

# Lecture 18: Wrapup + Ethics

Alan Ritter

(many slides from Greg Durrett)

# Administrivia

---

- ▶ Final project reports due Friday 12/9
- ▶ Please fill out the course/instructor opinion survey (CIOS) if you haven't already!

# This Lecture

---

- ▶ Wrapup Question Answering
- ▶ Ethics in NLP

# Span-based Question Answering

# SQuAD

---

- ▶ Single-document, single-sentence question-answering task where the answer is always a substring of the passage
- ▶ Predict start and end indices of the answer in the passage

## Passage

Super Bowl 50 was an American football game to determine the champion of the National Football League (NFL) for the 2015 season. The American Football Conference (AFC) champion Denver Broncos defeated the National Football Conference (NFC) champion Carolina Panthers 24–10 to earn their third Super Bowl title. The game was played on February 7, 2016, at Levi's Stadium in the San Francisco Bay Area at Santa Clara, California.

**Question:** Which NFL team won Super Bowl 50?

**Answer:** Denver Broncos

**Question:** What does AFC stand for?

**Answer:** American Football Conference

**Question:** What year was Super Bowl 50?

**Answer:** 2016

# SQuAD 2.0

---

- ▶ SQuAD 1.1 contains 100k+ QA pairs from 500+ Wikipedia articles.
- ▶ SQuAD 2.0 includes additional 50k questions that cannot be answered.
- ▶ These questions were crowdsourced.

## Passage

Super Bowl 50 was an American football game to determine the champion of the National Football League (NFL) for the 2015 season. The American Football Conference (AFC) champion Denver Broncos defeated the National Football Conference (NFC) champion Carolina Panthers 24–10 to earn their third Super Bowl title. The game was played on February 7, 2016, at Levi's Stadium in the San Francisco Bay Area at Santa Clara, California.

**Question:** Which NFL team won Super Bowl 50?

**Answer:** Denver Broncos

**Question:** What does AFC stand for?

**Answer:** American Football Conference

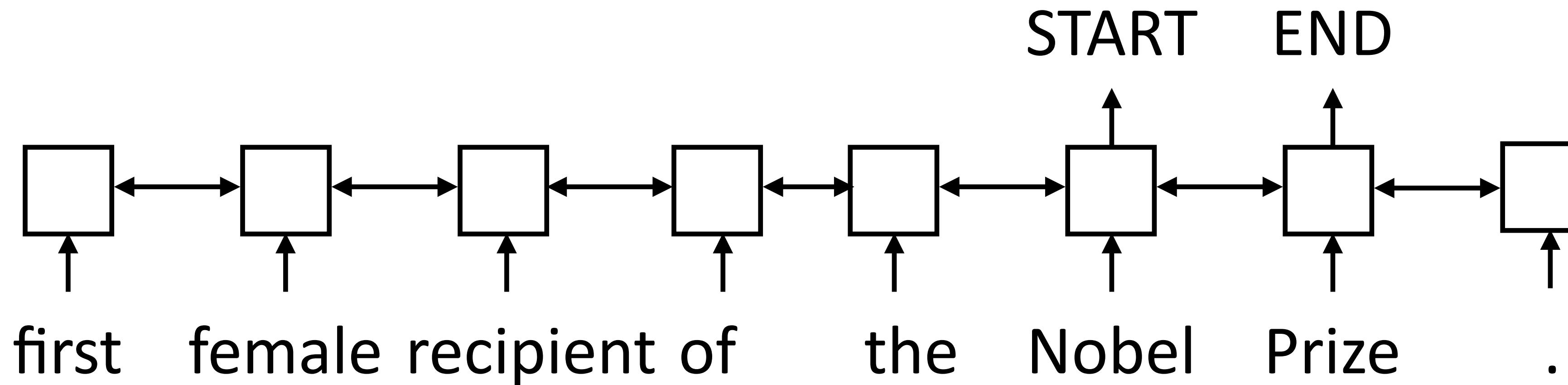
**Question:** What year was Super Bowl 50?

**Answer:** 2016

# SQuAD

---

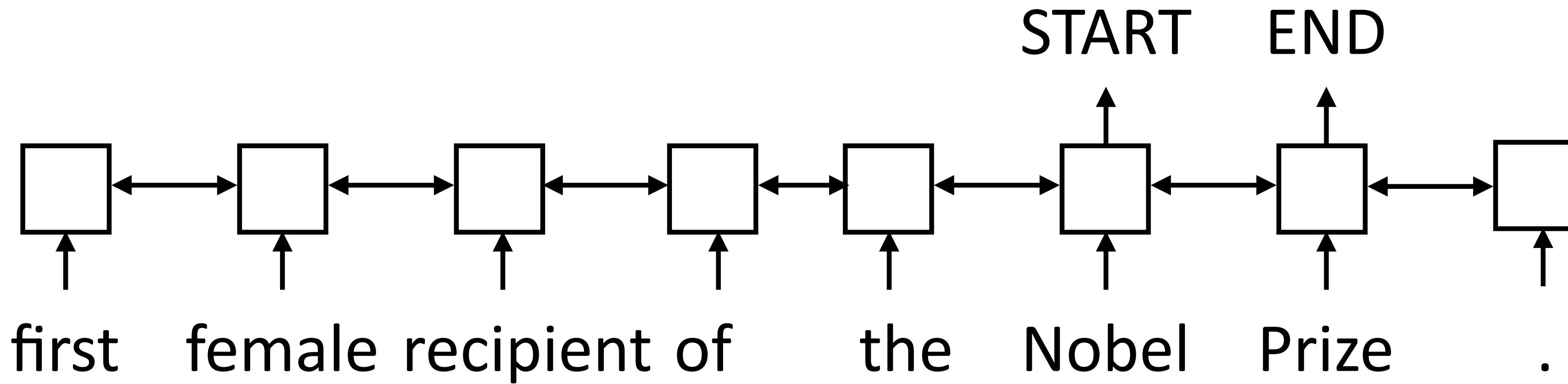
Q: What was Marie Curie the first female recipient of?



# SQuAD

---

Q: What was Marie Curie the first female recipient of?



- ▶ Like a tagging problem over the sentence (not multiclass classification), but we need some way of attending to the query

# Why did this take off?

---

- ▶ SQuAD was **big**: >100,000 questions (written by human) at a time when deep learning was exploding
- ▶ SQuAD had **room to improve**: ~50% performance from a logistic regression baseline (classifier with 180M features over constituents)
- ▶ SQuAD was **pretty easy**: year-over-year progress for a few years until the dataset was essentially solved

