

# Overflow + Case Study

LING 571 — Deep Processing Methods in NLP  
Shane Steinert-Threlkeld

# Announcements

- HW9:
  - Coarse-grained sense labels (e.g. “Comparison”, not “Comparison.Contrast”)
  - Pre-processing: not required, but lower-casing and removing punctuation are good ones
  - Unknown words (even after processing); ~two options when lookup fails
    - All zeros vector
    - Average of all GloVe vectors
  - For these last two points (and any other choices): no “wrong” answers, but tell us what you did in your readme!
- AMA for Wednesday discussion: <https://forms.gle/qgM8V5Pp1mQL6gCg8>

# Coreference Resolution Humor

# Coreference Resolution Humor



# Coreference Resolution Humor pt. 2

A young artist exhibits his work for the first time and a well known art critic is in attendance.

The critic says to the young artist, "would you like my opinion on your work?"

"Yes, " says the artist.

"It's worthless," says the critic

The artist replies, "I know, but tell me anyway."

# Coreference Resolution Humor pt. 2

A young artist exhibits his work for the first time and a well known art critic is in attendance.

The critic says to the young artist, "would you like **my opinion** on **your work**?"

"Yes, " says the artist.

"**It's** worthless," says the critic

The artist replies, "I know, but tell me anyway."

# Roadmap

- Overflow: coherence and cohesion
- Case study
  - deep vs. shallow processing in question answering
  - A bit on then vs now

# Question-Answering: A Case Study in Shallow vs. Deep Methods



# Question Answering: The Problem

- Grew out of information retrieval community

# Question Answering: The Problem

- Grew out of information retrieval community
- Document retrieval is great, but...
  - Sometimes you don't just want a ranked list of documents.
  - Sometimes you want an answer to a question
    - Short answer, possibly with supporting context

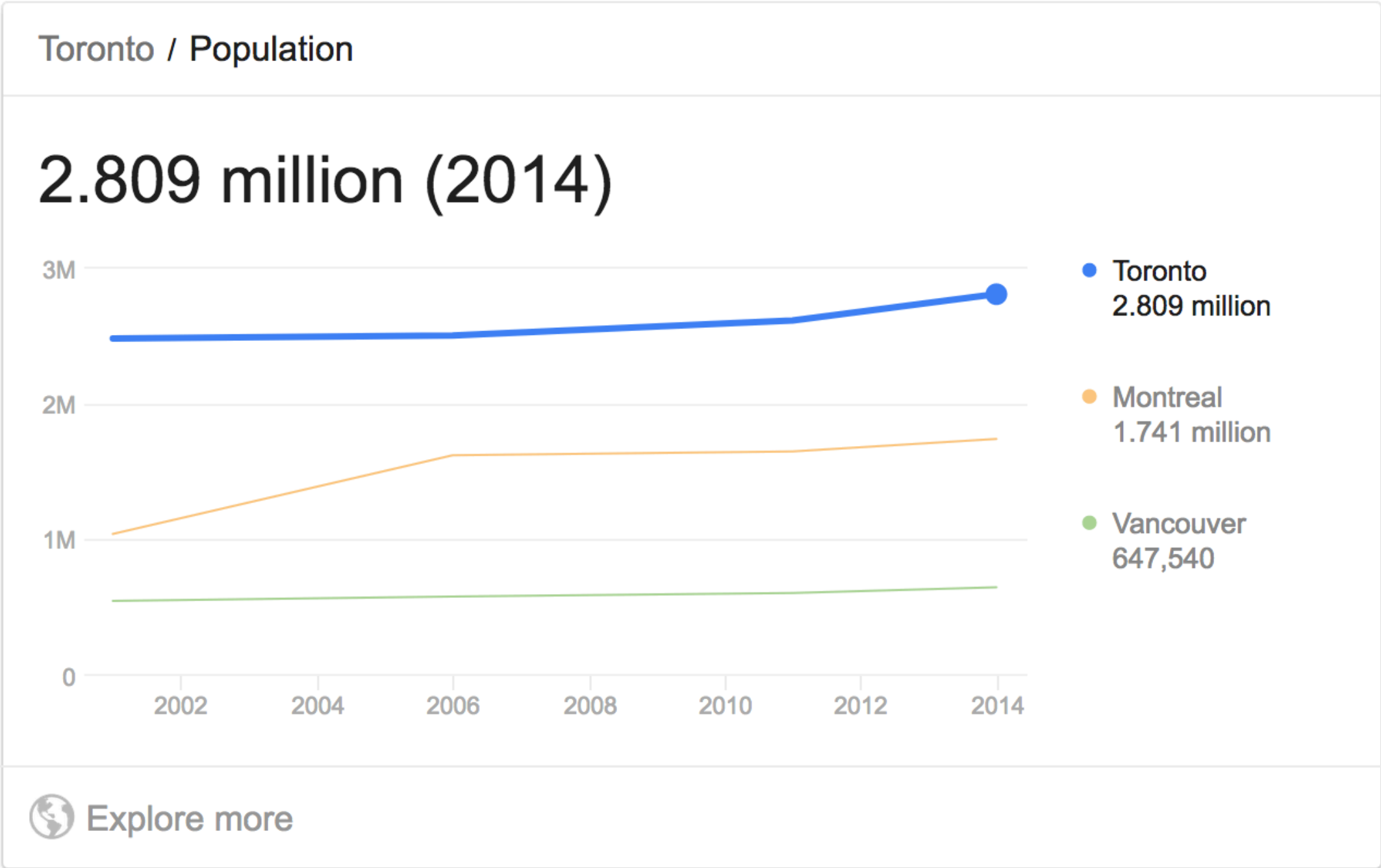
# Question Answering: The Problem

- Grew out of information retrieval community
- Document retrieval is great, but...
  - Sometimes you don't just want a ranked list of documents.
  - Sometimes you want an answer to a question
    - Short answer, possibly with supporting context
- People ask questions on the web
  - *Which English translation of the Bible is used in official Catholic liturgies?*
  - *Who invented surf music?*
  - *What are the seven wonders of the world?*
  - These account for 12–15% of web log queries

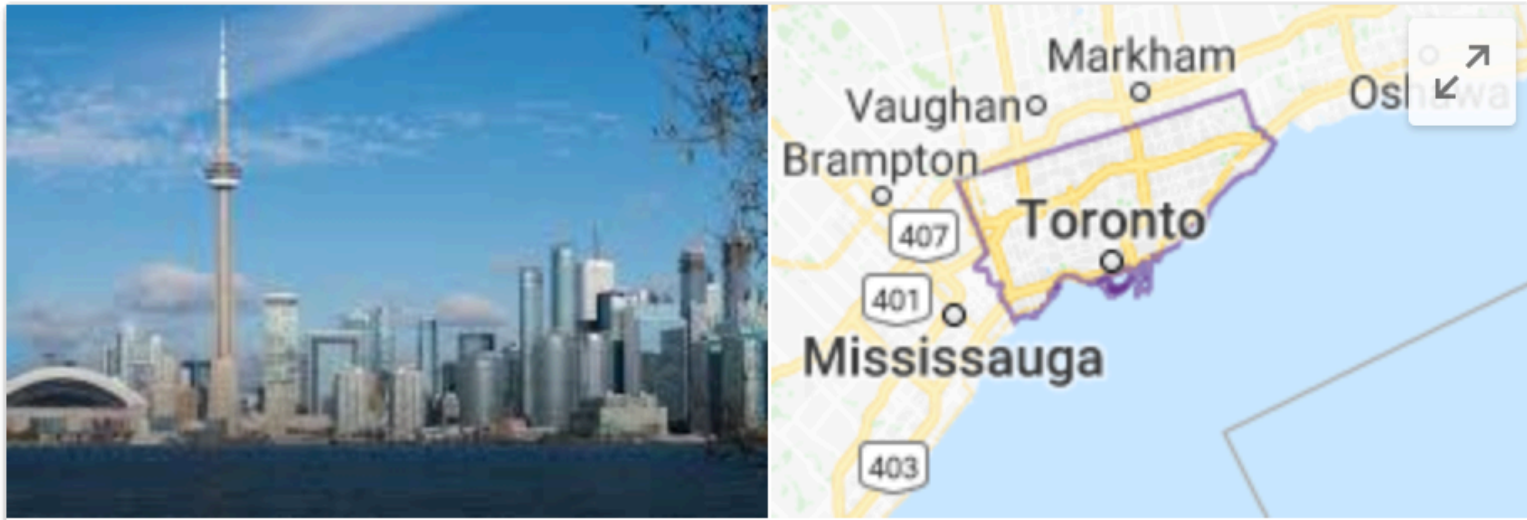
# Search Engines and Questions

- What do search engines do with questions?
  - Increasingly, try to answer questions
  - Especially for Wikipedia infobox types of info
  - Backoff to keyword search
- How well does this work?

# What Canadian city has the largest population?



- People also ask
- What are the 3 largest cities in Canada by population?
  - What are the 5 major cities in Canada?



## Toronto

City in Ontario, Canada

Toronto, the capital of the province of Ontario, is a major Canadian city along Lake Ontario's northwestern shore. It's a dynamic metropolis with a core of soaring skyscrapers, all dwarfed by the iconic, free-standing CN Tower. Toronto also has many green spaces, from the orderly oval of Queen's Park to 400-acre High Park and its trails, sports facilities and zoo.

### Population elsewhere

Canada	35.54 million (2014)
New York City	8.472 million (2014)
Chicago	2.719 million (2014)

Sources include: World Bank, United States Census Bureau

\* Answer from late 2022; more on this later

# What's the weather in Toronto?

Google search results for "what's the weather in toronto". The results show the current weather in Toronto, ON, Canada, and a forecast for the next eight days.

Results for **Toronto, ON, Canada** · Choose area

**Weather**  
Monday 1:00 PM  
Partly cloudy

Temperature | Precipitation | Wind

37° 29° (Mon) | 35° 31° (Tue) | 37° 30° (Wed) | 33° 24° (Thu) | 30° 27° (Fri) | 33° 33° (Sat) | 38° 35° (Sun) | 44° 38° (Mon)

Day	Mon	Tue	Wed	Thu	Fri	Sat	Sun	Mon
Time	2 PM	5 PM	8 PM	11 PM	2 AM	5 AM	8 AM	11 AM
Weather	Partly cloudy	Partly cloudy	Cloudy with rain	Cloudy with snow	Cloudy with snow	Cloudy with snow	Partly cloudy	Cloudy with rain
High/Low	37° / 29°	35° / 31°	37° / 30°	33° / 24°	30° / 27°	33° / 33°	38° / 35°	44° / 38°

Google Weather · Feedback

12/02/24

# What's the weather in the largest city in Canada?

The screenshot shows a Google search interface with the query "what's the weather in the largest city in Canada?". The search results are displayed in a dark theme. The top result is from "The Weather Channel" with the URL "https://weather.com › weather › tenday › La+Cité+Limo...". The title of the result is "10 Day Weather-La Cité-Limoilou, Quebec, Canada" and the snippet reads "Partly cloudy. Low 17F. Winds light and variable. UV Index0 of 11. Moonrise ...". Below the search results, there is a section titled "People also ask" with four questions: "What is the weather mostly in Canada?", "Is it hot or cold in Canada right now?", "How hot does Canada get in the summer?", and "What is the most extreme weather in Canada?". Each question has a downward arrow icon to its right. At the bottom right of the search results area, there is a "Feedback" link.

12/02/24

# *What is the total population of the ten largest capitals in the US?*

- Rank 1 snippet:
  - As of 2013, 61,669,629 citizens lived in ***America's 100 largest cities***, which was 19.48 percent of the nation's ***total population***.
  - See the top 50 ***U.S. cities by population*** and rank. ... The table below lists the *largest 50 cities in the*
  - The table below lists the *largest 10 cities in the United States...*



# Breaking QA Systems



<https://twitter.com/xkcd/status/1333529967079120896>

# Breaking QA Systems

## CREPE: Open-Domain Question Answering with False Presuppositions

Xinyan Velocity Yu<sup>†</sup> Sewon Min<sup>†</sup> Luke Zettlemoyer<sup>†</sup> Hannaneh Hajishirzi<sup>†,‡</sup>

<sup>†</sup>University of Washington <sup>‡</sup>Allen Institute for Artificial Intelligence

{xyu530, sewon, lsz, hannaneh}@cs.washington.edu

[333529967079120896](https://arxiv.org/abs/2305.10263)

### Abstract

Information seeking users often pose questions with false presuppositions, especially when asking about unfamiliar topics. Most existing question answering (QA) datasets, in



**Question:** If there's an equal and opposite reaction for everything, how does any action happen? Isn't it balanced out by the opposite reaction?

#### Newton's laws of motion

From Wikipedia, the free encyclopedia

Overly brief paraphrases of the third law, like "action



<https://doi.org/10.18653/v1/2023.acl-long.583>

# Search Engines and QA

- Search for exact question string
  - “Do I need a visa to go to Japan?”
    - Result: Exact match on Yahoo! Answers
    - Find “Best Answer” and return following chunk

# Search Engines and QA

- Search for exact question string
  - “Do I need a visa to go to Japan?”
    - Result: Exact match on Yahoo! Answers
    - Find “Best Answer” and return following chunk
- Works great... if the question matches exactly
  - Many websites are building archives
  - What happens if it doesn't match?
    - “Question mining” tries to learn paraphrases of questions to get answers.

