terror", such as targeting al Qaeda and building international counterterrorism alliances, remain in place. In **April 25, 2011**, **Sheryl Sandberg**, the General Counsel of the **ministry of education** stated that the military fight will be replaced by a law enforcement operation when speaking at **Ain Shams University** . . .

Answer: **Sheryl Sandberg**

Type Swaps

- Common SQuAD pitfall: pick wrong answer with right type (human, organization, etc.)
- Solution: increase typed-based QA pairs

Who said in **2012** that the fight would change from military to law enforcement?

... Basic objectives of the **British Empire** "war on terror", such as targeting al Qaeda and building international counterterrorism alliances, remain in place. In **2012**, **Genghis Khan**, the General Counsel of the **EMNLP** stated that the military fight will be replaced by a law enforcement operation when speaking at **George Washington University** ...

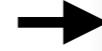
Answer: **Genghis Khan**

Type Swaps

... Basic objectives of the **Bush Administration** war on terror", such as targeting al Qaeda and building international counterterrorism alliances, remain in place. In **December 2012** **Jeh Johnson** the General Counsel of the **Department of Defense** stated that the military fight will be replaced by a law enforcement operation when speaking at **Oxford University** ...



- Bush Administration
- Department of Defense
- December 2012
- Jeh Johnson
- Oxford University

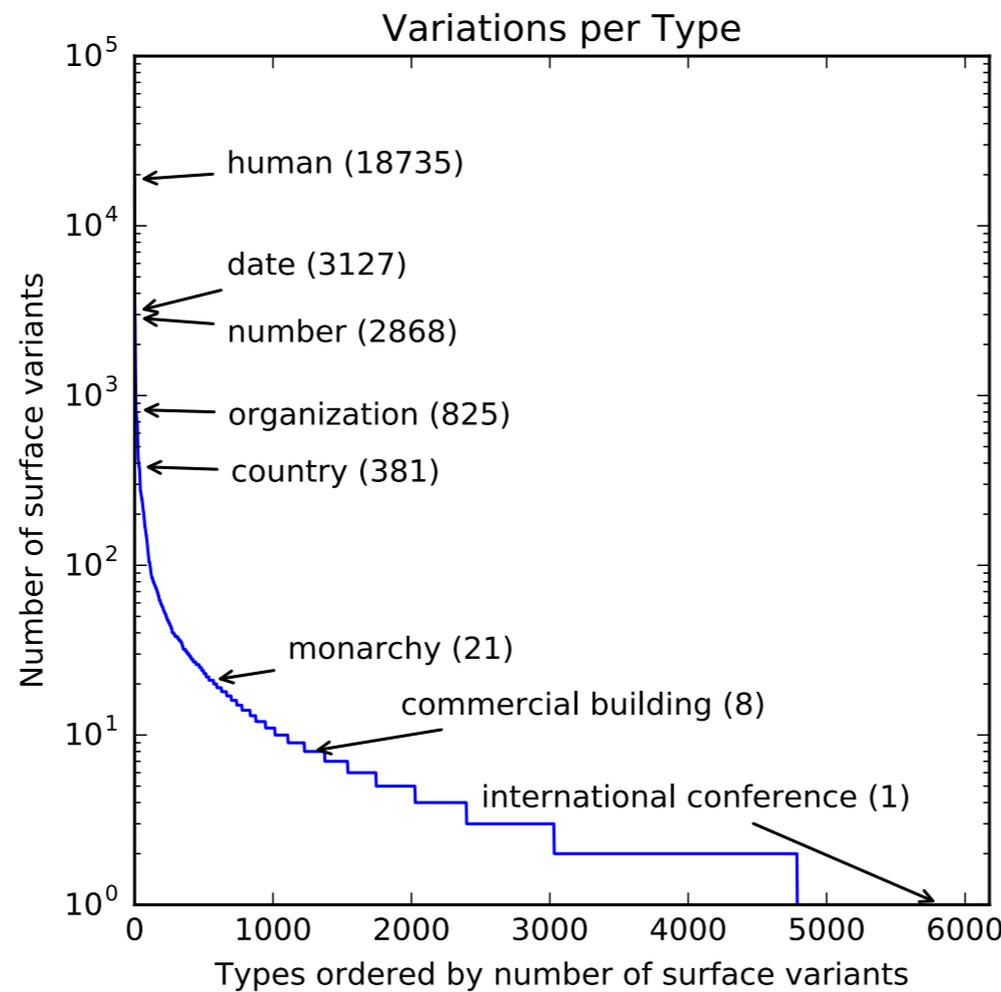


... Basic objectives of the **Cabinet of Japan** war on terror", such as targeting al Qaeda and building international counterterrorism alliances, remain in place. In **April 25, 2011** **Sheryl Sandberg** the General Counsel of the **ministry of education** stated that the military fight will be replaced by a law enforcement operation when speaking at **Ain Shams University** ...