# Bidirectional Attention Flow (BiDAF)

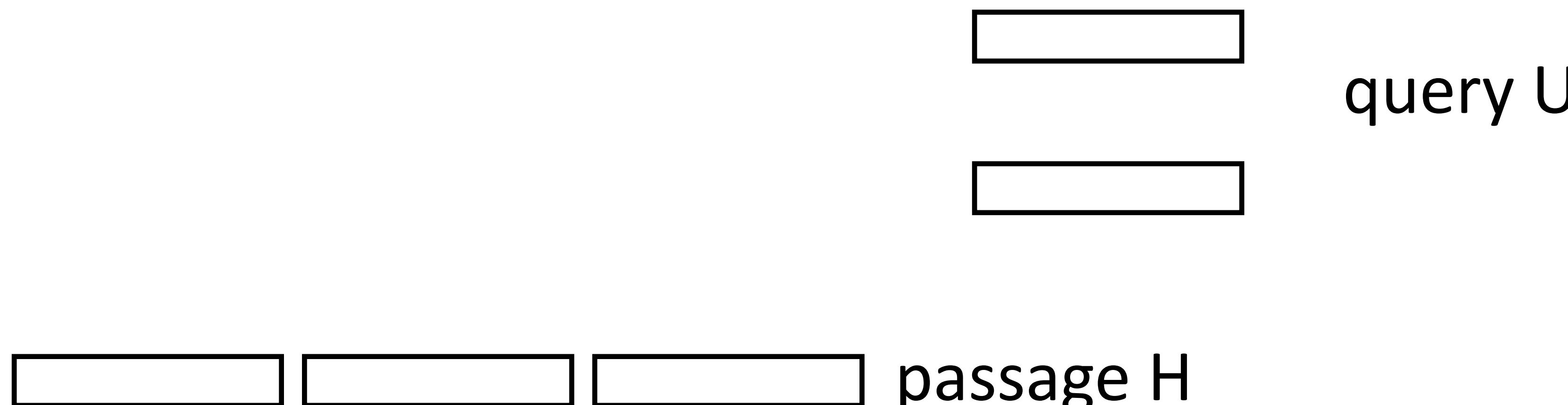
---

- ▶ Passage (context) and query are both encoded with BiLSTMs

# Bidirectional Attention Flow (BiDAF)

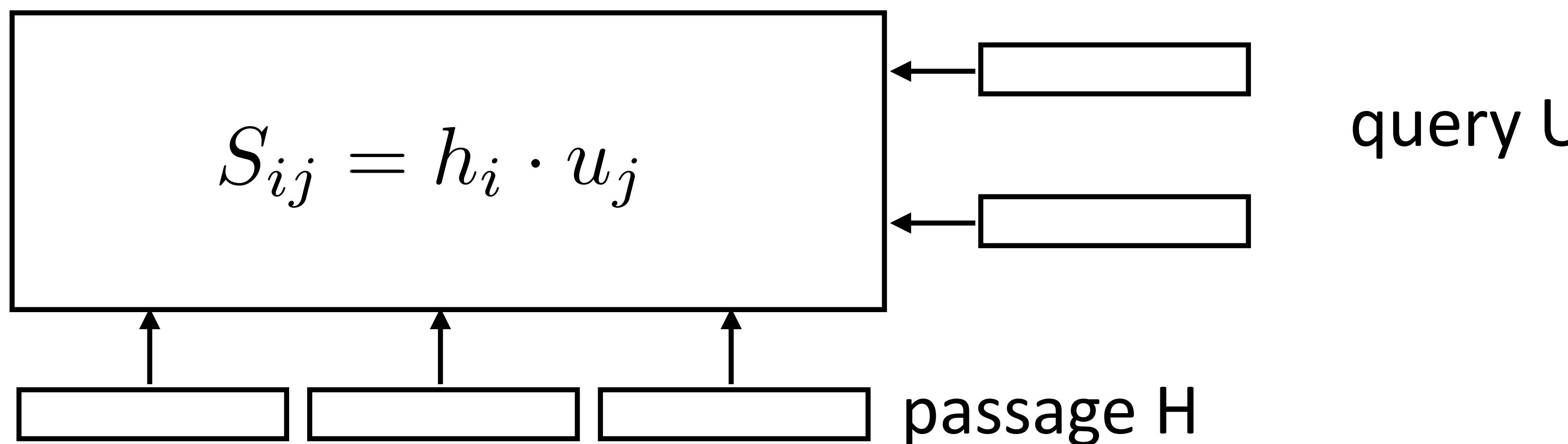
---

- ▶ Passage (context) and query are both encoded with BiLSTMs



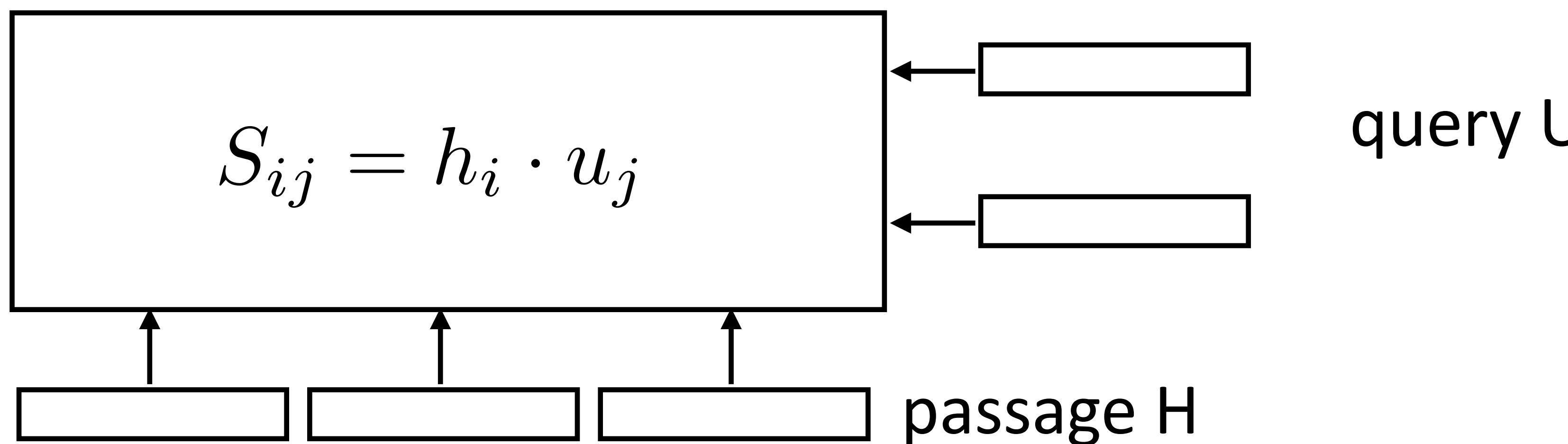
# Bidirectional Attention Flow (BiDAF)

- ▶ Passage (context) and query are both encoded with BiLSTMs



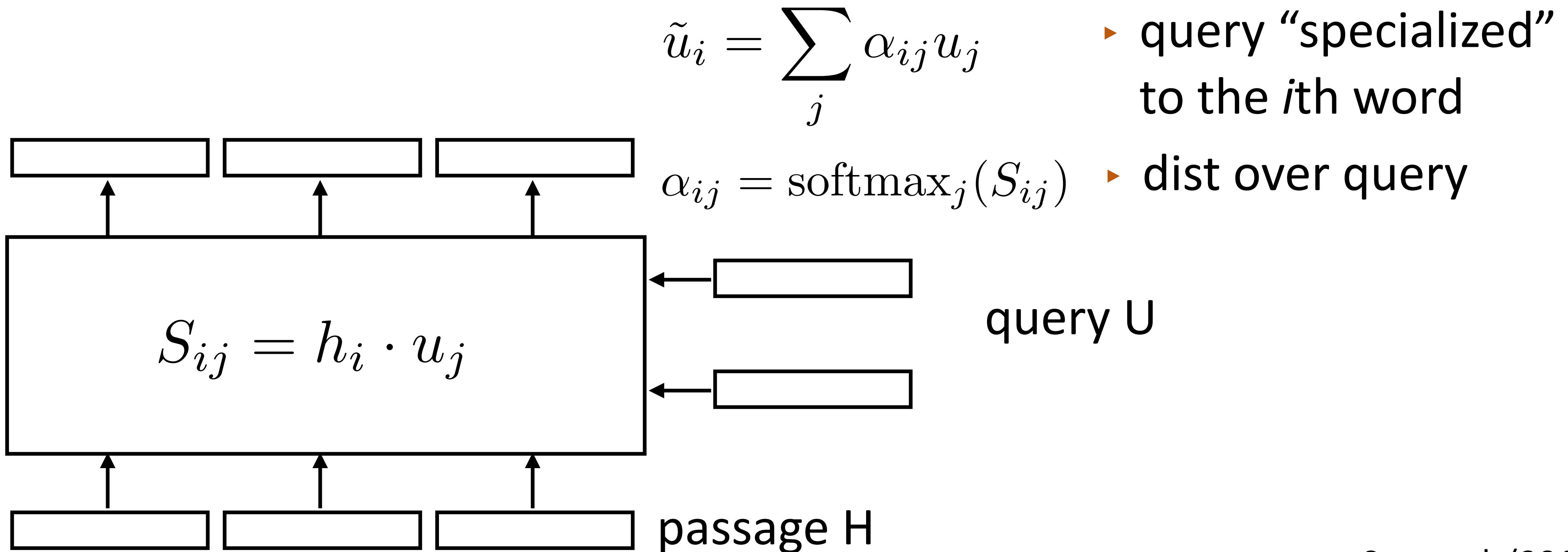
# Bidirectional Attention Flow (BiDAF)