# Perspectives on QA

- TREC QA track (~2000— )
  - Initially pure factoid questions, with fixed length answers
    - Based on large collection of fixed documents (news)
    - Increasing complexity: definitions, biographical info, etc
      - Single response

# Perspectives on QA

- TREC QA track (~2000— )
  - Initially pure factoid questions, with fixed length answers
    - Based on large collection of fixed documents (news)
    - Increasing complexity: definitions, biographical info, etc
      - Single response
- Reading comprehension ([Hirschman et al, 1999](#)— )
  - Think SAT/GRE
    - Short text or article (usually middle school level)
    - Answer questions based on text
  - Also, “Machine Reading”
  - [SQuAD](#)

# Perspectives on QA

- TREC QA track (~2000— )
  - Initially pure factoid questions, with fixed length answers
    - Based on large collection of fixed documents (news)
    - Increasing complexity: definitions, biographical info, etc
      - Single response
- Reading comprehension ([Hirschman et al, 1999](#)— )
  - Think SAT/GRE
    - Short text or article (usually middle school level)
    - Answer questions based on text
  - Also, “Machine Reading”
  - [SQuAD](#)
- And, of course, [Jeopardy! and Watson](#)

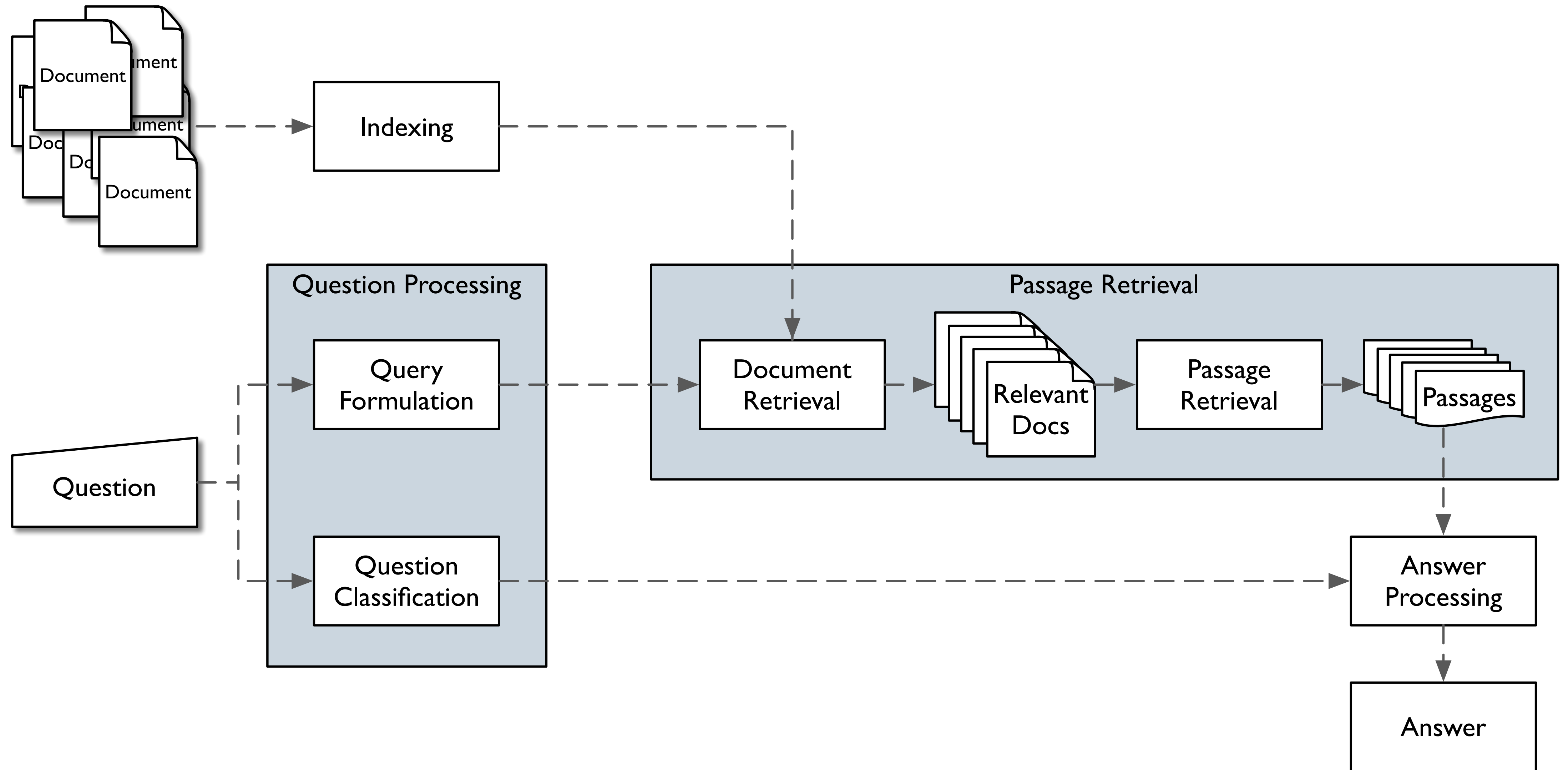
# Question Answering (*a la* TREC)

Question	Answer
Where is the Louvre Museum located?	in Paris, France
What's the abbreviation for limited partnership?	L.P.
What are the names of Odin's ravens?	Huginn and Muninn
What currency is used in China?	the yuan
What kind of nuts are used in marzipan?	almonds
What instrument does Max Roach play?	drums
What's the official language of Algeria?	Arabic
What is the telephone number for the University of Colorado, Boulder?	(303) 492-1411
How many pounds are there in a stone?	14



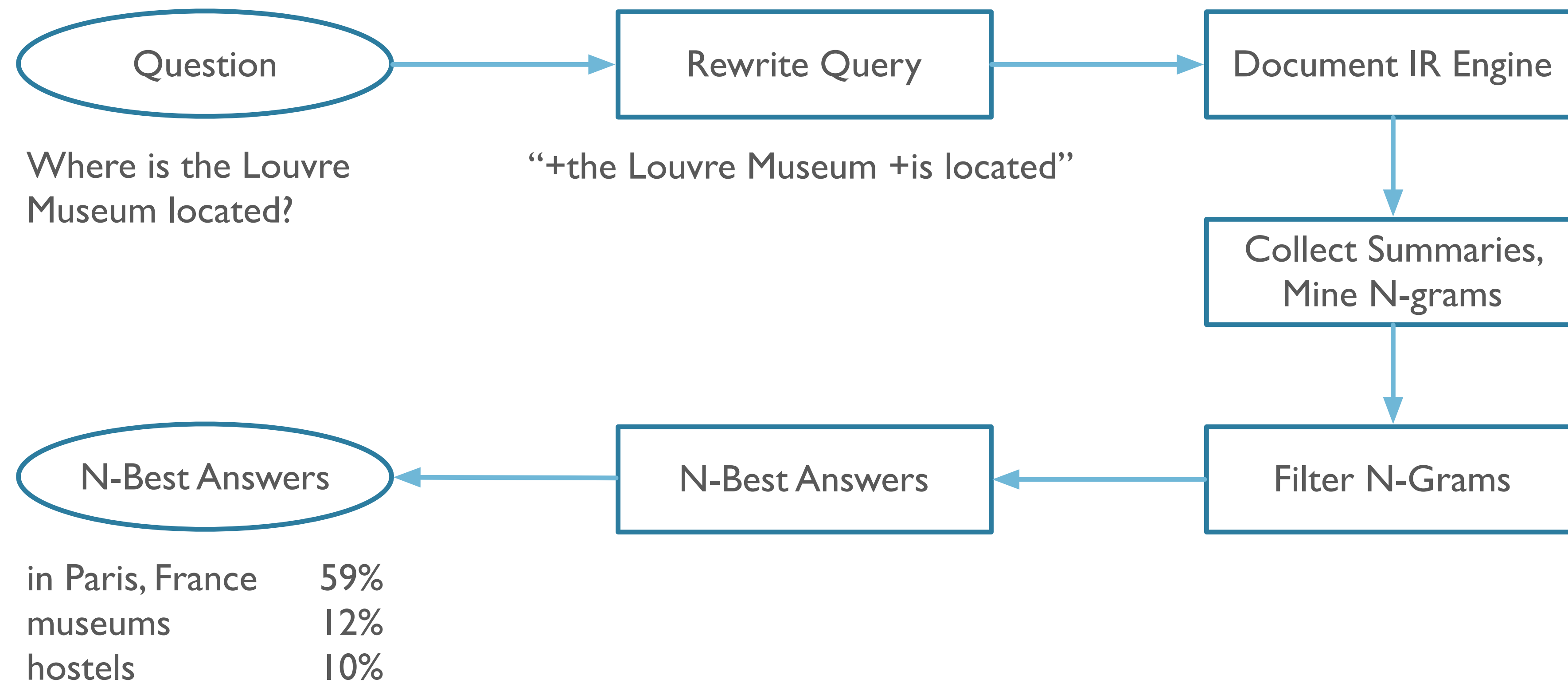
# Basic Strategy

- Given an indexed document collection...
- ...and a question...
- ...execute the following steps:
  - Query Formulation
  - Question Classification
  - Passage Retrieval
  - Answer Processing
  - Evaluation



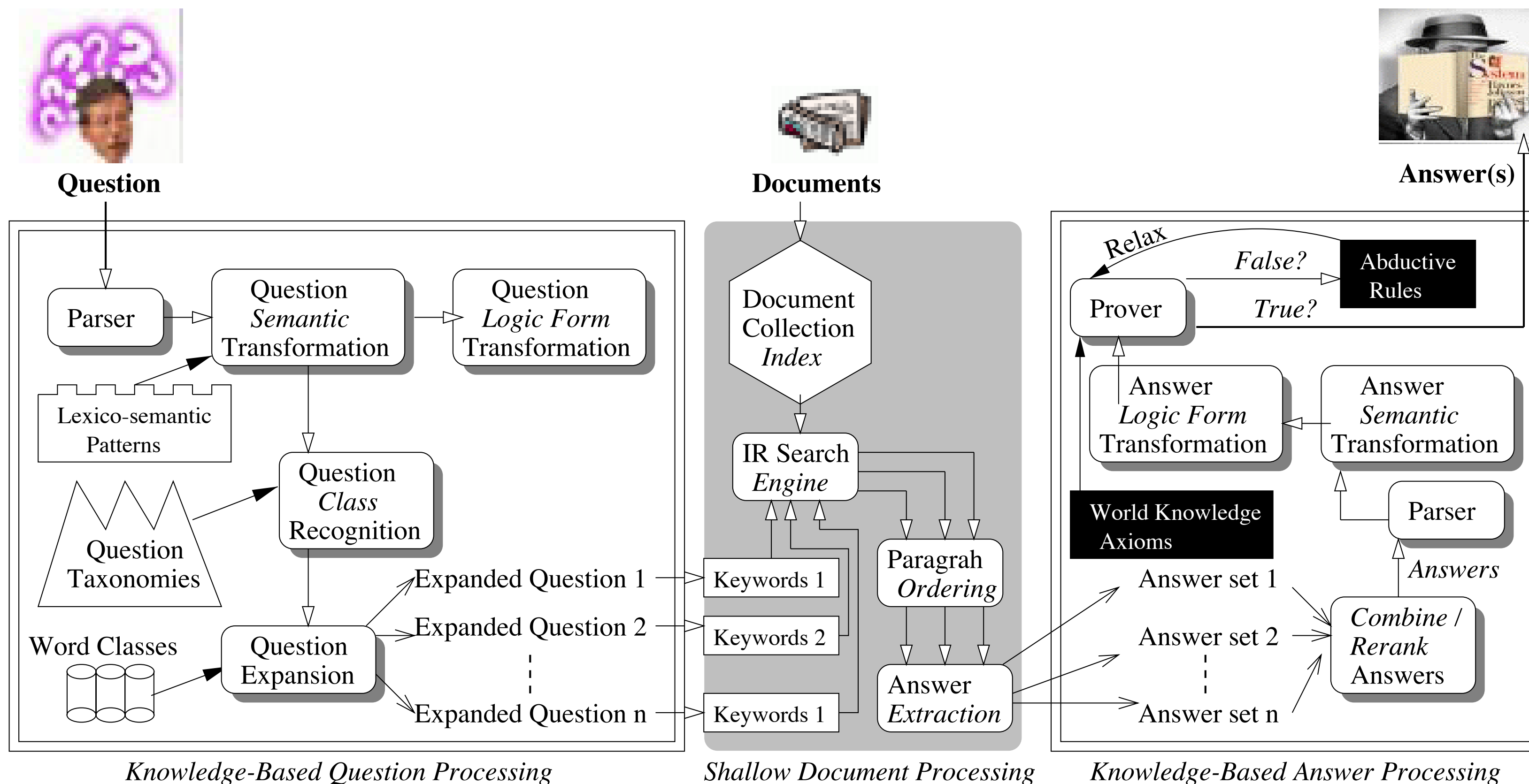
# AskMSR/Aranea (Lin, Brill)

- Shallow Processing for QA



# Deep Processing Technique for QA: LCC PowerAnswer

- Experiments with open-domain textual Question Answering, [Moldovan, Harabagiu, et al, 2000](#)



# A Victory for Deep Processing:

## TREC 2002 QA Track

Run Tag	Confidence weighted Score	Correct Answers		Number Inexact	NIL Accuracy	
		#	%		Prec	Recall
<b>LCCmain2002</b>	<b>0.856</b>	<b>415</b>	<b>83.0</b>	<b>8</b>	<b>0.578</b>	<b>0.804</b>
exactanswer	0.691	271	54.2	12	0.222	0.848
pris2002	0.610	290	58.0	17	0.241	0.891
IRST02DI	0.589	192	38.4	17	0.167	0.217
IBMPQSQACYC	0.588	179	35.8	9	0.196	0.630
uwmtB3	0.512	184	36.8	20	0.000	0.000
BBN2002C	0.499	142	28.4	18	0.182	0.087
isi02	0.498	149	29.8	15	0.385	0.109
limsiQalir2	0.497	133	26.6	11	0.188	0.196
ali2002b	0.496	181	36.2	15	0.156	0.848
ibmsqa02c	0.455	145	29.0	44	0.224	0.239
FDUTIIQAI	0.434	124	24.8	6	0.139	0.957
arana02a	0.433	152	30.4	36	0.235	0.174
nuslamp2002	0.396	105	21.0	17	0.000	0.000
pqas22	0.358	133	26.6	11	0.145	0.674