1) Extract Entities

2) Assign Types

3) Swap with same type variations



$2.92 \cdot 10^{369}$
unique documents

Experiments

- Evaluate GNR against baselines on SQuAD dev set (100,000 QA pairs)¹
- GNR Ablations:
 - Data Augmentation
 - Global Normalization
- Measure Speedup

¹ <https://rajpurkar.github.io/SQuAD-explorer/>

Comparison

Model	EM	F1
GNR	68.4	76.2
Bi-Attention-Flow (Seo et al., 2016)	67.7	77.3
Rasor (Lee et al., 2016)	66.4	74.9
DCN (Xiong et al., 2016)	65.4	75.6
FastQA (Weissenborn et al., 2017)	67.8	76.3
R-Net (Wang et al., 2017)	72.3	80.6

Model	EM	F1
GNR	68.4	76.2
GNR w/o Global Norm	67.21	76.0
GNR w/o Type Swaps	66.6	75.0

Data Augmentation

Impact of number of augmented examples:

Number of Swaps	EM	F1
0	66.6	75.0
1000	66.9	75,0
10 000	68.4	76.21
50 000	66.8	75.3
100 000	66.1	74.3

Impact of Type Swaps on the DCN+¹:

Number of Swaps	Train F1	Dev F1
0	81.3	78.1
50 000	72.5	78.2

¹ Updated DCN model, see <https://rajpurkar.github.io/SQuAD-explorer/>

Speedup

- Full dev set, batch size 32, average 5 runs, on Titan X:
 - Bi-Attention-Flow¹: 1260.23 ±17.26 seconds
 - GNR: **51.58 ±0.266 seconds**
- Key reasons:
 - Efficient batching of Beam Search
 - Only rank subset of spans
 - Factorize search with document structure

¹github.com/allenai/bi-att-flow

Conclusion

Key contributions:

- Learning-to-Search w/. early updates & global norm enables conditional computation.
- Data augmentation that improves performance
- 24.7x speedup over bi-attention-flow
- Achieve \geq results than bi-directional attention

Future Work

- Conditional computation for generative models & large search spaces
- Program induction/search with perfect simulator
- Model amplification (AlphaZero-style)
- Type Swaps on other NLP tasks? Grammar-aware type-swaps? Adversarial Type Swaps?