- ▶ Passage (context) and query are both encoded with BiLSTMs
- ▶ Context-to-query attention: compute softmax over columns of  $S$ , take weighted sum of  $u$  based on attention weights for each passage word



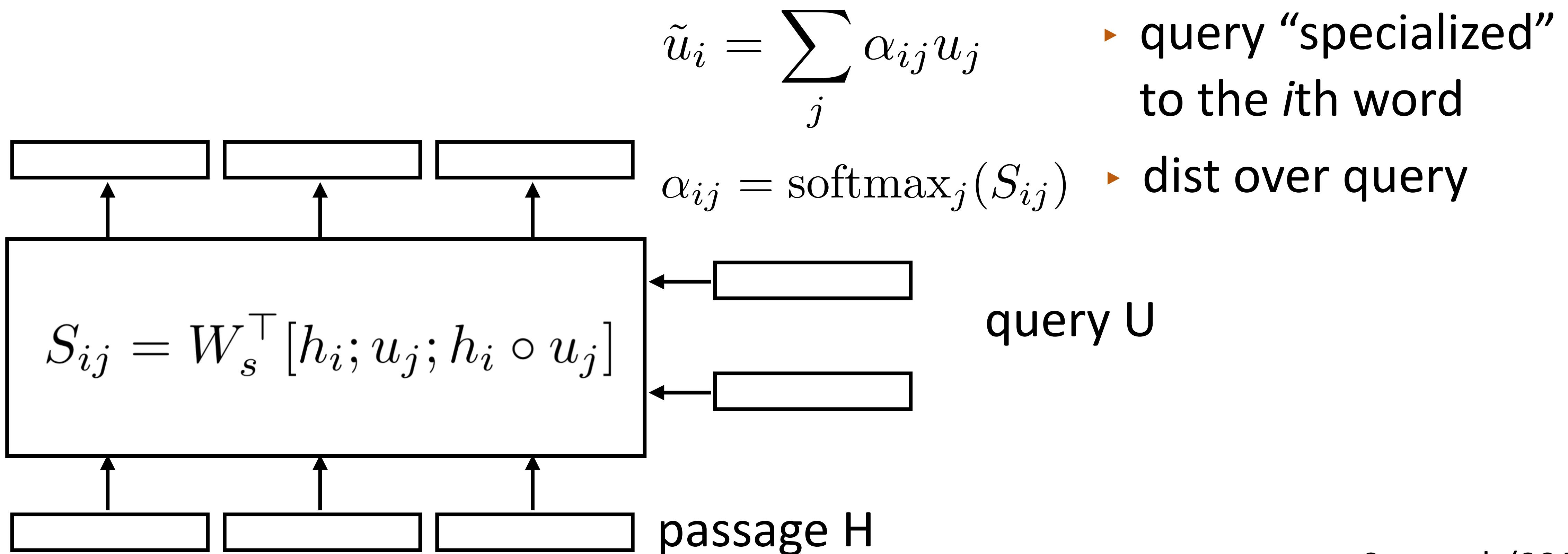
# Bidirectional Attention Flow (BiDAF)

- ▶ Passage (context) and query are both encoded with BiLSTMs
- ▶ Context-to-query attention: compute softmax over columns of  $S$ , take weighted sum of  $u$  based on attention weights for each passage word

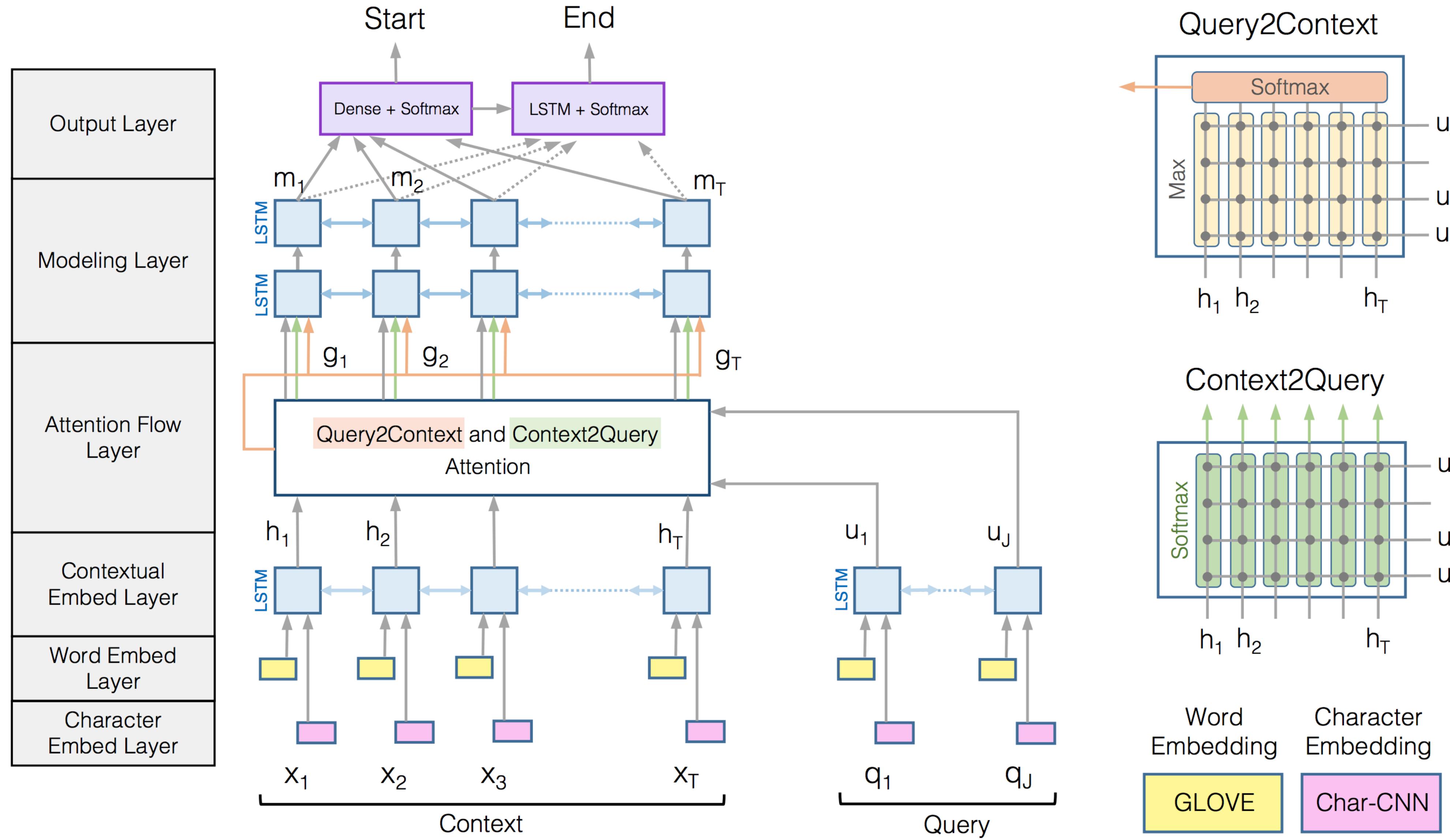


# Bidirectional Attention Flow (BiDAF)

- ▶ Passage (context) and query are both encoded with BiLSTMs
- ▶ Context-to-query attention: compute softmax over columns of  $S$ , take weighted sum of  $u$  based on attention weights for each passage word