# Example of Deep Processing in LLM era

## 🔗 LINC: A Neurosymbolic Approach for Logical Reasoning by Combining Language Models with First-Order Logic Provers

Theo X. Olausson\*<sup>1</sup> Alex Gu\*<sup>1</sup> Benjamin Lipkin\*<sup>2</sup> Cedegao E. Zhang\*<sup>2</sup>  
Armando Solar-Lezama<sup>1</sup> Joshua B. Tenenbaum<sup>1,2</sup> Roger Levy<sup>2</sup>  
{theo xo, gua, lipkinb, cedzhang}@mit.edu

<sup>1</sup>MIT CSAIL <sup>2</sup>MIT BCS

*\*Equal contribution.*

### Abstract

Logical reasoning, i.e., deductively inferring the truth value of a conclusion from a set of premises, is an important task for artificial intelligence with wide potential impacts on science, mathematics, and society. While many prompting-based strategies have been proposed to enable Large Language Models (LLMs) to do such reasoning more effectively, they still appear unsatisfactory, often failing in subtle and unpredictable ways. In this work, we investigate the validity of instead reformulating such tasks as modular neurosymbolic programming, which we call LINC: Logical Inference via Neurosymbolic Computation. In LINC, the LLM acts as a semantic parser, translating premises and conclusions from natural language to expressions in first-order logic. These expressions are then offloaded to an external theorem prover, which symbolically performs deductive inference. Leveraging this approach, we observe significant performance gains on FOLIO and a balanced subset of ProofWriter for three different models in nearly all experimental conditions we evaluate. On

### 1 Introduction

Widespread adoption of large language models (LLMs) such as GPT-3 (Brown et al., 2020), GPT-4 (OpenAI, 2023), and PaLM (Chowdhery et al., 2022) have led to a series of remarkable successes in tasks ranging from text summarization to program synthesis. Some of these successes have encouraged the hypothesis that such models are able to flexibly and systematically reason (Huang and Chang, 2022), especially when using prompting strategies that explicitly encourage verbalizing intermediate reasoning steps before generating the final answer (Nye et al., 2021; Wei et al., 2022; Kojima et al., 2022; Wang et al., 2023b). However, this reasoning ability appears to be unreliable for tasks that require reasoning out of domain (Liang et al., 2022; Saparov et al., 2023), understanding negation (Anil et al., 2022), and following long reasoning chains (Dziri et al., 2023). Furthermore, while the standard approach of “scaling up” seems to improve performance across some reasoning domains, other domains, e.g., reasoning involving use of Modus Tollens, show no such improvements

<https://openreview.net/forum?id=h00GHjWDEp>

# Example of Deep Processing in LLM era

## LINC: A Neurosymbolic Approach for Logical Reasoning by Language Models with First-Order Logic Provers

Theo X. Olausson\*<sup>1</sup> Alex Gu\*<sup>1</sup> Benjamin Lipkin\*<sup>2</sup> Cedegao E. Armando Solar-Lezama<sup>1</sup> Joshua B. Tenenbaum<sup>1,2</sup> Roger L.

{theoxo, gua, lipkinb, cedzhang}@mit.edu

<sup>1</sup>MIT CSAIL <sup>2</sup>MIT BCS

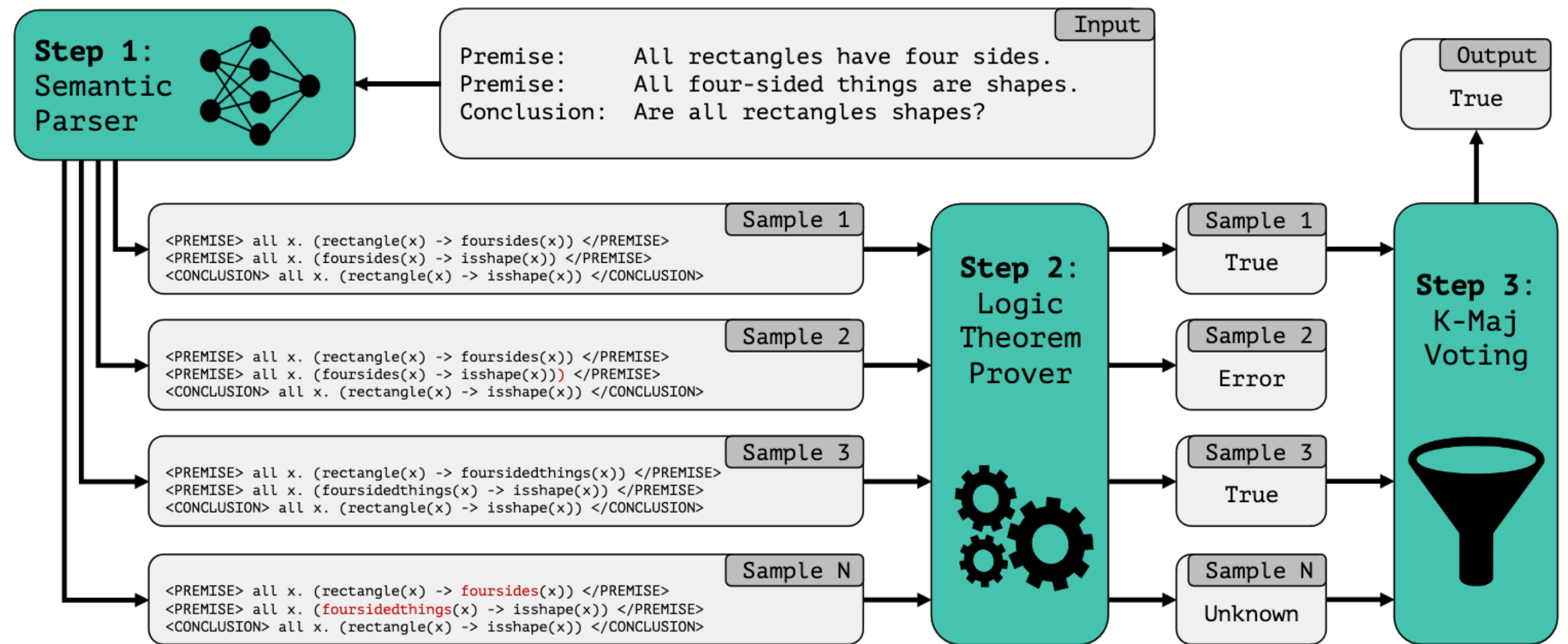
\*Equal contribution.

### Abstract

Logical reasoning, i.e., deductively inferring the truth value of a conclusion from a set of premises, is an important task for artificial intelligence with wide potential impacts on science, mathematics, and society. While many prompting-based strategies have been proposed to enable Large Language Models (LLMs) to do such reasoning more effectively, they still appear unsatisfactory, often failing in subtle and unpredictable ways. In this work, we investigate the validity of instead reformulating such tasks as modular neurosymbolic programming, which we call LINC: Logical Inference via Neurosymbolic Computation. In LINC, the LLM acts as a semantic parser, translating premises and conclusions from natural language to expressions in first-order logic. These expressions are then offloaded to an external theorem prover, which symbolically performs deductive inference. Leveraging this approach, we observe significant performance gains on FOLIO and a balanced subset of ProofWriter for three different models in nearly all experimental conditions we evaluate. On

### 1 Introduction

Widespread adoption of large LLMs (LLMs) such as GPT-3 (Brown et al., 2020), and PaLM (Chang, 2022), especially when using prompting strategies that explicitly encourage verbalizing intermediate reasoning steps before generating the final answer (Nye et al., 2021; Wei et al., 2022; Kojima et al., 2022; Wang et al., 2023b). However, this reasoning ability appears to be unreliable for tasks that require reasoning out of domain (Liang et al., 2022; Saparov et al., 2023), understanding negation (Anil et al., 2022), and following long reasoning chains (Dziri et al., 2023). Furthermore, while the standard approach of “scaling up” seems to improve performance across some reasoning domains, other domains, e.g., reasoning involving use of Modus Tollens, show no such improvements



# Example of Deep Processing in LLM era

## LINC: A Neurosymbolic Approach for Logical Reasoning by Language Models with First-Order Logic Provers

Theo X. Olausson\*<sup>1</sup> Alex Gu\*<sup>1</sup> Benjamin Lipkin\*<sup>2</sup> Cedegao E. Armando Solar-Lezama<sup>1</sup> Joshua B. Tenenbaum<sup>1,2</sup> Roger L.

{theoxo, gua, lipkinb, cedzhang}@mit.edu

<sup>1</sup>MIT CSAIL <sup>2</sup>MIT BCS

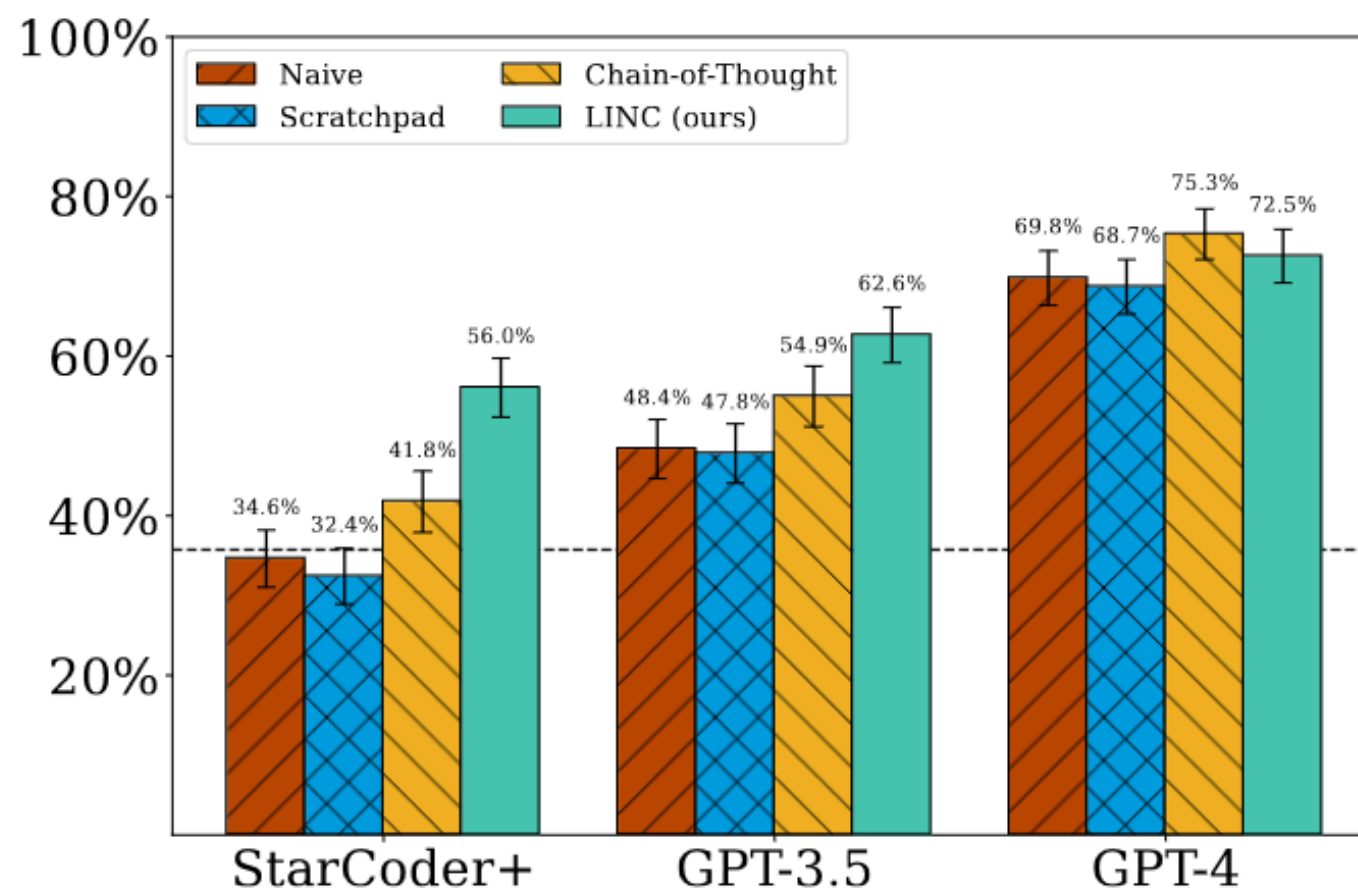
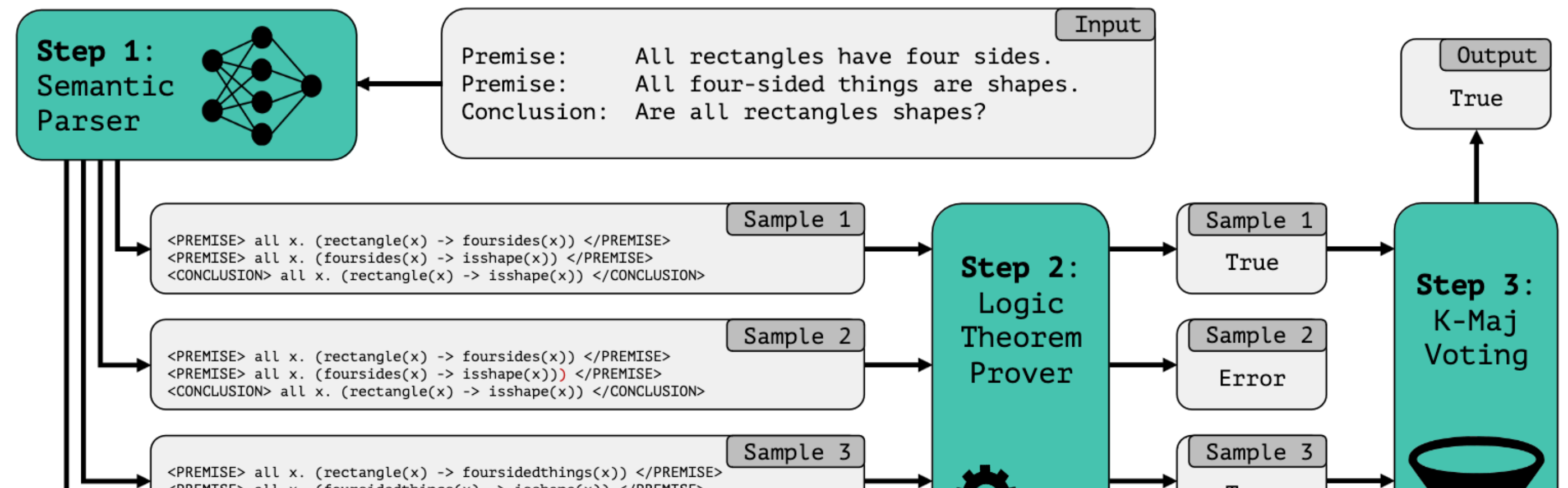
\*Equal contribution.