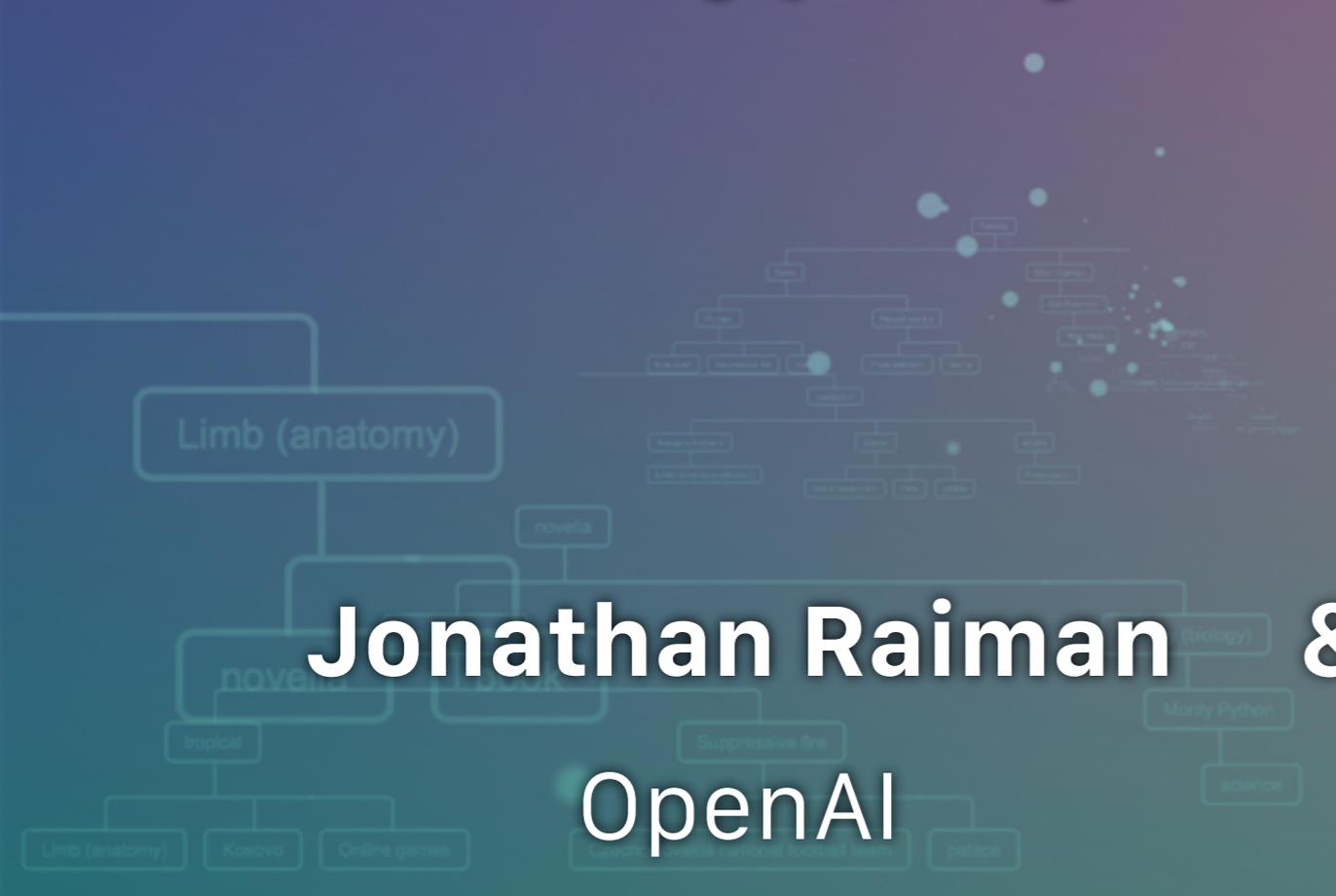
Code & Dataset:

[github.com/baidu-research/
GloballyNormalizedReader](https://github.com/baidu-research/GloballyNormalizedReader)

Thank You!

DeepType

Multilingual Entity Linking by Neural
Type System Evolution



Jonathan Raiman &

OpenAI

Olivier Raiman

Agilience

Entity Linking

The man saw a Jaguar speed on the highway.



Animal



Vehicle

Entity Linking

The prey saw a Jaguar cross the jungle.



Animal



Vehicle

The man saw a Jaguar speed on the highway.

vehicle

With types accuracy reaches **98.6-99%**

(type oracle on TAC KBP 2010/CoNLL YAGO)

The prey saw a Jaguar cross the jungle.