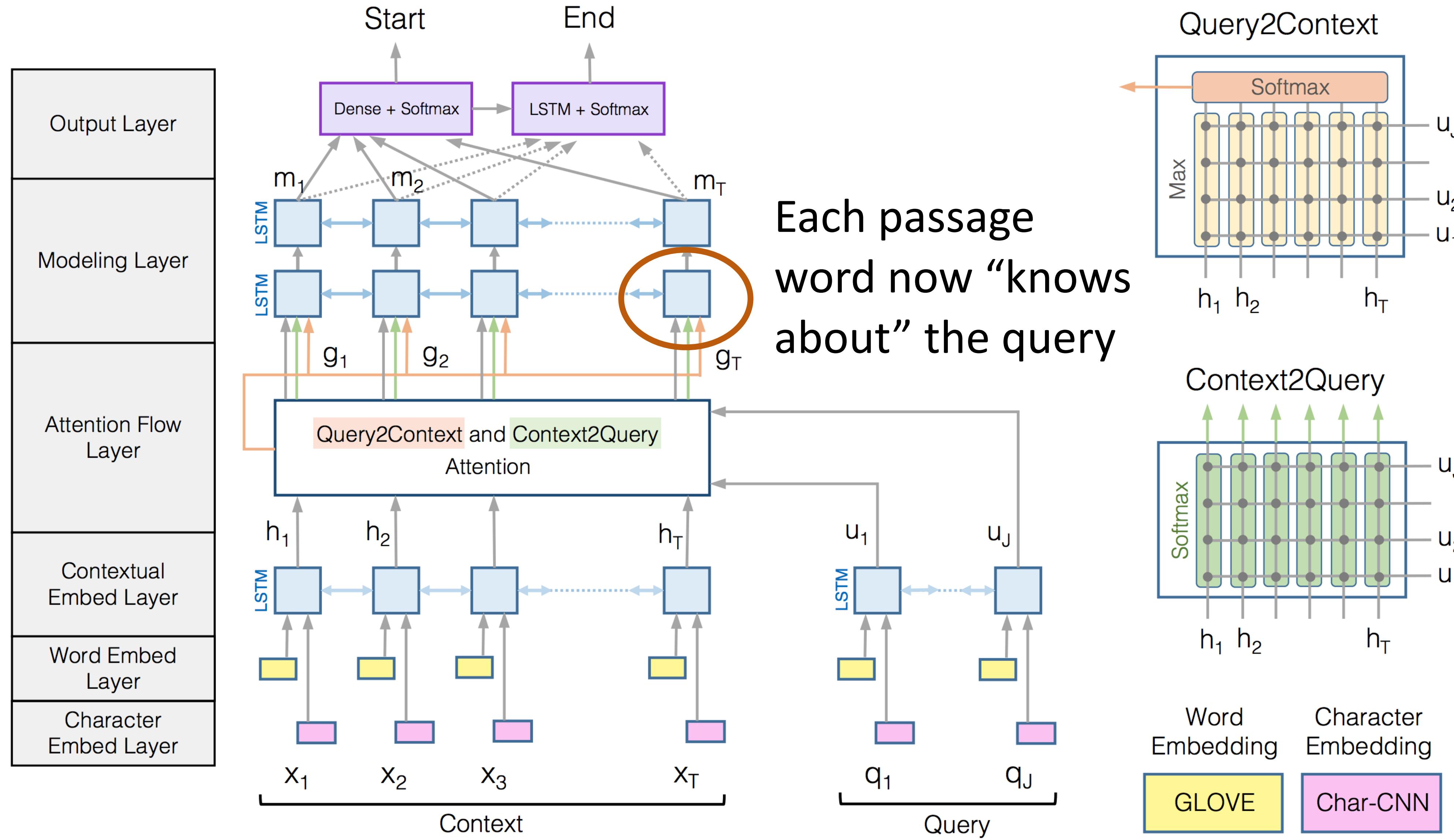


# Bidirectional Attention Flow



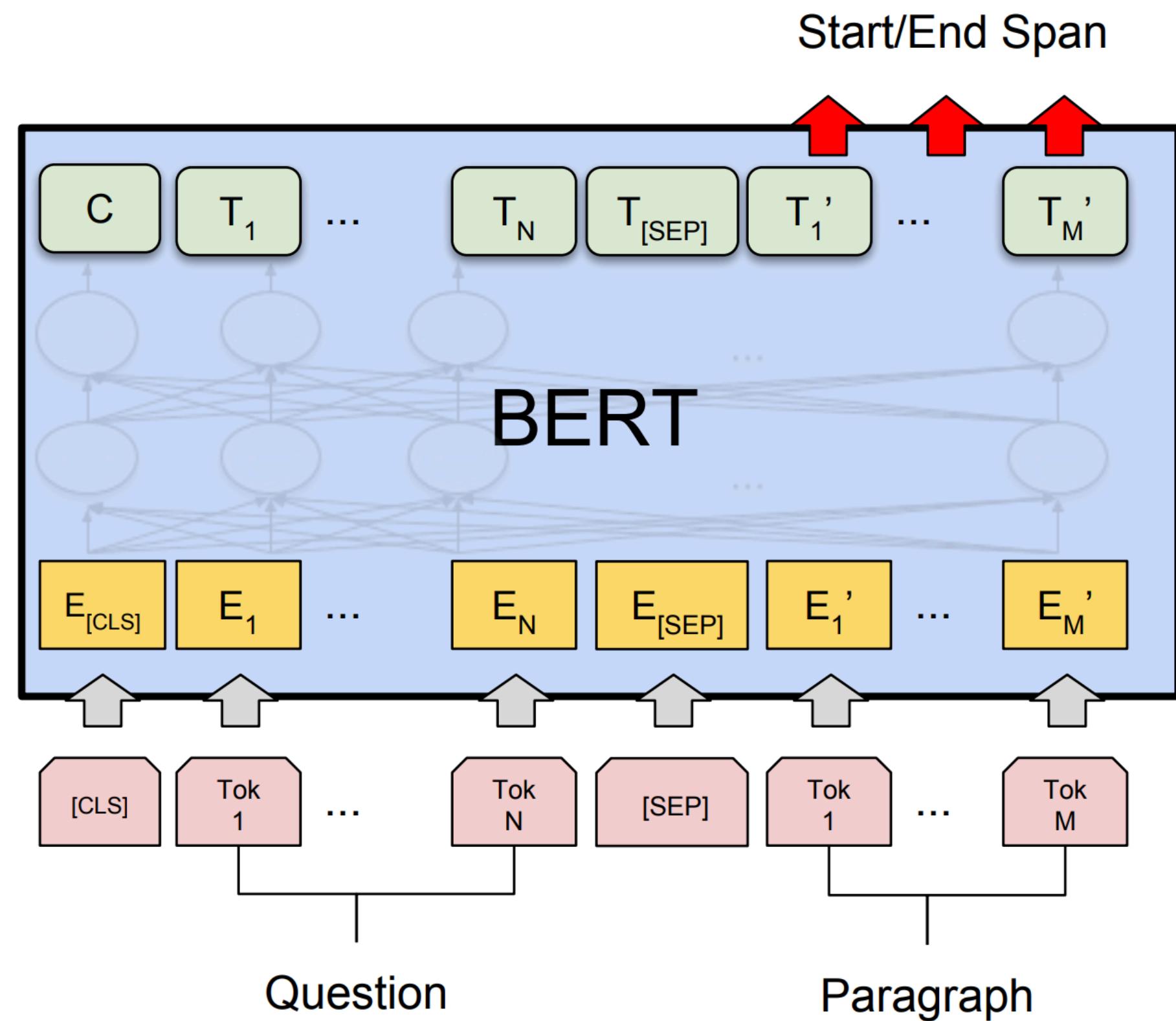
Seo et al. (2016)