### Abstract

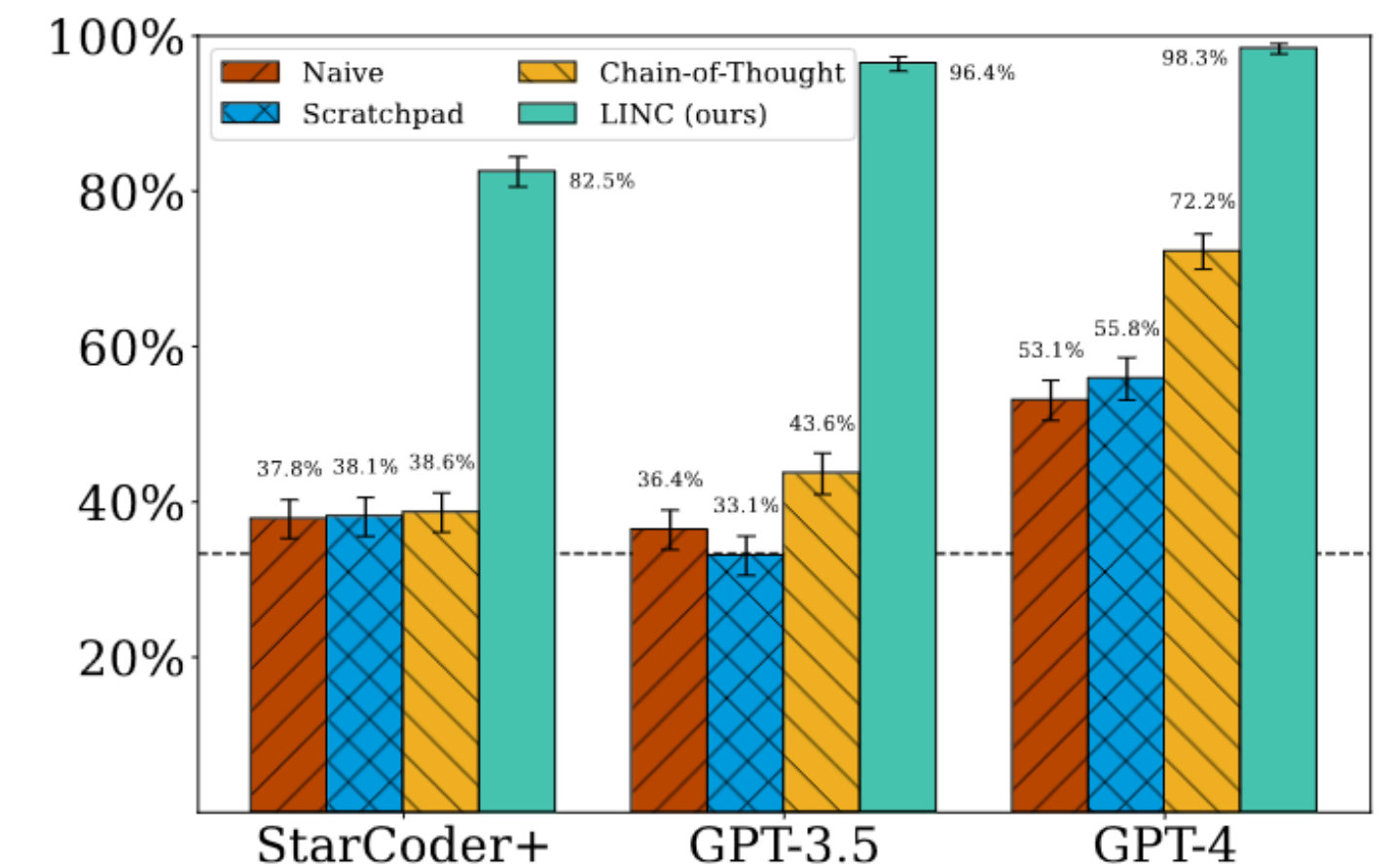
Logical reasoning, i.e., deductively inferring the truth value of a conclusion from a set of premises, is an important task for artificial intelligence with wide potential impacts on science, mathematics, and society. While many prompting-based strategies have been proposed to enable Large Language Models (LLMs) to do such reasoning more effectively, they still appear unsatisfactory, often failing in subtle and unpredictable ways. In this work, we investigate the validity of instead reformulating such tasks as modular neurosymbolic programming, which we call LINC: Logical Inference via Neurosymbolic Computation. In LINC, the LLM acts as a semantic parser, translating premises and conclusions from natural language to expressions in first-order logic. These expressions are then offloaded to an external theorem prover, which symbolically performs deductive inference. Leveraging this approach, we observe significant performance gains on FOLIO and a balanced subset of ProofWriter for three different models in nearly all experimental conditions we evaluate. On

### 1 Introduction

Widespread adoption of large language models (LLMs) such as GPT-3 (Brown et al., 2020), GPT-4 (OpenAI, 2023), and Gemini (Google, 2023) have led to a series of breakthroughs in tasks ranging from natural language generation to program synthesis. Some of these breakthroughs have encouraged the hypothesis that LLMs can flexibly and systematically learn to perform a wide range of tasks. Strategies that explicitly model intermediate reasoning steps, such as chain-of-thought (Kojima et al., 2022) and scratchpad (Chang, 2022), encourage LLMs to use this reasoning ability as a tool to solve tasks that require reasoning. However, these strategies show no such improvements in domains, e.g., Modus Tollens, where the standard approach



(a) FOLIO.



(b) ProofWriter.



# Re-visiting some questions

What Canadian city has the largest population?

By area Toronto 2020

AI Overview

Toronto is Canada's most populated city, with a population of around 6.2 million.

City	Population
Toronto	6.2 million
Montreal	4.2 million
Vancouver	2.7 million
Calgary	1.6 million
Edmonton	1.5 million
Ottawa	1.4 million
Winnipeg	825,000
Quebec City	800,000

Generative AI is experimental.

**Toronto**

1. **Toronto** (6.2 million residents) As Canada's largest city and part of the Greater Toronto Area, Toronto is a bustling metropolis sitting on the shores of Lake Ontario. Known for its diverse population, Toronto is home to world-class museums, theaters and sports teams. Oct 25, 2024

HowStuffWorks  
https://animals.howstuffworks.com > ... > Snakes

10 Largest Cities in Canada by Land Area - Yahoo  
Oct 25, 2024 — 10 Most Populated Cities in Canada  
Toronto (6.2 million residents)

The Capital Cities | Legislative Assembly of Ontario  
Ontario's capital city, Toronto, is a large metropolitan municipality with nearly 3 million residents.

# Re-visiting some questions

What Canadian city has the largest population?

By area Toronto 2020

AI Overview

Toronto is Canada's most populated city, with a population of around 6.2 million.

City	Population
Toronto	6.2 million
Montreal	4.2 million
Vancouver	2.7 million
Calgary	1.6 million
Edmonton	1.5 million
Ottawa	1.4 million
Winnipeg	825,000
Quebec City	800,000

Generative AI is experimental.

Toronto

1. **Toronto** (6.2 million residents) As Canada's largest city and part of the Greater Toronto Area, Toronto is a bustling metropolis sitting on the shores of Lake Ontario. Known for its diverse population, Toronto is home to world-class museums, theaters and sports teams. Oct 25, 2024

HowStuffWorks  
https://animals.howstuffworks.com > ... > Snakes

10 Largest Cities in Canada by Land Area - Yahoo  
Oct 25, 2024 — 10 Most Populated Cities in Canada  
Toronto (6.2 million residents)

The Capital Cities | Legislative Assembly of Ontario  
Ontario's capital city, Toronto, is a large metropolitan municipality with nearly 3 million residents.

## Notes:

- Metro area population, not city
- Because of a document it found
- Not connected to structured data source

# Re-visiting some questions

Google

What is the total population of the ten largest capitals in the US?

All Images News Shopping Videos Web Forums More Tools

AI Overview

The total population of the ten largest capital cities in the United States is over 5 million people: [↗](#)

Capital	Population
Phoenix, Arizona	1,660,272
Austin, Texas	964,254
Columbus, Ohio	892,553
Indianapolis, Indiana	867,125
Denver, Colorado	716,492
Boston, Massachusetts	694,583
Nashville, Tennessee	669,053
Oklahoma City, Oklahoma	649,021
Sacramento, California	508,529
Atlanta, Georgia	498,044

Phoenix is the largest capital city in the United States. [↗](#)

Generative AI is experimental.

The 10 Largest Capital Cities  
ThoughtCo  
Nov 18, 2020 — 1. Phoenix, Pop  
Arizona. Largest City: Phoenix. 2

List of the largest U.S. cities  
Britannica  
Nov 17, 2024 — list of the larges  
\* New York City, New York. 8,33

Biggest Cities in the U.S. (Ar  
State)  
Jun 27, 2024 — 10 Biggest Citi  
New York, New York (8,258,035

ThoughtCo  
Britannica  
Yahoo News UK

# Re-visiting some questions

**Left Screenshot:** Search: "What is the total population of the ten largest capitals in the US?"

**AI Overview**

The total population of the ten largest capital cities over 5 million people: [↗](#)

**Capital**

- Phoenix, Arizona
- Austin, Texas
- Columbus, Ohio
- Indianapolis, Indiana
- Denver, Colorado
- Boston, Massachusetts
- Nashville, Tennessee
- Oklahoma City, Oklahoma
- Sacramento, California
- Atlanta, Georgia

Phoenix is the largest capital city in the United States. [↗](#)

Generative AI is experimental. [👍](#) [🗨️](#)

**Right Screenshot:** Search: "What is the total population of the twelve largest capitals in the US?"

**WorldAtlas**  
<https://www.worldatlas.com> › World Facts

### Most Populated US State Capitals

Sep 3, 2020 — Census counts in 2020 reveal Pheonix as the **largest capital city** in the States at more than 1.7 million, followed by Austin, Columbus, Indianapolis, and Denver.  
Missing: twelve | Show results with: twelve

**People also ask** :

- What is the population of the capital of the United States? [▼](#)
- What are the largest US cities by population? [▼](#)
- What are the top 10 capitals by population? [▼](#)
- Which of the 50 state capitals is the largest by population? [▼](#)

[Feedback](#)

# Re-visiting some questions

The screenshot shows a Google search interface with the query "what year did tom hanks land on the moon". The search results are displayed in a dark theme. The top navigation bar includes "All", "Images", "News", "Videos", "Shopping", "Forums", "Web", "More", and "Tools". The main content area features an "AI Overview" section with a blue star icon. The overview text states: "Tom Hanks has not landed on the moon, but he has been involved in several projects related to the Apollo space program:". Below this, there are three sections: "Apollo 13 (1995)", "From the Earth to the Moon", and "The Moonwalkers: A Journey with Studio". The "Apollo 13 (1995)" section includes a snippet: "Hanks played one of the three astronauts on NASA's Apollo 13 spacecraft in this movie based on the true story of the near-disaster. The movie is known for the line 'Houston, we have a problem'.". The "The Moonwalkers: A Journey with Studio" section includes a snippet: "Home media. Apollo 13 was released in 1995 and on LaserDisc the following year". A Wikipedia link is also visible in the right-hand column.

Google

what year did tom hanks land on the moon

All Images News Videos Shopping Forums Web More Tools

AI Overview

Tom Hanks has not landed on the moon, but he has been involved in several projects related to the Apollo space program:

**Apollo 13 (1995)**

Hanks played one of the three astronauts on NASA's Apollo 13 spacecraft in this movie based on the true story of the near-disaster. The movie is known for the line "Houston, we have a problem".

**From the Earth to the Moon**

Apollo 13 (film) - Wikipedia

Home media. Apollo 13 was released in 1995 and on LaserDisc the following year.

Wikipedia

The Moonwalkers: A Journey with Studio

# Re-visiting some questions

The image displays two overlapping screenshots of Google search results. The left screenshot shows a search for "what year did tom hanks land on the moon". The AI Overview section states: "Tom Hanks has not landed on the moon, but he has been involved in several projects related to the Apollo space program:". It lists "Apollo 13 (1995)" where Hanks played an astronaut, and "From the Earth to the Moon". The right screenshot shows a search for "statue of liberty into the ocean". It lists results for "Apollo 13 (film) - Wikipedia", "The Moonwalkers: A Journey with Studio", "Statue of Liberty into the ocean, but it is taller than the Statue of Liberty.", "Saturn V Rocket - Kennedy Space Center", and "NASA Artemis | Taller than the Statue of Liberty Instagram".