animal

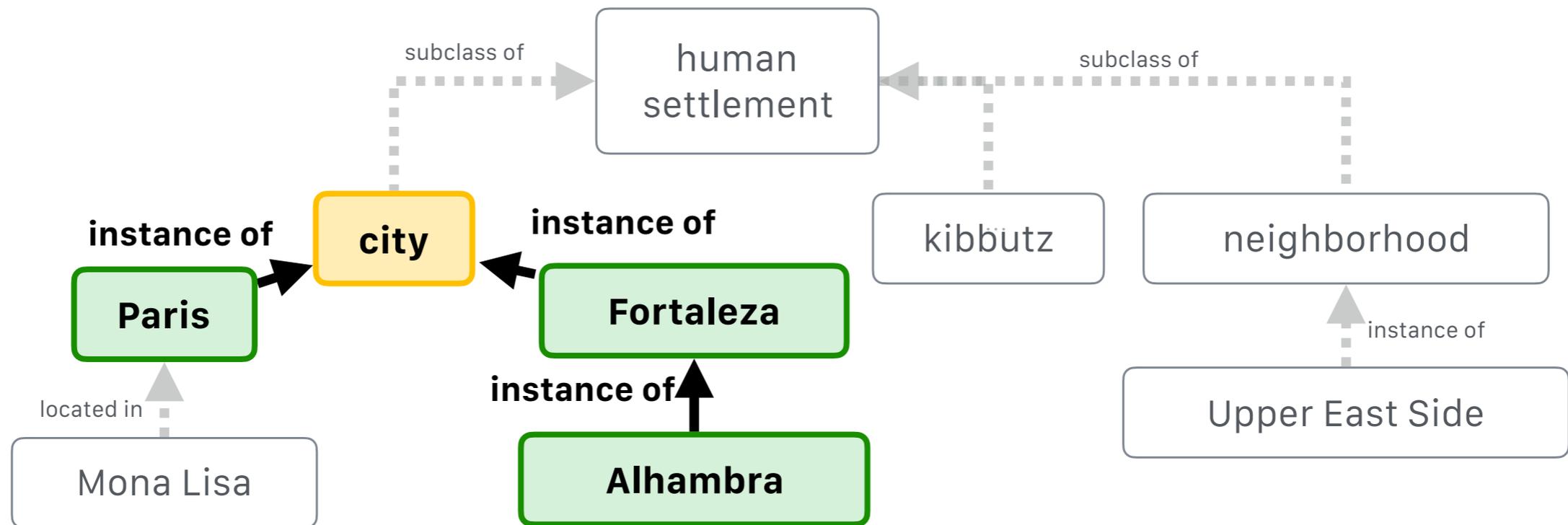
Summary

- Design a neural type system
- Results
- Contributions

Type generation

Wikidata is a graph with 40M+ entities

`isCity = child(city, instance of)`
root relation



● Root ● isCity ○ non-member

Design a Type System

$$\max_{\mathcal{A}} \max_{\theta} S_{\text{model}}(\mathcal{A}, \theta) = \frac{\sum_{(m, e_{\text{GT}}, \mathcal{E}_m) \in M} \mathbb{1}_{e_{\text{GT}}}(e^*)}{|M|}$$