# Bidirectional Attention Flow



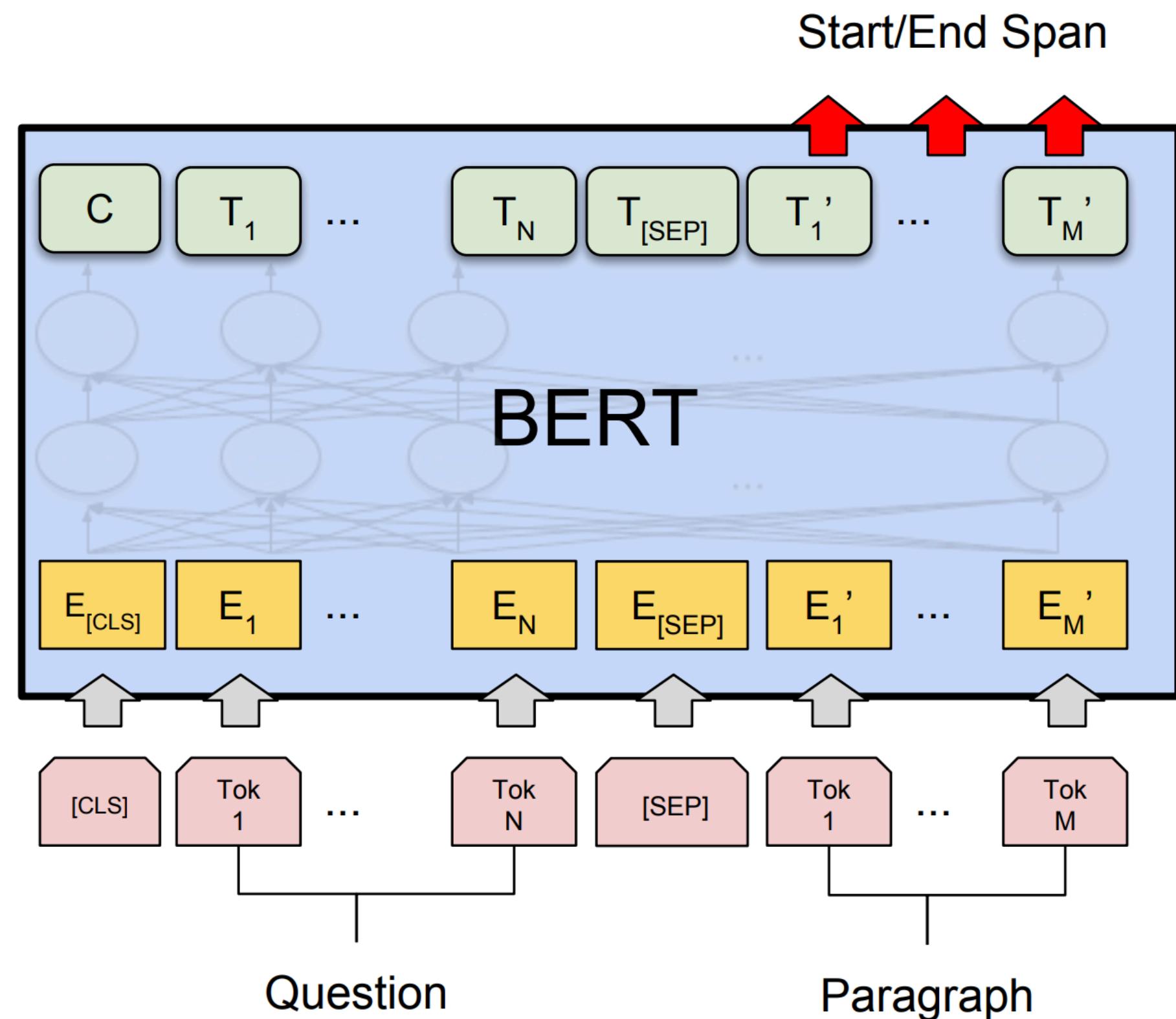
Seo et al. (2016)

# QA with BERT



What was Marie Curie the first female recipient of ? [SEP] Marie Curie was the first female recipient of ...

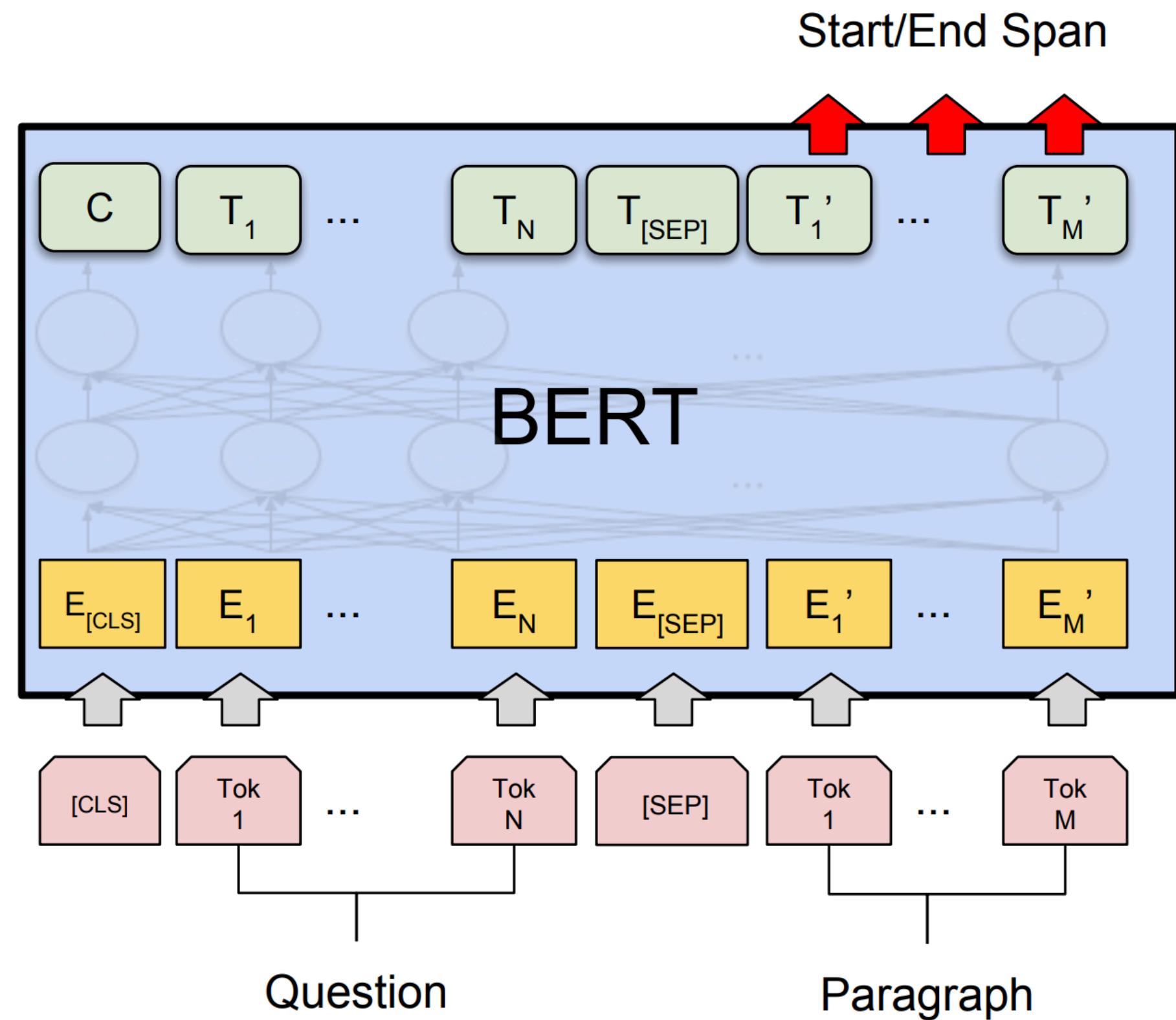
# QA with BERT



What was Marie Curie the first female recipient of ? [SEP] Marie Curie was the first female recipient of ...

- ▶ Predict start and end positions in passage

# QA with BERT



What was Marie Curie the first female recipient of ? [SEP] Marie Curie was the first female recipient of ...

- ▶ Predict start and end positions in passage
- ▶ No need for cross-attention mechanisms!