# Re-visiting some questions

The image displays three overlapping screenshots of Google search results, each featuring an AI Overview section. The top-left screenshot shows a search for "what year did tom hanks land on the moon". The AI Overview states that Tom Hanks has not landed on the moon but was involved in the Apollo 13 mission. The top-right screenshot shows a search for "Statue of Liberty into the ocean". The AI Overview notes that the statue was not thrown into the ocean. The bottom screenshot shows a search for "when did the university of washington give an honorary degree to taylor swift". The AI Overview reports that NYU awarded Taylor Swift an honorary degree on May 18, 2022. Each screenshot also shows search filters and various search results snippets.

**Search 1: what year did tom hanks land on the moon**

AI Overview: **Tom Hanks has not landed on the moon**, but he has been involved in several projects related to the Apollo space program:

**Apollo 13 (1995)**  
Hanks played one of the three astronauts on NASA's Apollo 13 spacecraft in this movie based on the true story of the near-disaster. The movie is known for the line "Houston, we have a problem."

**From the Earth to the Moon**

**Statue of Liberty into the ocean**

AI Overview: **Statue of Liberty into the ocean**, but it was never thrown into the ocean.

**Search 2: when did the university of washington give an honorary degree to taylor swift**

AI Overview: **New York University (NYU) awarded Taylor Swift an honorary Doctor of Fine Arts degree on May 18, 2022:**

**When**  
The ceremony took place at Yankee Stadium as part of NYU's commencement for the class of 2022.

# Re-visiting some questions

Google search for "what year did tom hanks land on the moon".

AI Overview: Tom Hanks has not landed on the moon, but he has been involved in several projects related to the Apollo space program.

**Apollo 13 (1995)**  
Hanks played one of the three astronauts on NASA's Apollo 13 spacecraft in this movie based on the true story of the near-disaster. The movie is known for the line "Houston, we have a problem".

**From the Earth to the Moon**

Google search for "Statue of Liberty into the ocean".

AI Overview: Statue of Liberty into the ocean, but it is the tallest vehicle ever used.

**Apollo 13 (film) - Wikipedia**  
Home media. Apollo 13 was released in 1995 and on LaserDisc the following year.

**Saturn V Rocket - Kennedy Space Center**  
Description of the Saturn V Rocket ... It could carry up to 143,000 pounds to the Earth's orbit and about 50 tons to the Moon.

**Statue of Liberty is | Taller than the Statue of Liberty.**  
The Statue of Liberty is taller than the Statue of Liberty. In 1969, the Saturn V rocket for the Apollo 13 mission was taller than the Statue of Liberty.

Google search for "when did the university of washington give an honorary degree to Taylor Swift".

AI Overview: New York University (NYU) awarded Taylor Swift an honorary Doctor of Fine Arts degree.

Google search for "when did the university of washington give an honorary degree to Barack Obama".

AI Overview: I rolled my eyes at Taylor Swift at NYU.

May 20, 2022 — Looking at Taylor Swift's honorary degree at NYU is a reminder that there's more to her than just her music.

NBC News

Wikipedia: List of awards and honors received by Barack Obama

^ Barack Obama Receives Honorary Degree at Wesleyan University. spazeboydotnet. May 26, 2008. Archived from the original on October 18, 2020. Retrieved ...



# Conclusions

- Deep processing for QA
  - Exploits parsing, semantics, anaphora, reasoning
  - Computationally expensive
    - More tractable when applied only to questions and passages
- Systems trending toward greater use of:
  - Web resources: Wikipedia, answer repositories
  - Machine learning! Possible over-reliance on LLMs vs structured data (see also [Shah and Bender 2022](#))
- But still: real use of deep representations and processing thereof, even in the LLM era

# Next Time

# Next Time

- More on current directions (e.g. unsupervised learning)

# Next Time

- More on current directions (e.g. unsupervised learning)
- Summary / wrap-up

# Next Time

- More on current directions (e.g. unsupervised learning)
- Summary / wrap-up
- AMA / general discussion

# Next Time

- More on current directions (e.g. unsupervised learning)
- Summary / wrap-up
- AMA / general discussion
  - Submit questions here!

# Next Time

- More on current directions (e.g. unsupervised learning)
- Summary / wrap-up
- AMA / general discussion
  - Submit questions here!
  - <https://forms.gle/qiZ22sedhWXtyoGS9>

# Next Time

- More on current directions (e.g. unsupervised learning)
- Summary / wrap-up
- AMA / general discussion
  - Submit questions here!
  - <https://forms.gle/qiZ22sedhWXtyoGS9>
- Course evaluation!



# Bonus Slides: Neural Approaches to Coreference and WSD

# End-to-End Neural Coreference Resolution

[Lee et al., 2017](#)

# End-to-End Neural Coreference Resolution

[Lee et al, 2017](#)

- Begin with dataset with gold mention clusters (aka chains)

*“General Electric said the Postal Service contacted the company.”*

# End-to-End Neural Coreference Resolution

[Lee et al, 2017](#)

- Begin with dataset with gold mention clusters (aka chains)

*“General Electric said the Postal Service contacted the company.”*



# End-to-End Neural Coreference Resolution

Lee et al, 2017

- Can think of the coref problem as finding the maximally likely distribution:

$$P(y_1, \dots, y_N | D) = \prod_{i=1}^N \frac{\exp(s(i, y_i))}{\sum_{y' \in y(i)} \exp(s(i, y'))}$$

# End-to-End Neural Coreference Resolution

[Lee et al, 2017](#)

- Can think of the coref problem as finding the maximally likely distribution:

$$P(y_1, \dots, y_N | D) = \prod_{i=1}^N \frac{\exp(s(i, y_i))}{\sum_{y' \in y(i)} \exp(s(i, y'))}$$

Where

$$s(i, j) = \begin{cases} 0 & j = \epsilon \\ s_m(i) + s_m(j) + s_a(i, j) & j \neq \epsilon \end{cases}$$

# End-to-End Neural Coreference Resolution

[Lee et al, 2017](#)

- Can think of the coref problem as finding the maximally likely distribution:

$$P(y_1, \dots, y_N | D) = \prod_{i=1}^N \frac{\exp(s(i, y_i))}{\sum_{y' \in y(i)} \exp(s(i, y'))}$$

Where

$$s(i, j) = \begin{cases} 0 & j = \epsilon \\ s_m(i) + s_m(j) + s_a(i, j) & j \neq \epsilon \end{cases}$$

Coref Score

# End-to-End Neural Coreference Resolution

Lee et al, 2017

- Can think of the coref problem as finding the maximally likely distribution:

$$P(y_1, \dots, y_N | D) = \prod_{i=1}^N \frac{\exp(s(i, y_i))}{\sum_{y' \in y(i)} \exp(s(i, y'))}$$

Where

$$s(i, j) = \begin{cases} 0 & j = \epsilon \\ s_m(i) + s_m(j) + s_a(i, j) & j \neq \epsilon \end{cases}$$

Coref Score

Mention Score



# End-to-End Neural Coreference Resolution

Lee et al, 2017

- Can think of the coref problem as finding the maximally likely distribution:

$$P(y_1, \dots, y_N | D) = \prod_{i=1}^N \frac{\exp(s(i, y_i))}{\sum_{y' \in y(i)} \exp(s(i, y'))}$$

Where

$$s(i, j) = \begin{cases} 0 & j = \epsilon \\ s_m(i) + s_m(j) + s_a(i, j) & j \neq \epsilon \end{cases}$$

Coref Score

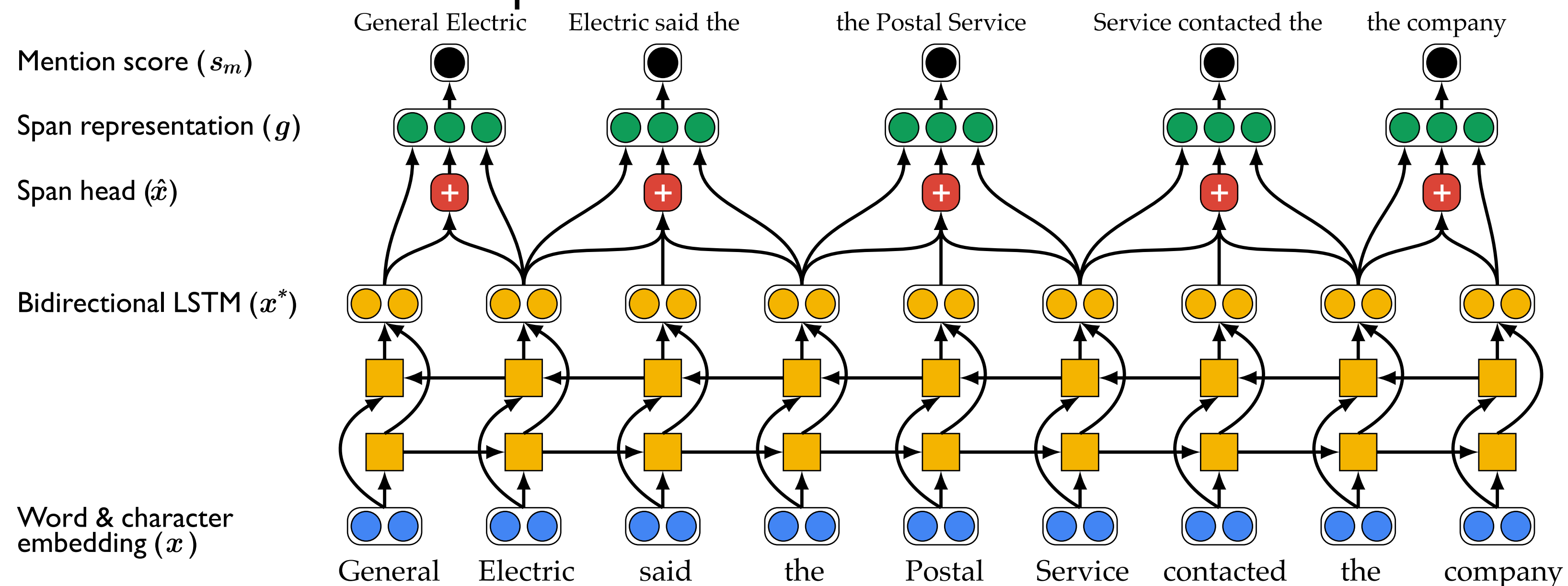
Mention Score

Antecedent Score

# End-to-End Neural Coreference Resolution

[Lee et al, 2017](#)

- **Step 1** — Train model to identify spans based on gold span labels
  - Use bi-LSTMs to model sequential information preceding/following/within spans
  - Include “headedness” of span with a learned *attention* mechanism



# End-to-End Neural Coreference Resolution

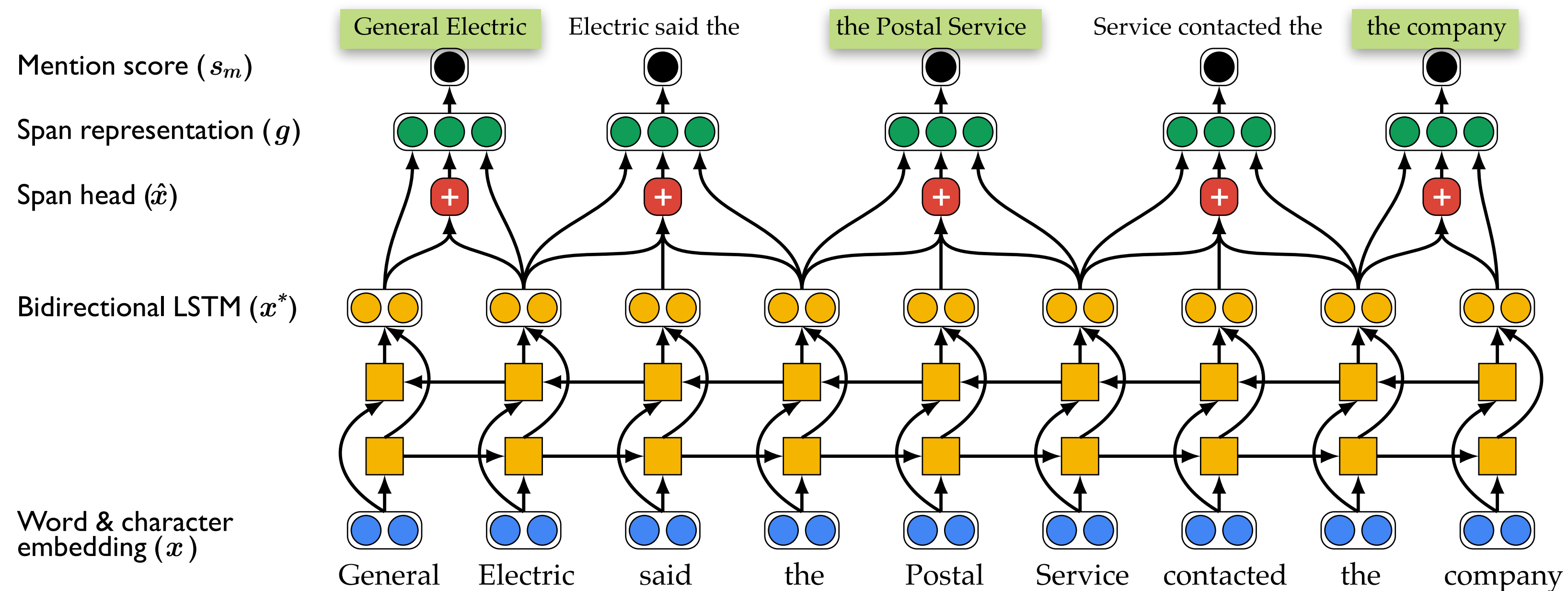
Lee et al, 2017

- **Attention** can be visualized by heatmap over spans:
  - (The flight attendants) have until 6:00 today to ratify labor concessions. (The pilots') union and ground crew did so yesterday.
  - (Prince Charles and his new wife Camilla) have jumped across the pond and are touring the United States making (their) first stop today in New York. It's Charles' first opportunity to showcase his new wife, but few Americans seem to care. Here's Jeanie Mowth. What a difference two decades make. (Charles and Diana) visited a JC Penney's on the prince's last official US tour. Twenty years later, here's the prince with his new wife.