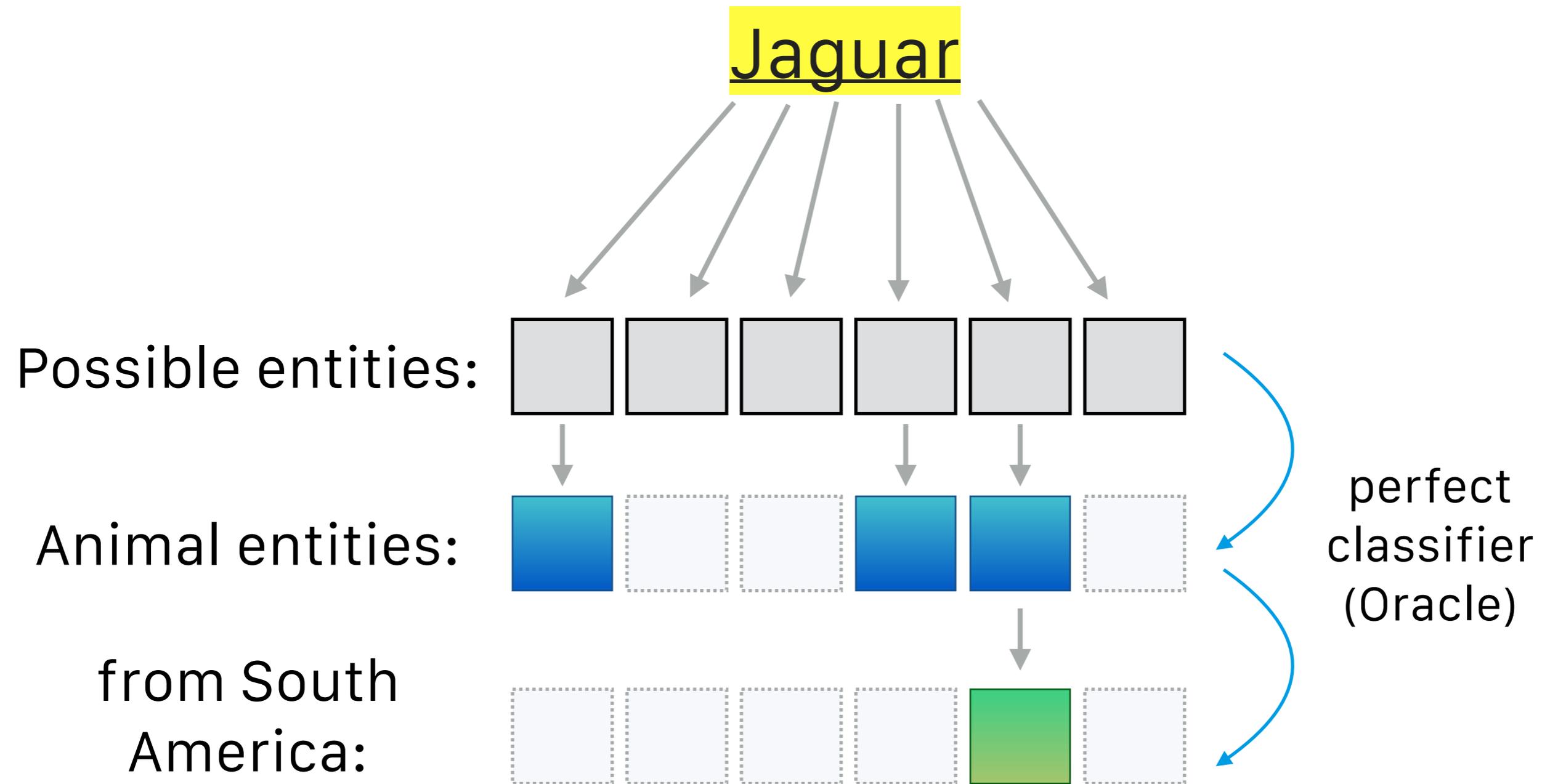
Type System **Disambiguation accuracy**

Intractable mixed integer problem

Subproblems

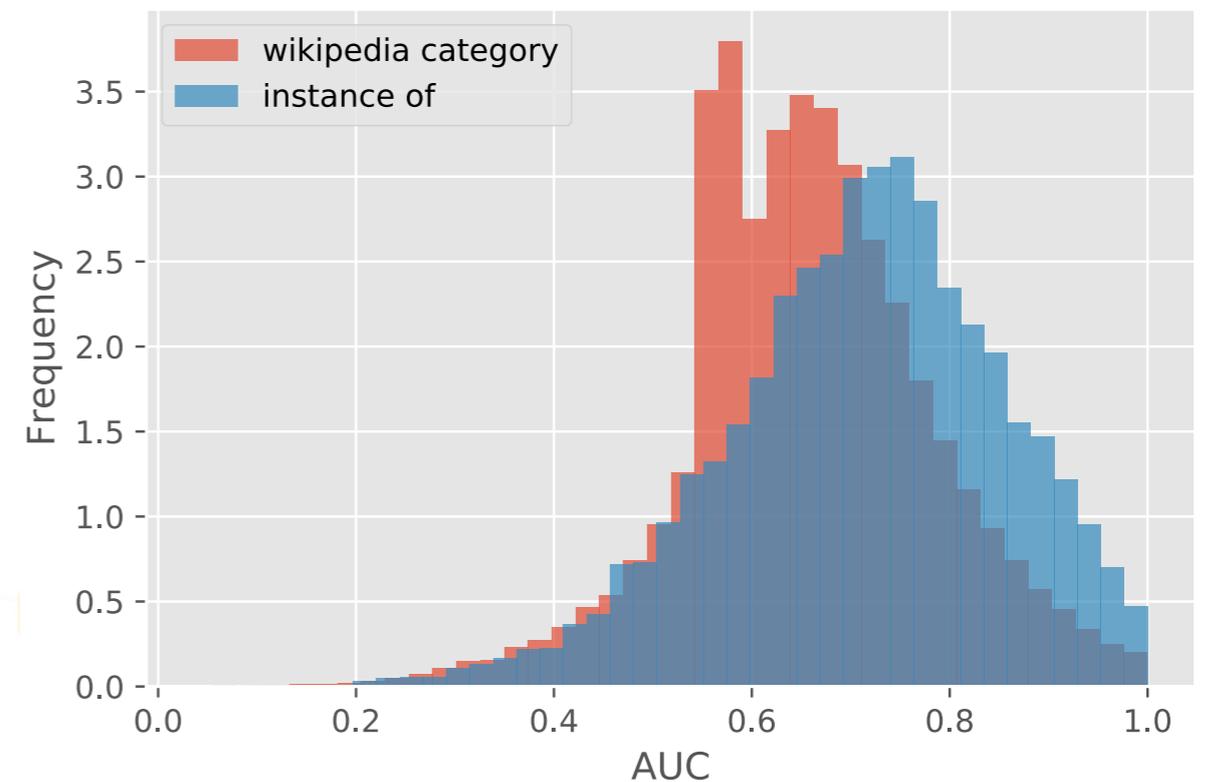
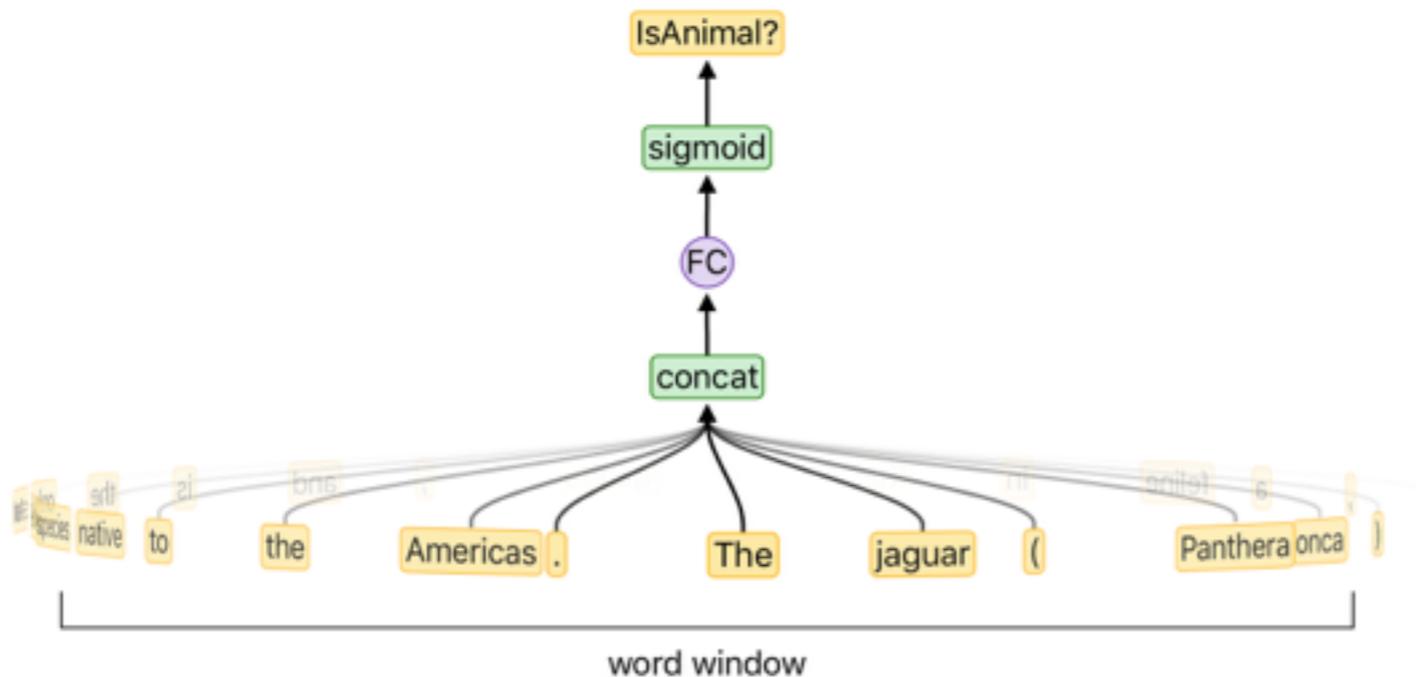
1. Stochastic optimization/heuristic search to design type system
2. Gradient descent to train a type classifier

Oracle Accuracy



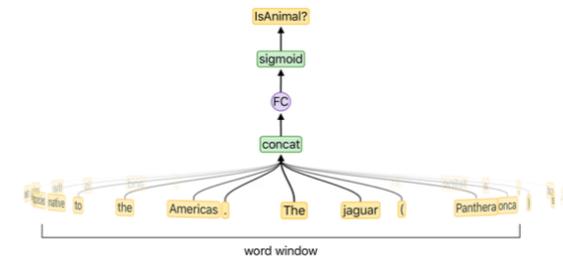
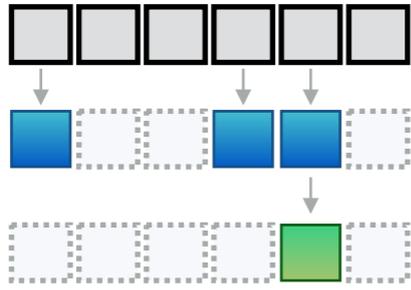
Type Learnability

- Which types are predictable from context?
- Train a proxy binary classifier for each type
- AUC* of classifier is an estimate of Learnability