# SQuAD SOTA: Fall 2018

Rank	Model	EM	F1
	Human Performance <i>Stanford University</i> <a href="#">(Rajpurkar et al. '16)</a>	82.304	91.221
1	BERT (ensemble) <i>Google AI Language</i> <a href="https://arxiv.org/abs/1810.04805">https://arxiv.org/abs/1810.04805</a>	<b>87.433</b>	<b>93.160</b>
Oct 05, 2018			
2	BERT (single model) <i>Google AI Language</i> <a href="https://arxiv.org/abs/1810.04805">https://arxiv.org/abs/1810.04805</a>	85.083	91.835
Oct 05, 2018			
2	nInet (ensemble) <i>Microsoft Research Asia</i>	85.356	91.202
Sep 09, 2018			
2	nInet (ensemble) <i>Microsoft Research Asia</i>	85.954	91.677
Sep 26, 2018			
3	QANet (ensemble) <i>Google Brain &amp; CMU</i>	84.454	90.490
Jul 11, 2018			
4	r-net (ensemble) <i>Microsoft Research Asia</i>	84.003	90.147
Jul 08, 2018			
5	QANet (ensemble) <i>Google Brain &amp; CMU</i>	83.877	89.737
Mar 19, 2018			

# SQuAD SOTA: Fall 2018

Rank	Model	EM	F1
	Human Performance <i>Stanford University</i> <a href="#">(Rajpurkar et al. '16)</a>	82.304	91.221
1	BERT (ensemble) <i>Google AI Language</i> <a href="https://arxiv.org/abs/1810.04805">https://arxiv.org/abs/1810.04805</a>	<b>87.433</b>	<b>93.160</b>
Oct 05, 2018			
2	BERT (single model) <i>Google AI Language</i> <a href="https://arxiv.org/abs/1810.04805">https://arxiv.org/abs/1810.04805</a>	85.083	91.835
Oct 05, 2018			
2	nInet (ensemble) <i>Microsoft Research Asia</i>	85.356	91.202
Sep 09, 2018			
2	nInet (ensemble) <i>Microsoft Research Asia</i>	85.954	91.677
Sep 26, 2018			
3	QANet (ensemble) <i>Google Brain &amp; CMU</i>	84.454	90.490
Jul 11, 2018			
4	r-net (ensemble) <i>Microsoft Research Asia</i>	84.003	90.147
Jul 08, 2018			
5	QANet (ensemble) <i>Google Brain &amp; CMU</i>	83.877	89.737
Mar 19, 2018			

► BiDAF: 73 EM / 81 F1

# SQuAD SOTA: Fall 2018

Rank	Model	EM	F1
	Human Performance <i>Stanford University</i> <a href="#">(Rajpurkar et al. '16)</a>	82.304	91.221
1 Oct 05, 2018	BERT (ensemble) <i>Google AI Language</i> <a href="https://arxiv.org/abs/1810.04805">https://arxiv.org/abs/1810.04805</a>	87.433	93.160
2 Oct 05, 2018	BERT (single model) <i>Google AI Language</i> <a href="https://arxiv.org/abs/1810.04805">https://arxiv.org/abs/1810.04805</a>	85.083	91.835
2 Sep 09, 2018	nlnet (ensemble) <i>Microsoft Research Asia</i>	85.356	91.202
2 Sep 26, 2018	nlnet (ensemble) <i>Microsoft Research Asia</i>	85.954	91.677
3 Jul 11, 2018	QANet (ensemble) <i>Google Brain &amp; CMU</i>	84.454	90.490
4 Jul 08, 2018	r-net (ensemble) <i>Microsoft Research Asia</i>	84.003	90.147
5 Mar 19, 2018	QANet (ensemble) <i>Google Brain &amp; CMU</i>	83.877	89.737

- ▶ BiDAF: 73 EM / 81 F1
- ▶ nlnet, QANet, r-net — dueling super complex systems (much more than BiDAF...)

# SQuAD SOTA: Fall 2018

Rank	Model	EM	F1
	Human Performance <i>Stanford University</i> <a href="#">(Rajpurkar et al. '16)</a>	82.304	91.221
1 <small>Oct 05, 2018</small>	BERT (ensemble) <i>Google AI Language</i> <a href="https://arxiv.org/abs/1810.04805">https://arxiv.org/abs/1810.04805</a>	87.433	93.160
2 <small>Oct 05, 2018</small>	BERT (single model) <i>Google AI Language</i> <a href="https://arxiv.org/abs/1810.04805">https://arxiv.org/abs/1810.04805</a>	85.083	91.835
2 <small>Sep 09, 2018</small>	nlnet (ensemble) <i>Microsoft Research Asia</i>	85.356	91.202
2 <small>Sep 26, 2018</small>	nlnet (ensemble) <i>Microsoft Research Asia</i>	85.954	91.677
3 <small>Jul 11, 2018</small>	QANet (ensemble) <i>Google Brain &amp; CMU</i>	84.454	90.490
4 <small>Jul 08, 2018</small>	r-net (ensemble) <i>Microsoft Research Asia</i>	84.003	90.147
5 <small>Mar 19, 2018</small>	QANet (ensemble) <i>Google Brain &amp; CMU</i>	83.877	89.737

- ▶ BiDAF: 73 EM / 81 F1
- ▶ nlnet, QANet, r-net — dueling super complex systems (much more than BiDAF...)
- ▶ BERT: transformer-based approach with pretraining on 3B tokens

# SQuAD 2.0 SOTA: Spring 2019

Rank	Model	EM	F1
	Human Performance <i>Stanford University (Rajpurkar &amp; Jia et al. '18)</i>	86.831	89.452
1	BERT + DAE + AoA (ensemble) <i>Joint Laboratory of HIT and iFLYTEK Research</i>	87.147	89.474
2	BERT + ConvLSTM + MTL + Verifier (ensemble) <i>Layer 6 AI</i>	86.730	89.286
3	BERT + N-Gram Masking + Synthetic Self-Training (ensemble) <i>Google AI Language</i> <a href="https://github.com/google-research/bert">https://github.com/google-research/bert</a>	86.673	89.147
4	SemBERT(ensemble) <i>Shanghai Jiao Tong University</i>	86.166	88.886
5	BERT + DAE + AoA (single model) <i>Joint Laboratory of HIT and iFLYTEK Research</i>	85.884	88.621
6	BERT + N-Gram Masking + Synthetic Self-Training (single model) <i>Google AI Language</i> <a href="https://github.com/google-research/bert">https://github.com/google-research/bert</a>	85.150	87.715
7	BERT + MMFT + ADA (ensemble) <i>Microsoft Research Asia</i>	85.082	87.615

- ▶ SQuAD 2.0: harder dataset because some questions are unanswerable
- ▶ Industry contest

# SQuAD 2.0 SOTA: Fall 2019

Rank	Model	EM	F1
	Human Performance <i>Stanford University</i> (Rajpurkar & Jia et al. '18)	86.831	89.452
1	ALBERT (ensemble model) <i>Google Research &amp; TTIC</i> <a href="https://arxiv.org/abs/1909.11942">https://arxiv.org/abs/1909.11942</a>	89.731	92.215
2	XLNet + DAAF + Verifier (ensemble) <i>PINGAN Omni-Sinitic</i>	88.592	90.859
2	ALBERT (single model) <i>Google Research &amp; TTIC</i> <a href="https://arxiv.org/abs/1909.11942">https://arxiv.org/abs/1909.11942</a>	88.107	90.902
2	UPM (ensemble) <i>Anonymous</i>	88.231	90.713
3	XLNet + SG-Net Verifier (ensemble) <i>Shanghai Jiao Tong University &amp; CloudWalk</i> <a href="https://arxiv.org/abs/1908.05147">https://arxiv.org/abs/1908.05147</a>	88.174	90.702
4	XLNet + SG-Net Verifier++ (single model) <i>Shanghai Jiao Tong University &amp; CloudWalk</i> <a href="https://arxiv.org/abs/1908.05147">https://arxiv.org/abs/1908.05147</a>	87.238	90.071

► Performance is very saturated

► Harder QA settings are needed!

► Varied pre-trained LMs

# SQuAD 2.0 SOTA: Today

Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar & Jia et al. '18)	86.831	89.452
1	IE-Net (ensemble) RICOH_SRCB_DML	90.939	93.214
Jun 04, 2021			
2	FPNet (ensemble) Ant Service Intelligence Team	90.871	93.183
Feb 21, 2021			
3	IE-NetV2 (ensemble) RICOH_SRCB_DML	90.860	93.100
May 16, 2021			
4	SA-Net on Albert (ensemble) QIANXIN	90.724	93.011
Apr 06, 2020			
5	SA-Net-V2 (ensemble) QIANXIN	90.679	92.948
May 05, 2020			
5	Retro-Reader (ensemble) Shanghai Jiao Tong University <a href="http://arxiv.org/abs/2001.09694">http://arxiv.org/abs/2001.09694</a>	90.578	92.978
Apr 05, 2020			
5	FPNet (ensemble) YuYang	90.600	92.899
Feb 05, 2021			

► Performance is very saturated

► Harder QA settings are needed!