# End-to-End Neural Coreference Resolution

[Lee et al, 2017](#)

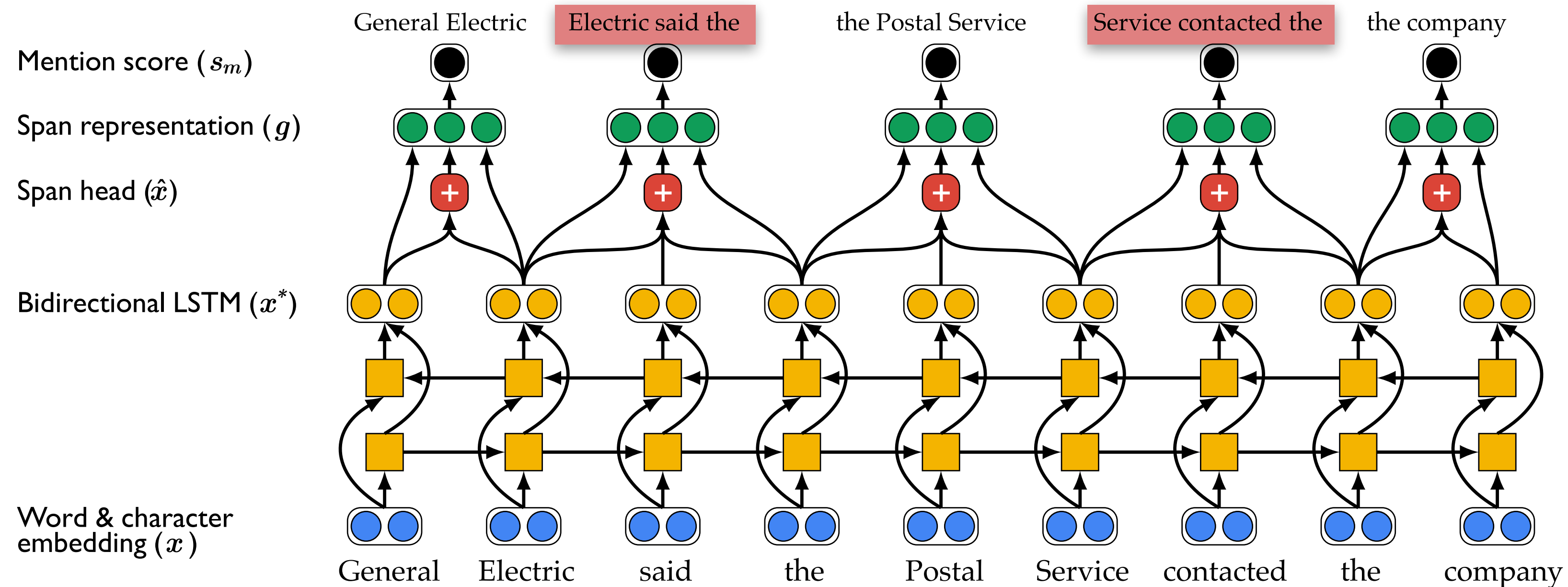
- **These** are valid gold mentions (network gets “reward” for getting these right)



# End-to-End Neural Coreference Resolution

[Lee et al, 2017](#)

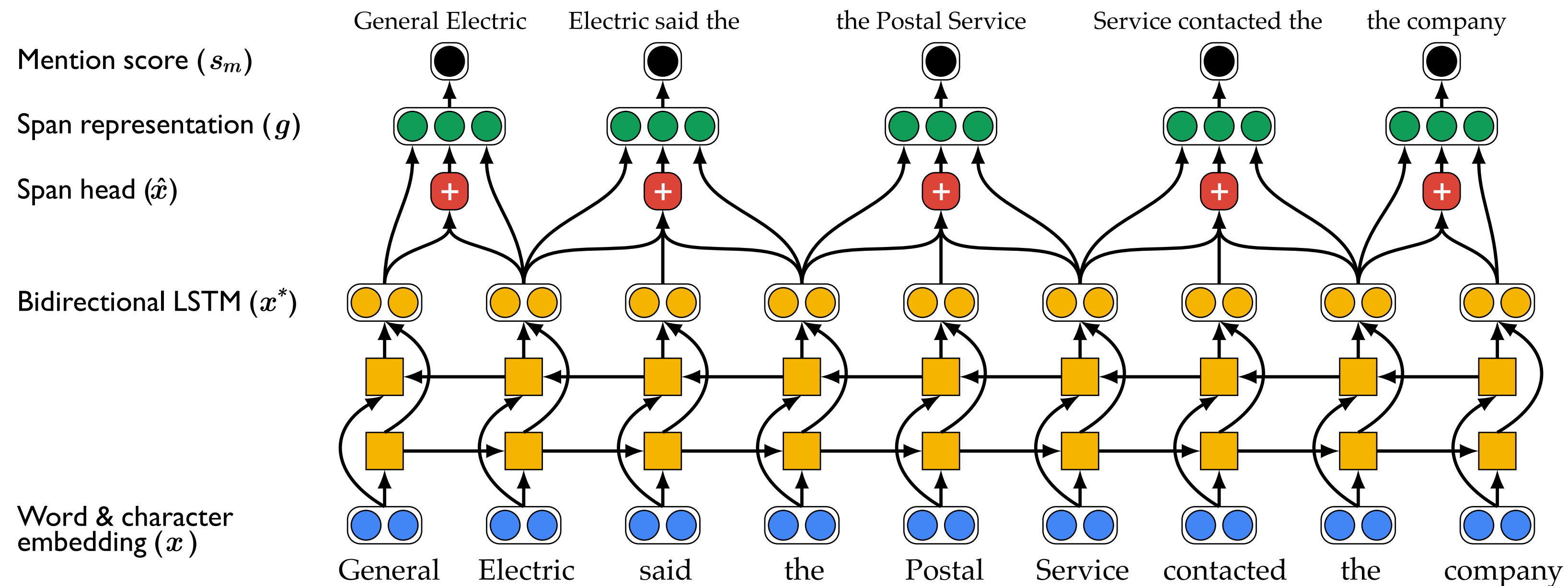
- **These** are invalid mentions (network accumulates error if these are selected)



# End-to-End Neural Coreference Resolution

[Lee et al, 2017](#)

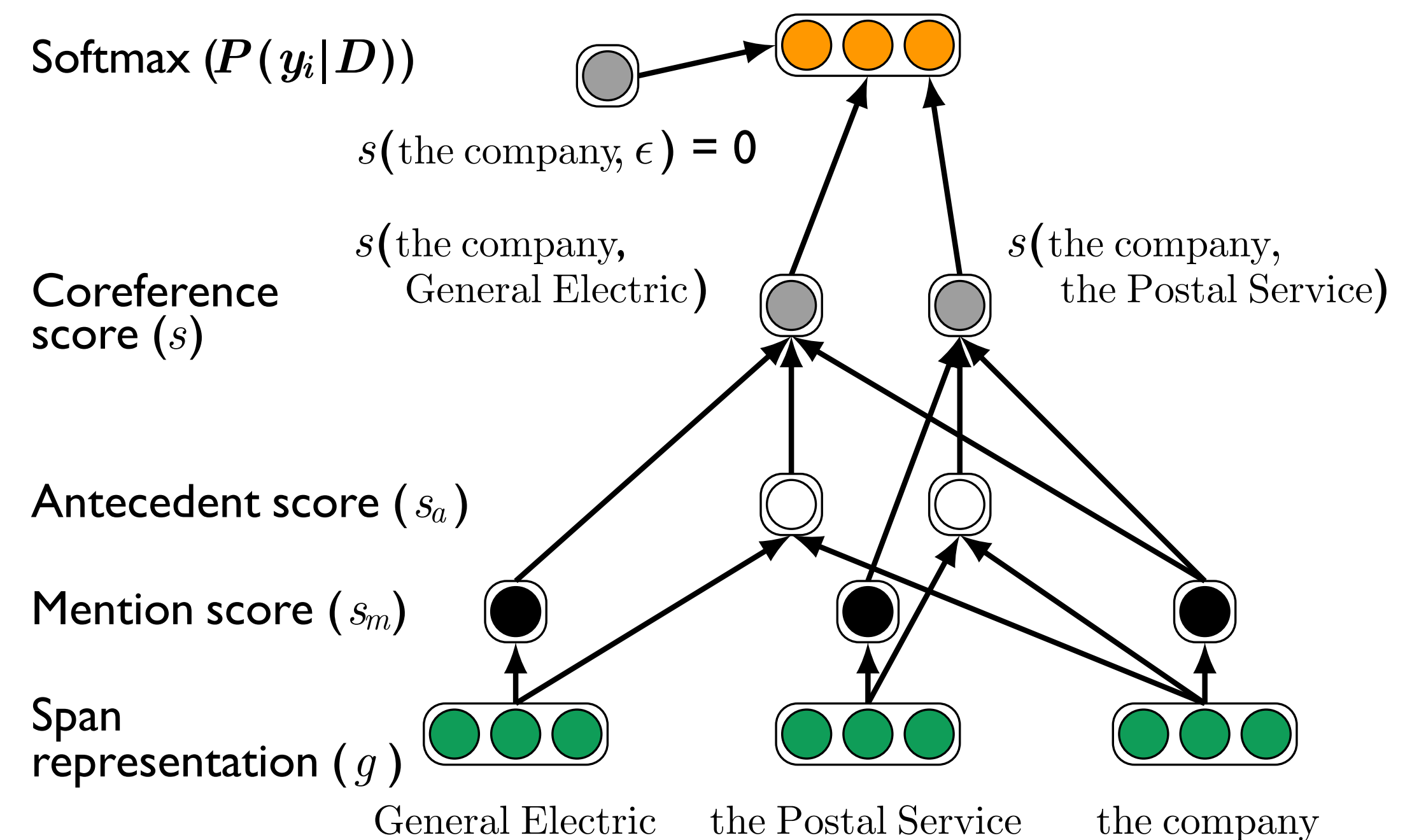
- Network thus learns to identify features from (embeddings → sequence) + head
- As more or less likely to identify a span of words as a mention



# End-to-End Neural Coreference Resolution

[Lee et al, 2017](#)

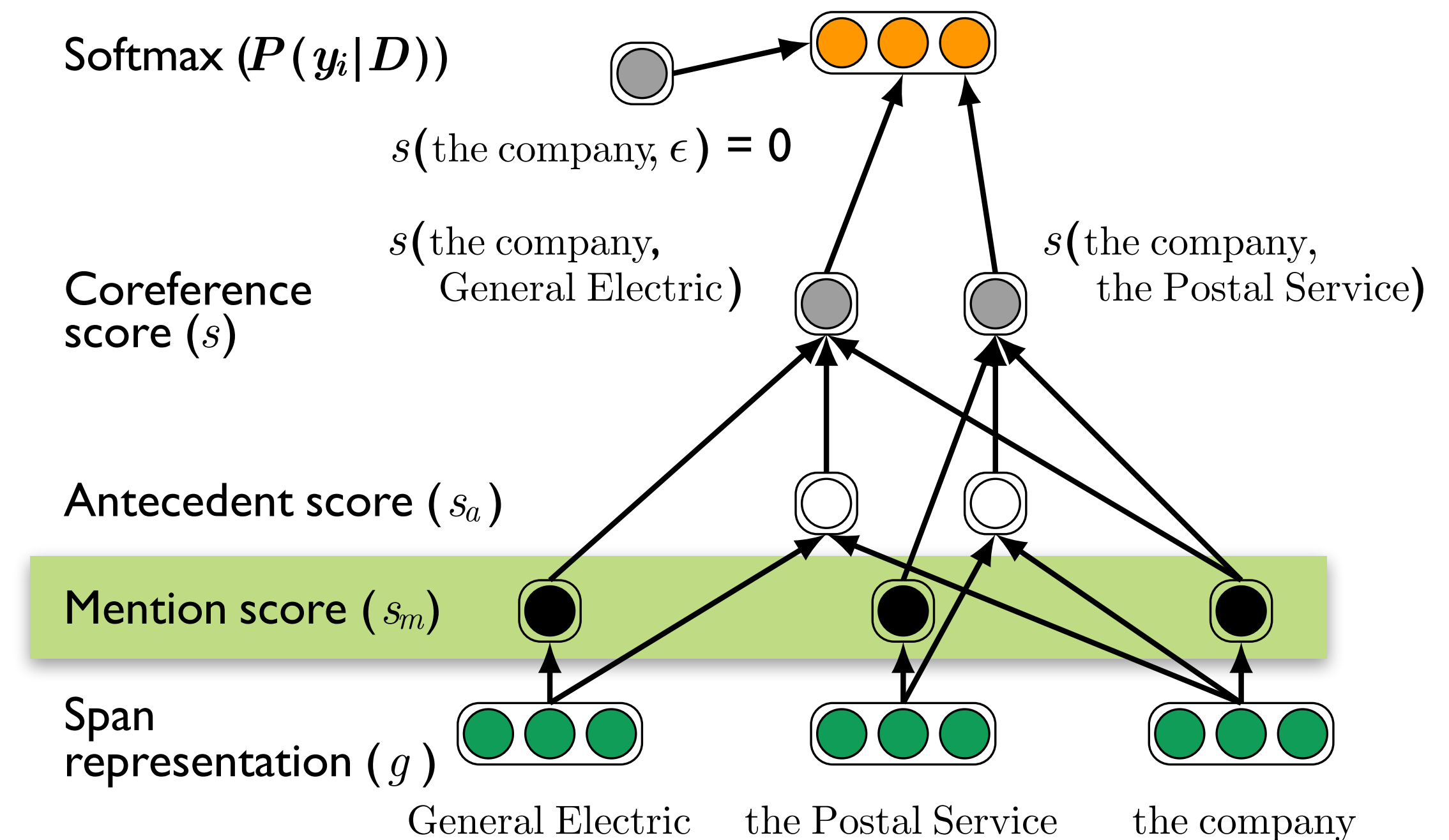
- **Step 2** — Learn Coref Clusters



# End-to-End Neural Coreference Resolution

[Lee et al, 2017](#)