* average AUC over 4 training runs

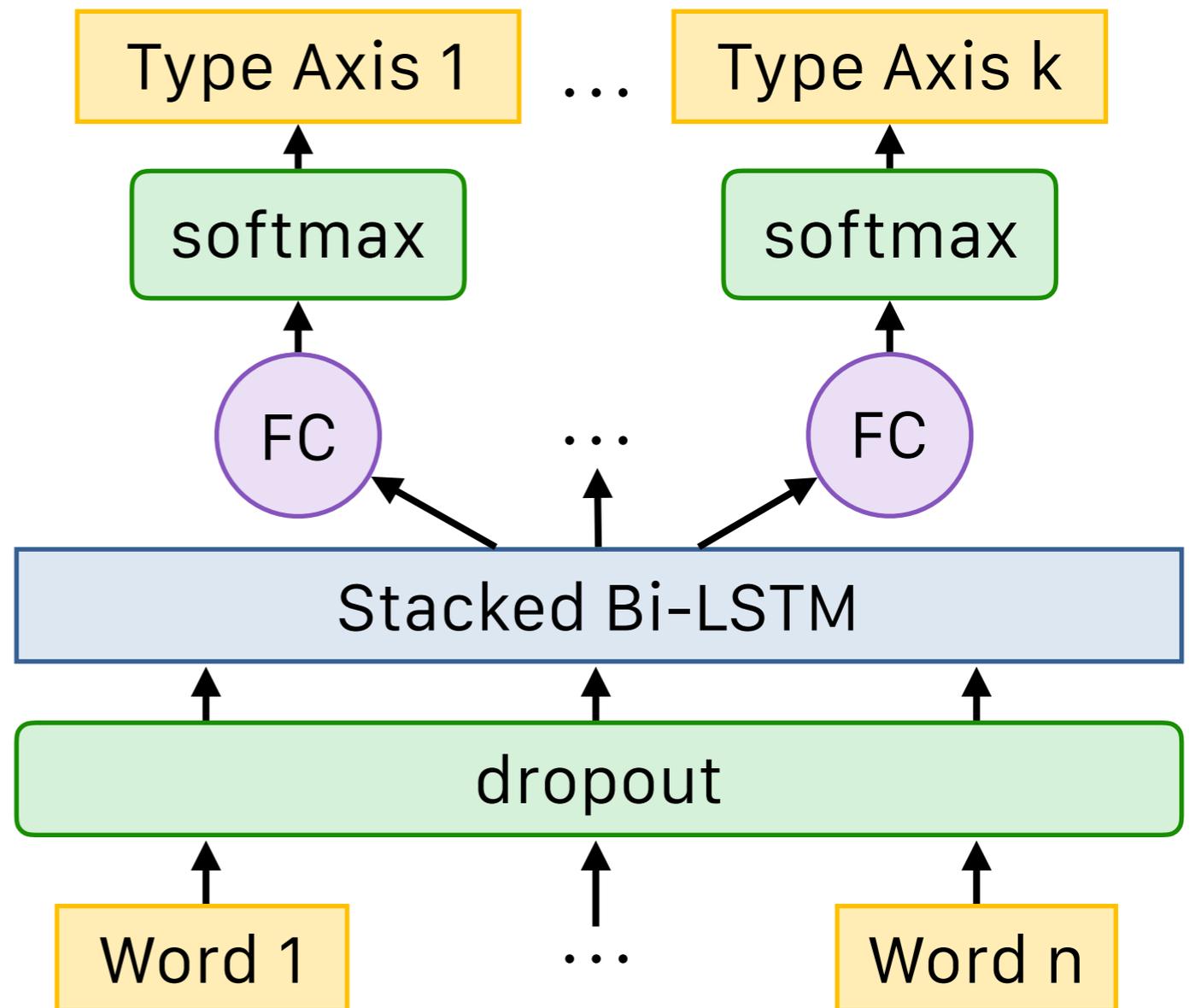
Type System Evolution



$$J(\mathcal{A}) = (\text{Accuracy}_{\text{oracle}}(\mathcal{A}) - \text{Accuracy}_{\text{greedy}}) \cdot \text{Learnability}(\mathcal{A}) - \lambda \cdot |\mathcal{A}|$$