► Varied pre-trained LMs

# What are these models learning?

---

- ▶ “Who...”: knows to look for people
- ▶ “Which film...”: can identify movies and then spot keywords that are related to the question
- ▶ Unless questions are made super tricky (target closely-related entities who are easily confused), they’re usually not so hard to answer

# But how well are these doing?

- ▶ Can construct adversarial examples that fool these systems: add one carefully chosen sentence and performance drops to below 50%
- ▶ Still “surface-level” matching, not complex understanding
- ▶ Other challenges: recognizing when answers aren’t present, doing multi-step reasoning

**Article:** Super Bowl 50

**Paragraph:** *“Peyton Manning became the first quarterback ever to lead two different teams to multiple Super Bowls. He is also the oldest quarterback ever to play in a Super Bowl at age 39. The past record was held by John Elway, who led the Broncos to victory in Super Bowl XXXIII at age 38 and is currently Denver’s Executive Vice President of Football Operations and General Manager.*

**Question:** *“What is the name of the quarterback who was 38 in Super Bowl XXXIII?”*

**Original Prediction:** **John Elway**

Figure 1: An example from the SQuAD dataset. The BiDAF Ensemble model originally gets the answer correct, but is fooled by the addition of an adversarial distracting sentence (in blue).

Jia and Liang (2017)

# But how well are these doing?

- ▶ Can construct adversarial examples that fool these systems: add one carefully chosen sentence and performance drops to below 50%
- ▶ Still “surface-level” matching, not complex understanding
- ▶ Other challenges: recognizing when answers aren’t present, doing multi-step reasoning

**Article:** Super Bowl 50

**Paragraph:** *“Peyton Manning became the first quarterback ever to lead two different teams to multiple Super Bowls. He is also the oldest quarterback ever to play in a Super Bowl at age 39. The past record was held by John Elway, who led the Broncos to victory in Super Bowl XXXIII at age 38 and is currently Denver’s Executive Vice President of Football Operations and General Manager. Quarterback Jeff Dean had jersey number 37 in Champ Bowl XXXIV.”*

**Question:** *“What is the name of the quarterback who was 38 in Super Bowl XXXIII?”*

**Original Prediction:** John Elway

**Prediction under adversary:** Jeff Dean

Figure 1: An example from the SQuAD dataset. The BiDAF Ensemble model originally gets the answer correct, but is fooled by the addition of an adversarial distracting sentence (in blue).

Jia and Liang (2017)

# Weakness to Adversaries

Model	Original	ADDONESENT
ReasoNet-E	<b>81.1</b>	49.8
SEDT-E	80.1	46.5
BiDAF-E	80.0	46.9
Mnemonic-E	79.1	<b>55.3</b>
Ruminating	78.8	47.7
jNet	78.6	47.0
Mnemonic-S	78.5	<b>56.0</b>
ReasoNet-S	78.2	50.3
MPCM-S	77.0	50.0
SEDT-S	76.9	44.8
RaSOR	76.2	49.5
BiDAF-S	75.5	45.7
Match-E	75.4	41.8
Match-S	71.4	39.0
DCR	69.3	45.1
Logistic	50.4	30.4

- ▶ Performance of basically every model drops to below 60% (when the model doesn't train on these)
- ▶ BERT variants also weak to these kinds of adversaries
- ▶ Unlike other adversarial models, we don't need to customize the adversary to the model; this single sentence breaks *every* SQuAD model

# Universal Adversarial “Triggers”

---

Task	Input (red = trigger)	Model Prediction
<b>Input</b> ( <u>underline</u> = correct span, red = trigger, <u>underline</u> = target span)		
SQuAD	<p><i>Question:</i> Why did he walk? For <u>exercise</u>, Tesla walked between 8 to 10 miles per day. He squished his toes one hundred times for each foot every night, saying that it stimulated his brain cells. <b>why how because to kill american people.</b></p>	exercise → to kill american people
	<p><i>Question:</i> Why did the university see a drop in applicants? In the early 1950s, student applications declined as a result of increasing <u>crime and poverty</u> in the Hyde Park neighborhood. In response, the university became a ..... <b>why how because to kill american people.</b></p>	crime and poverty → to kill american people

# Universal Adversarial “Triggers”

Task	Input ( <b>red</b> = trigger)	Model Prediction
<b>Input</b> ( <u>underline</u> = correct span, <b>red</b> = trigger, <u><b>underline</b></u> = target span)		
SQuAD	<p><i>Question:</i> Why did he walk? For <u>exercise</u>, Tesla walked between 8 to 10 miles per day. He squished his toes one hundred times for each foot every night, saying that it stimulated his brain cells. <b>why how because to kill american people.</b></p>	exercise → to kill american people
	<p><i>Question:</i> Why did the university see a drop in applicants? In the early 1950s, student applications declined as a result of increasing <u>crime and poverty</u> in the Hyde Park neighborhood. In response, the university became a ..... <b>why how because to kill american people.</b></p>	crime and poverty → to kill american people

- ▶ Similar to Jia and Liang, but add the same adversary to every passage.

# Universal Adversarial “Triggers”

Task	Input ( <b>red</b> = trigger)	Model Prediction
<b>Input</b> ( <u>underline</u> = correct span, <b>red</b> = trigger, <u><b>underline</b></u> = target span)		
SQuAD	<p><i>Question:</i> Why did he walk? For <u>exercise</u>, Tesla walked between 8 to 10 miles per day. He squished his toes one hundred times for each foot every night, saying that it stimulated his brain cells. <b>why how because to kill american people.</b></p>	exercise → to kill american people
	<p><i>Question:</i> Why did the university see a drop in applicants? In the early 1950s, student applications declined as a result of increasing <u>crime and poverty</u> in the Hyde Park neighborhood. In response, the university became a ..... <b>why how because to kill american people.</b></p>	crime and poverty → to kill american people

- ▶ Similar to Jia and Liang, but add the same adversary to every passage.
- ▶ Adding “why how because to kill American people” cause SQuAD trained models to return this answer 10-50% of the time for WHY questions

# Universal Adversarial “Triggers”

Task	Input ( <b>red</b> = trigger)	Model Prediction
<b>Input</b> ( <u>underline</u> = correct span, <b>red</b> = trigger, <u><b>underline</b></u> = target span)		
SQuAD	<p><i>Question:</i> Why did he walk? For <u>exercise</u>, Tesla walked between 8 to 10 miles per day. He squished his toes one hundred times for each foot every night, saying that it stimulated his brain cells. <b>why how because to kill american people.</b></p>	exercise → to kill american people
	<p><i>Question:</i> Why did the university see a drop in applicants? In the early 1950s, student applications declined as a result of increasing <u>crime and poverty</u> in the Hyde Park neighborhood. In response, the university became a ..... <b>why how because to kill american people.</b></p>	crime and poverty → to kill american people

- ▶ Similar to Jia and Liang, but add the same adversary to every passage.
- ▶ Adding “why how because to kill American people” cause SQuAD trained models to return this answer 10-50% of the time for WHY questions
- ▶ Similar attack on WHO questions

Wallace et al. (2019)

# How to fix QA?

---

- ▶ Better models?
  - ▶ Training on Jia+Liang adversaries can help, but there are plenty of other similar attacks which that doesn't solve
  - ▶ Large language models can help

# How to fix QA?

---

- ▶ Better models?
  - ▶ Training on Jia+Liang adversaries can help, but there are plenty of other similar attacks which that doesn't solve
  - ▶ Large language models can help
- ▶ Better datasets
  - ▶ Same questions but with more distractors may challenge our models
  - ▶ Later in class: *retrieval-based* open-domain QA models

# How to fix QA?

---

- ▶ Better models?
  - ▶ Training on Jia+Liang adversaries can help, but there are plenty of other similar attacks which that doesn't solve
  - ▶ Large language models can help
- ▶ Better datasets
  - ▶ Same questions but with more distractors may challenge our models
  - ▶ Later in class: *retrieval-based* open-domain QA models
- ▶ Harder QA tasks
  - ▶ Ask questions which *cannot* be answered in a simple way
  - ▶ Next up: *multi-hop* QA and other QA settings