- **Step 2** — Learn Coref Clusters
- **Mention Scores**
  - Likelihood a given span is a mention
  - Unary over spans

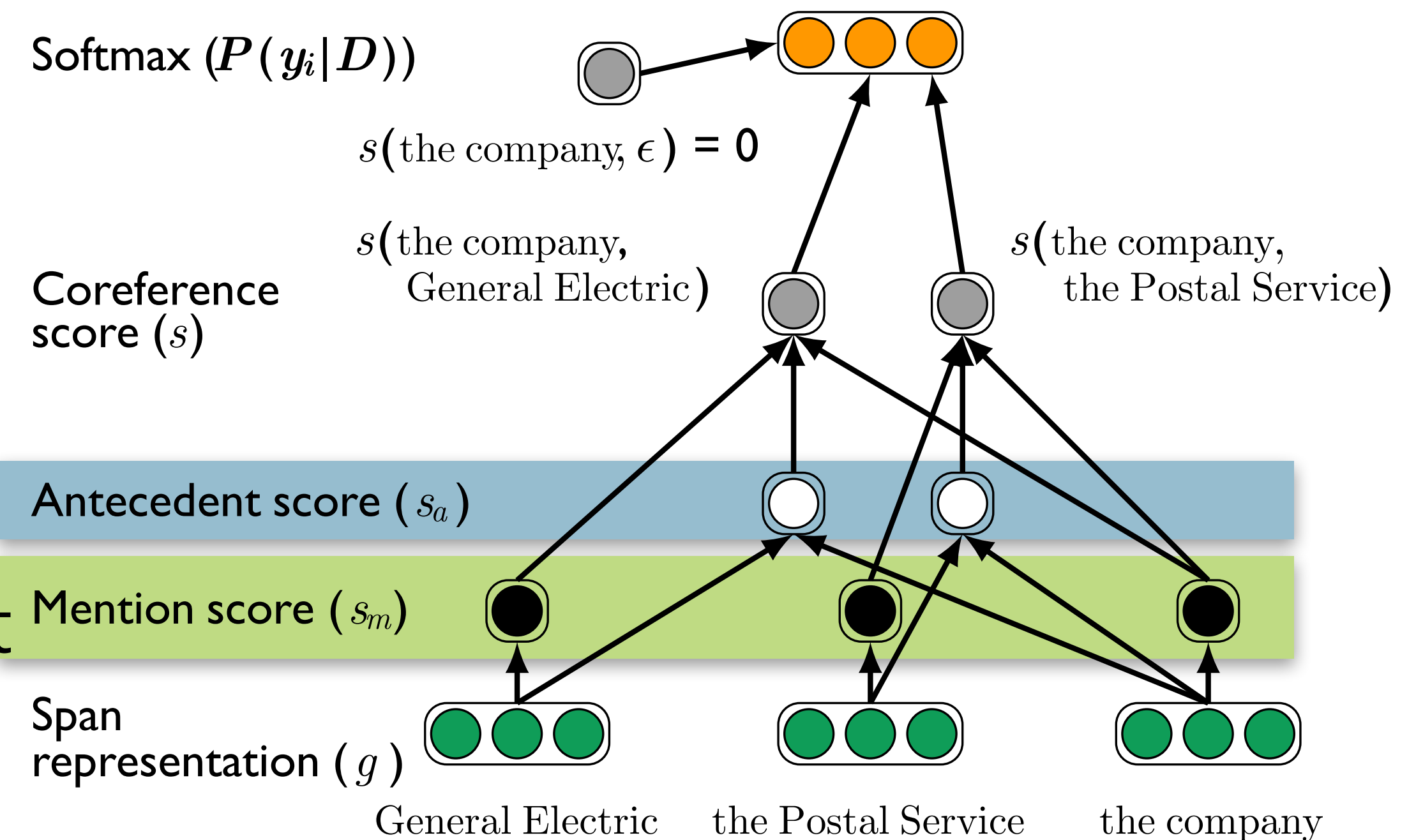




# End-to-End Neural Coreference Resolution

[Lee et al, 2017](#)

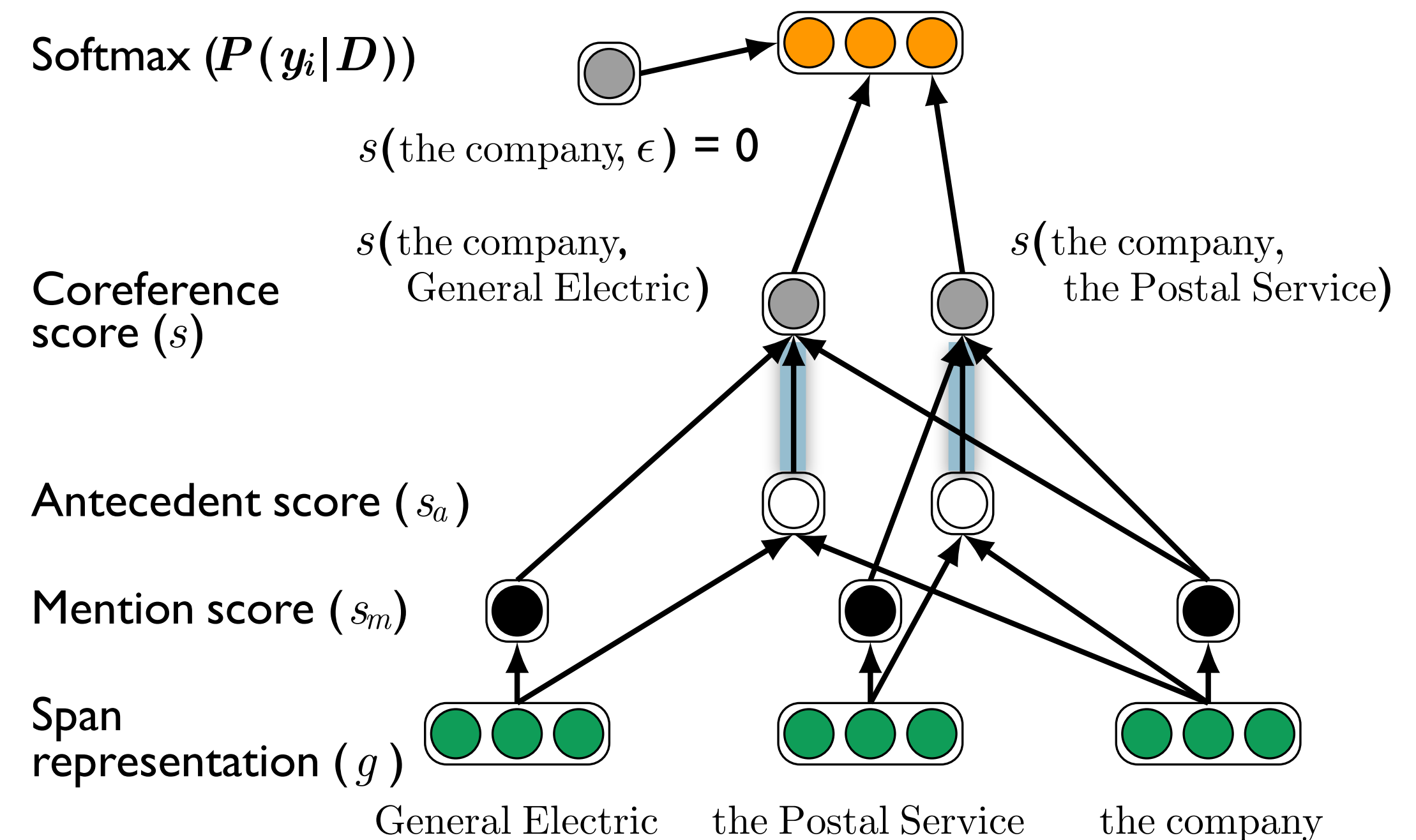
- **Step 2** — Learn Coref Clusters
- **Mention Scores**
  - Likelihood a given span is a mention
  - Unary over spans
- **Antecedent scores**
  - Likelihood another mention is antecedent
  - Pairwise between spans



# End-to-End Neural Coreference Resolution

[Lee et al, 2017](#)

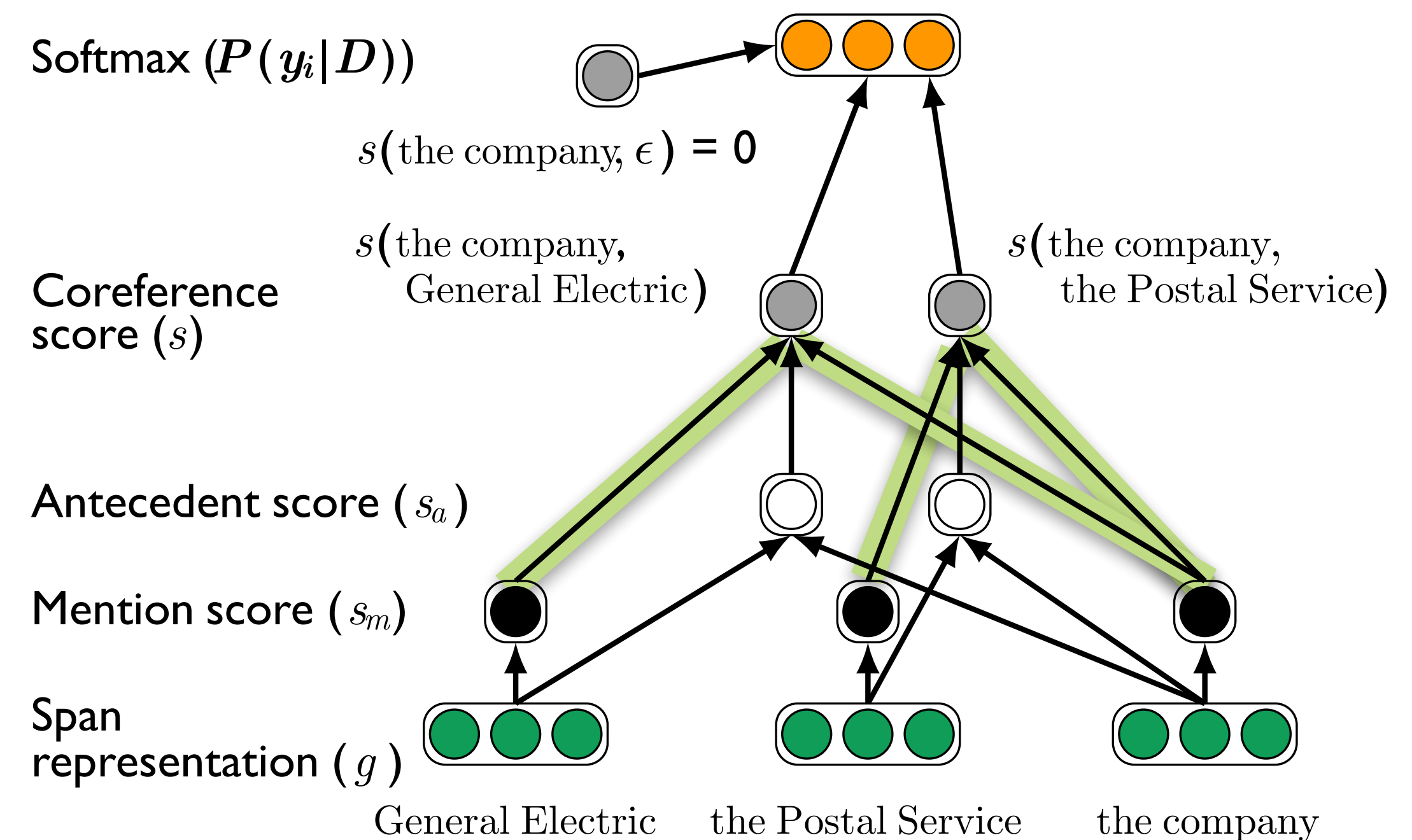
- The coref score is a combination of:
  - antecedent scores



# End-to-End Neural Coreference Resolution

[Lee et al, 2017](#)

- The coref score is a combination of:
  - antecedent scores
  - mention scores



# End-to-End Neural Coreference Resolution

[Lee et al, 2017](#)

- Other info:
  - Also implement pruning to avoid dealing with *all* spans
  - Also encode metadata, such as speaker and genre in mention representation

# End-to-End Neural Coreference Resolution

Lee et al, 2017

- Data:
  - [CoNLL-2012 Shared Task](#) (Coref on OntoNotes)
  - **2802** training docs
  - **343** development docs
  - **348** test docs
  - 454 words/doc average

# End-to-End Neural Coreference Resolution

[Lee et al, 2017](#)

- Positive:
  - State-of-the-art on CoNLL-2012 Test Data
- Errors:
  - Word embeddings tend to conflate paraphrasing with relatedness
    - e.g. (The flight attendants) have until 6:00 today to ratify labor concessions. (The pilots') union and ground crew did so yesterday.
  - (Prince Charles and his new wife Camilla) have jumped across the pond ... What a difference two decades make. (Charles and Diana) visited a JC Penney's on the Prince's last official US tour. ...