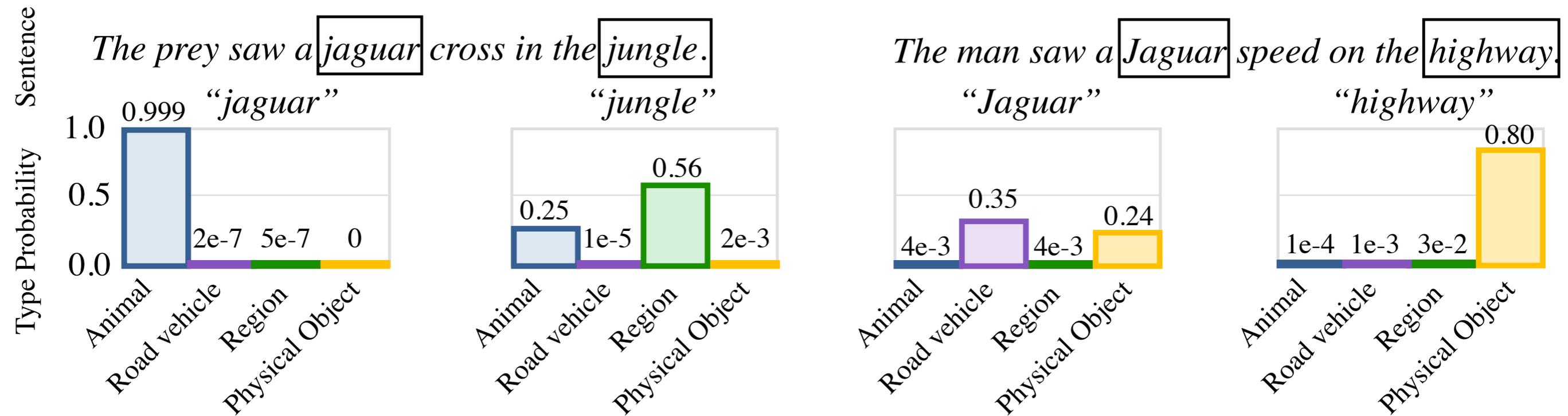
Diminishing Returns

Train a type classifier



Intra-wiki links are type labels!
(in any Wikipedia language)

Inference



Entity	jaguar	Jaguar	jungle	jungle	jaguar	Jaguar	highway	Highway
Type	Animal	Road vehicle	Region	Music	Animal	Road vehicle	Physical Object	Film
only link Prob.	0.29	0.60	0.35	0.17	0.29	0.60	0.85	0.04
Prob. w/. types	1.0	0.0	1.0	0.0	0.0	1.0	1.0	0.0

Results

	Model	CoNLL (YAGO)	TAC 2010
no types	Globerson et al. 2016	91,70 %	87,20 %
	Yamada et al. 2016	91,50 %	85,20 %
	NTEE (Yamada et al. 2017)	-	87,70 %
types	DeepType (human types)	93,11 %	90,74 %
	DeepType (greedy)	94,15 %	90,85 %
	DeepType (GA)	94,88 %	90,31 %
	DeepType (CEM)	93,96 %	90,30 %

Contributions

- Outperform state of the art on several entity linking benchmark datasets
- Add entities without retraining by specifying their types
- Design & integrate symbolic structure to constrain neural network outputs

Code:

github.com/openai/deeptype

Objective

- Given:
 - ambiguous mentions $m \in M$
 - The ground truth entity e_{GT} for each m
 - Model prediction e^*
 - Model accuracy S_{model}

$$\max_{\mathcal{A}} \max_{\theta} S_{\text{model}}(\mathcal{A}, \theta) = \frac{\sum_{(m, e_{GT}, \mathcal{E}_m) \in M} \mathbb{1}_{e_{GT}}(e^*)}{|M|}$$

- Find type system \mathcal{A} , and parameters θ to maximise disambiguation accuracy
- \mathcal{A} are discrete variables selecting the types to use

Inference

- $\mathbb{P}(\text{types}(e)|c)$ = compute type probabilities per token
- $\mathbb{P}_{\text{link}}(e|c)$ = # intra-wiki links from anchor \rightarrow articles
- Baye's rule, entity e , context c :
 - $\mathbb{P}(e|c) = \mathbb{P}_{\text{link}}(e|c) * \mathbb{P}(\text{types}(e)|c)$

Type system objective

$$J(\mathcal{A}) = (\text{ACC}_{\text{oracle}}(\mathcal{A}) - \text{ACC}_{\text{greedy}}) \cdot \text{Learnability}(\mathcal{A}) - \lambda \cdot |\mathcal{A}|$$