# Neural Sequence Learning Models for Word Sense Disambiguation

[Raganato et. al \(2017b\)](#)

# Neural Sequence Learning Models for WSD

[Raganato et. al \(2017b\)](#)

- Authors propose several models for encoding words and senses
  - **bi-LSTM**
  - **bi-LSTM + Attention**
  - **Sequence to Sequence**
- All approaches are encoding sequential information
- All approaches use sense-tagged corpus

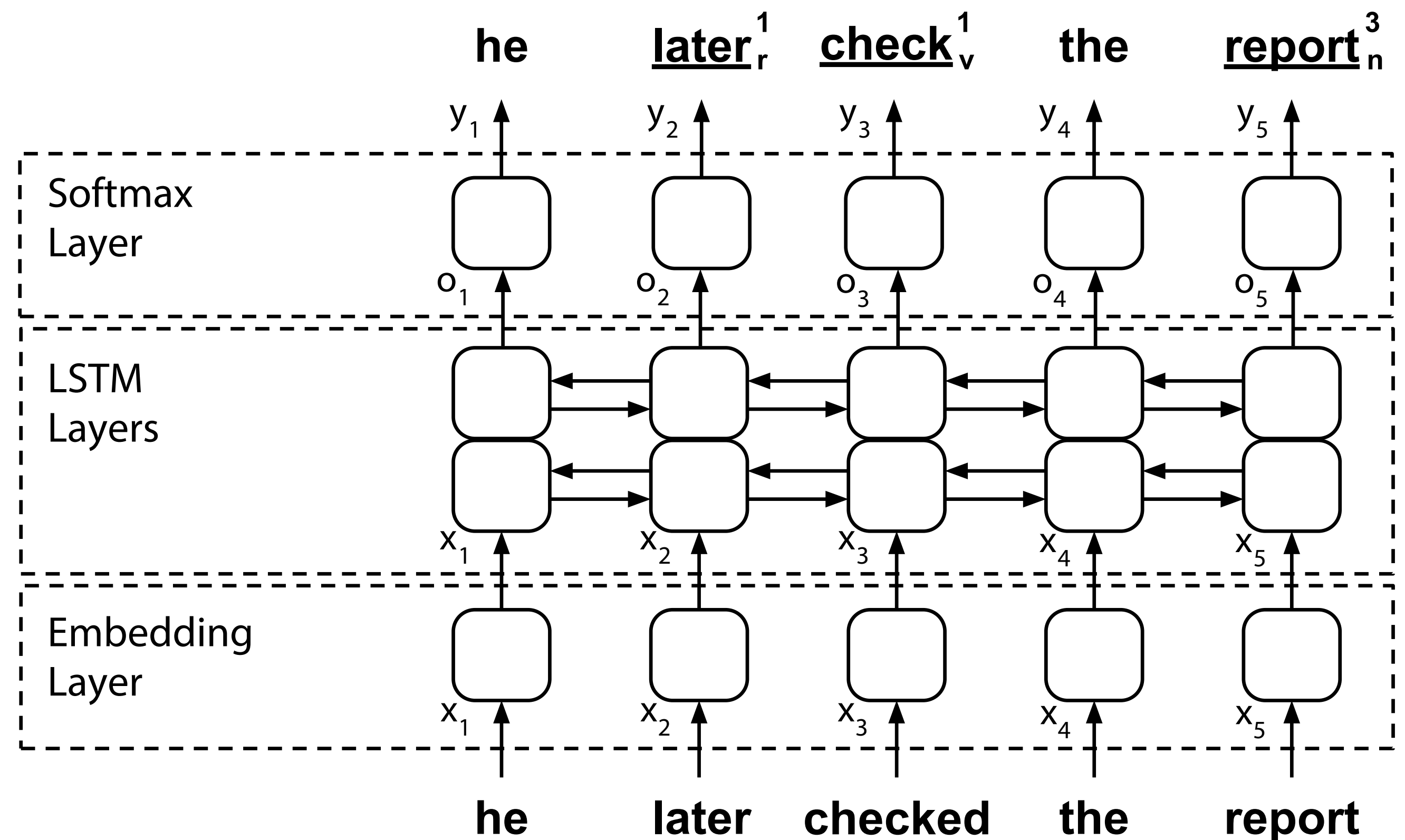


# Neural Sequence Learning Models for WSD

[Raganato et. al \(2017b\)](#)

- **bi-LSTM**

- Learn to label proper sense given word embedding and context (LSTM)



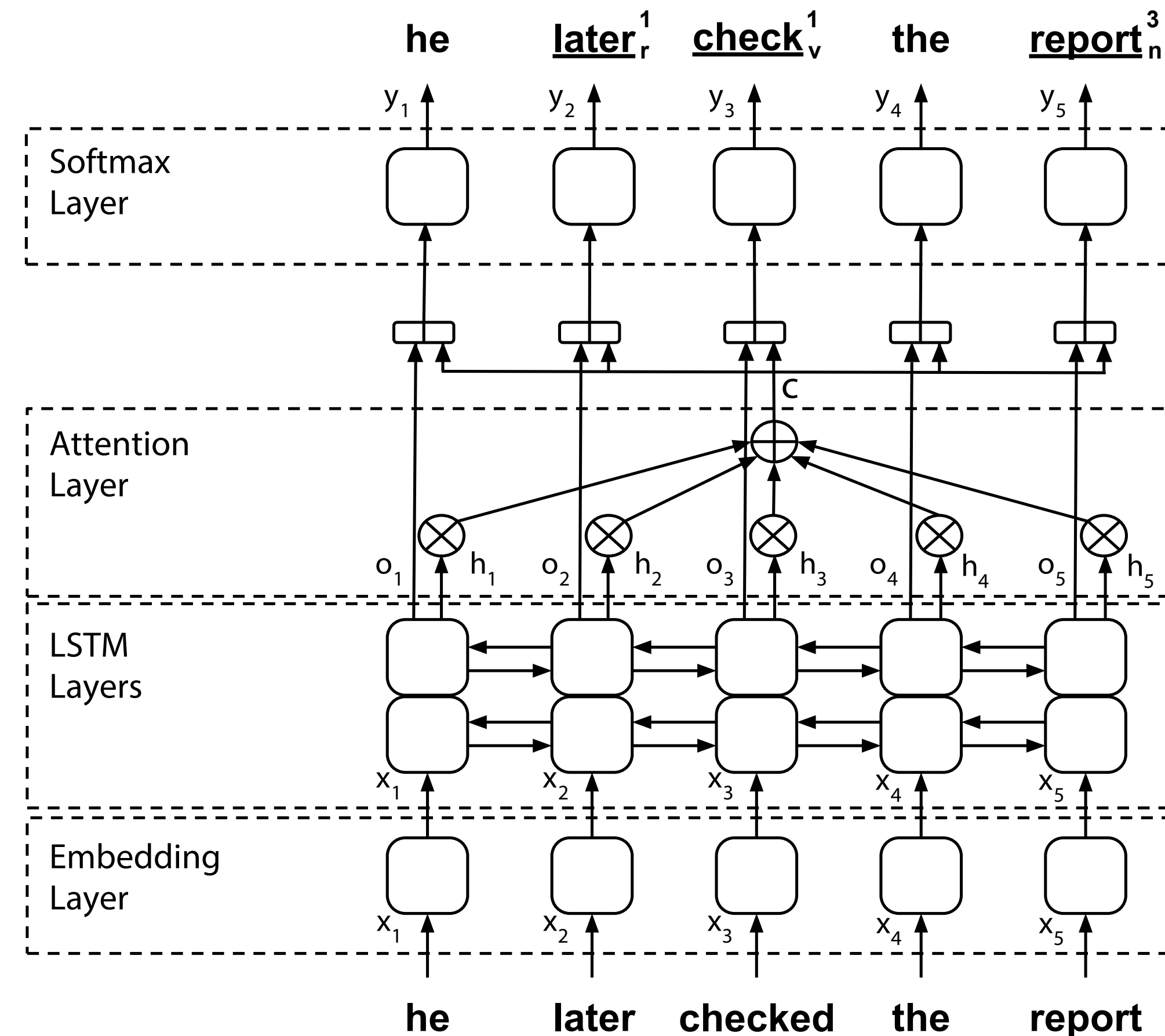
# Neural Sequence Learning Models for WSD

[Raganato et. al \(2017b\)](#)

- **bi-LSTM + Attention**

- Attention layer adds sentence-level representation  $c$  to guide the labels generate at each sequence time step by focusing on what part of the sentence may be relevant

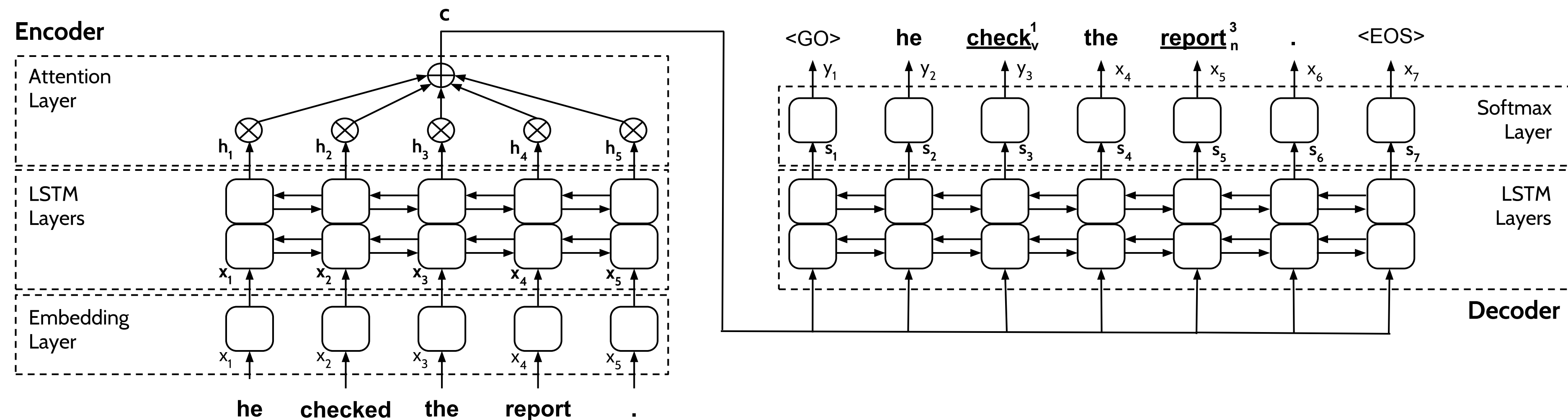
- (e.g. with *wicket* in focus, *match* might be influenced toward the game sense, rather than firestarter)



# Neural Sequence Learning Models for WSD

[Raganato et. al \(2017b\)](#)

- **seq2seq**
  - Two-step task:
    - Memorization — Model is trained to replicate input token-by-token
    - Disambiguation — Model learns to replace surface forms with appropriate senses



# Neural Sequence Learning Models for WSD

[Raganato et. al \(2017b\)](#)

- Also try models that jointly learn WSD and:
  - coarse semantic labels
    - e.g. *noun.location, verb.motion*
  - POS tags
  - Both

# Neural Sequence Learning Models for WSD

[Raganato et. al \(2017b\)](#)

- Data:
  - Use SemCor 3.0 for training/evaluating word senses

# Neural Sequence Learning Models for WSD

Raganato et. al (2017b)

- Results:

	Dev	Test Datasets				Concatenation of All Test Datasets				
	SE07	SE2	SE3	SE13	SE15	Nouns	Verbs	Adj.	Adv.	All
BLSTM	61.8	71.4	68.8	65.6	69.2	70.2	56.3	75.2	<b>84.4</b>	68.9
BLSTM + att.	62.4	71.4	<b>70.2</b>	66.4	70.8	71.0	<b>58.4</b>	75.2	83.5	69.7
BLSTM + att. $\mathbb{L}EX$	63.7	<b>72.0</b>	69.4	66.4	<b>72.4</b>	<b>71.6</b>	57.1	<b>75.6</b>	83.2	<b>69.9</b>
BLSTM + att. $\mathbb{L}EX + POS$	<b>64.8</b>	<b>72.0</b>	69.1	<b>66.9</b>	71.5	71.5	57.5	75.0	83.8	<b>69.9</b>
Seq2Seq	60.9	68.5	67.9	65.3	67.0	68.7	54.5	74.0	81.2	67.3
Seq2Seq + att.	62.9	69.9	69.6	65.6	67.7	69.5	57.2	74.5	81.8	68.4
Seq2Seq + att. $\mathbb{L}EX$	64.6	70.6	67.8	66.5	68.7	70.4	55.7	73.3	82.9	68.5
Seq2Seq + att. $\mathbb{L}EX + POS$	63.1	70.1	68.5	66.5	69.2	70.1	55.2	75.1	84.4	68.6
IMS	61.3	70.9	69.3	65.3	69.5	70.5	55.8	75.6	82.9	68.9
IMS+emb	<b>62.6</b>	<b>72.2</b>	<b>70.4</b>	65.9	71.5	<b>71.9</b>	56.6	<b>75.9</b>	<b>84.7</b>	<b>70.1</b>
Context2Vec	61.3	71.8	69.1	65.6	<b>71.9</b>	71.2	<b>57.4</b>	75.2	82.7	69.6
Lesk <sub>ext</sub> + emb	56.7	63.0	63.7	66.2	64.6	70.0	51.1	51.7	80.6	64.2
UKB <sub>gloss</sub> w2w	42.9	63.5	55.4	62.9	63.3	64.9	41.4	69.5	69.7	61.1
Babelfy	51.6	67.0	63.5	<b>66.4</b>	70.3	68.9	50.7	73.2	79.8	66.4
MFS	54.5	65.6	66.0	63.8	67.1	67.7	49.8	73.1	80.5	65.5

# Neural Sequence Learning Models for WSD

[Raganato et. al \(2017b\)](#)

- Analysis:
  - Comparable to other supervised systems
  - Adding coarse-grained lexical tags appears to help
  - POS did not seem to help
- ***None of these systems substantially better than using the Most Frequent Sense***