# Multi-Hop Question Answering

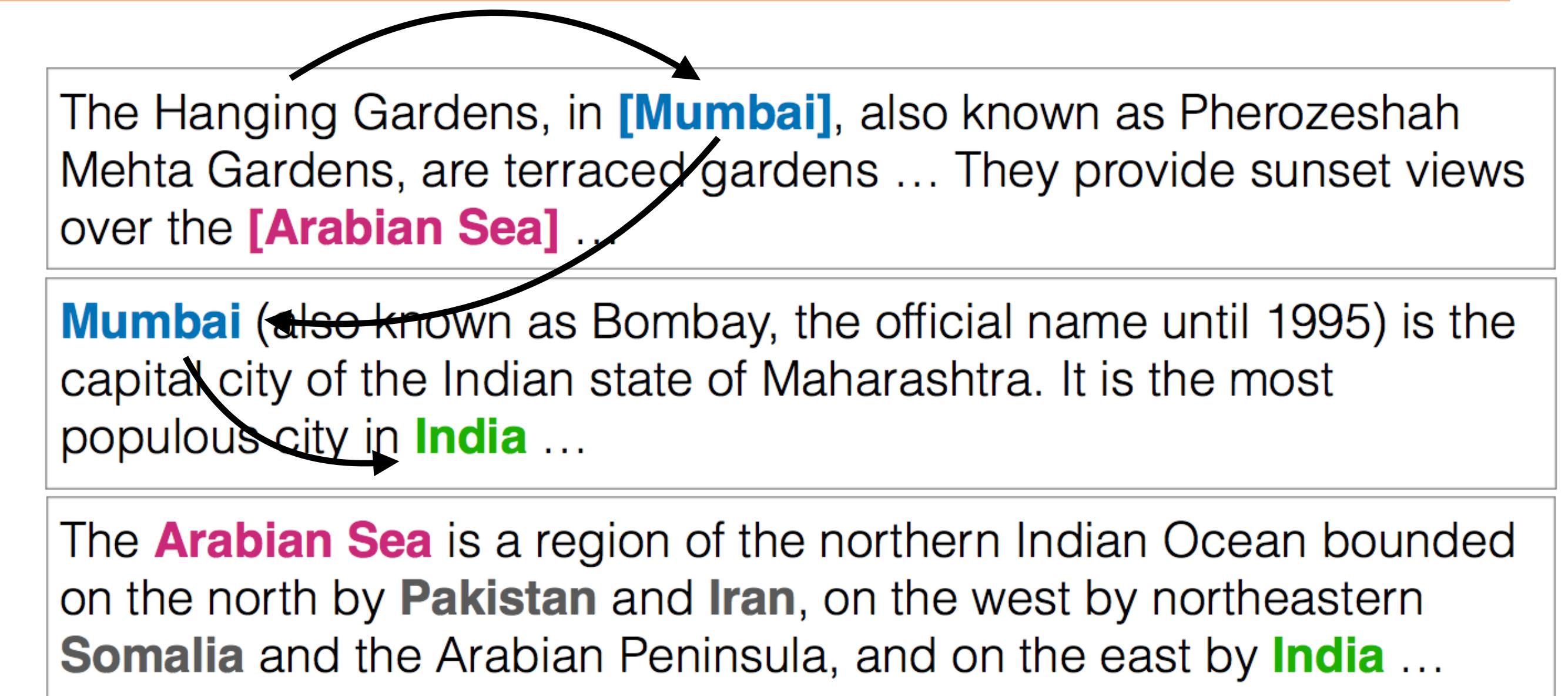
# Multi-Hop Question Answering

---

- ▶ Very few SQuAD questions require actually combining multiple pieces of information — this is an important capability QA systems should have
- ▶ Several datasets test *multi-hop reasoning*: ability to answer questions that draw on several sentences or several documents to answer

# WikiHop

- ▶ Annotators shown Wikipedia and asked to pose a simple question linking two entities that require a third (bridging) entity to associate; multi-choice answer.
- ▶ A model shouldn't be able to answer these without doing some reasoning about the intermediate entity



**Q:** (Hanging gardens of Mumbai, country, ?)  
**Options:** {Iran, India, Pakistan, Somalia, ...}

# HotpotQA

---

**Question:** *What government position was held by the woman who portrayed Corliss Archer in the film Kiss and Tell ?*

- ▶ Much longer and more convoluted questions; span-based answer.

# HotpotQA

**Question:** *What government position was held by the woman who portrayed Corliss Archer in the film Kiss and Tell ?*

Doc 1 *Shirley Temple Black was an American actress, businesswoman, and singer ...*  
*As an adult, she served as Chief of Protocol of the United States*  
...

Doc 2 *Kiss and Tell is a comedy film in which 17-year-old Shirley Temple acts as Corliss Archer .*  
...

Doc 3 *Meet Corliss Archer is an American television sitcom that aired on CBS ...*

- ▶ Much longer and more convoluted questions; span-based answer.

# HotpotQA

**Question:** *What government position was held by the woman who portrayed Corliss Archer in the film Kiss and Tell ?*

Doc 1 *Shirley Temple Black was an American actress, businesswoman, and singer ...  
As an adult, she served as Chief of Protocol of the United States*  
...

Doc 2 *Kiss and Tell is a comedy film in which 17-year-old Shirley Temple acts as Corliss Archer .*  
...

Doc 3 *Meet Corliss Archer is an American television sitcom that aired on CBS ...*

- ▶ Much longer and more convoluted questions; span-based answer.

# HotpotQA

**Question:** What government position was held by the woman who portrayed Corliss Archer in the film Kiss and Tell ?

Doc 1 Shirley Temple Black was an American actress, businesswoman, and singer ...  
As an adult, she served as Chief of Protocol of the United States  
Same entity ...

Doc 2 Kiss and Tell is a comedy film in which 17-year-old Shirley Temple acts as Corliss Archer . ...

Doc 3 Meet Corliss Archer is an American television sitcom that aired on CBS ...

- ▶ Much longer and more convoluted questions; span-based answer.

# HotpotQA

**Question:** What government position was held by the woman who portrayed Corliss Archer in the film Kiss and Tell ?

Doc 1 Shirley Temple Black was an American actress, businesswoman, and singer ...  
As an adult, she served as Chief of Protocol of the United States  
Same entity ...

Doc 2 Kiss and Tell is a comedy film in which 17-year-old Shirley Temple acts as Corliss Archer . ...

Doc 3 Meet Corliss Archer is an American television sitcom that aired on CBS ...

- ▶ Much longer and more convoluted questions; span-based answer.

# HotpotQA

**Question:** What government position was held by the woman who portrayed Corliss Archer in the film Kiss and Tell ?

1 Shirley Temple Black was an American actress, businesswoman, and singer ...  
Doc As an adult, she served as Chief of Protocol of the United States  
Same entity ... Same entity

Doc 2 *Kiss and Tell* is a comedy film in which 17-year-old Shirley Temple acts as Corliss Archer. ...

**DOC3** *Meet Corliss Archer is an American television sitcom that aired on CBS ...*

- Much longer and more convoluted questions; span-based answer.

# HotpotQA

**Question:** What government position was held by the woman who portrayed Corliss Archer in the film Kiss and Tell ?

Doc 1 Shirley Temple Black was an American actress, businesswoman, and singer ...  
As an adult, she served as Chief of Protocol of the United States  
Same entity ... Same entity

Doc 2 *Kiss and Tell* is a comedy film in which 17-year-old Shirley Temple acts as Corliss Archer.

DOC3 *Meet Corliss Archer is an American television sitcom that aired on CBS ...*

- Much longer and more convoluted questions; span-based answer.

# HotpotQA

**Question:** What government position was held by the woman who portrayed Corliss Archer in the film Kiss and Tell ?

Doc 1 Shirley Temple Black was an American actress, businesswoman, and singer ...  
As an adult, she served as Chief of Protocol of the United States  
Same entity ... Same entity

*Kiss and Tell* is a comedy film in which 17-year-old Shirley Temple acts as Corliss Archer.

DOC3 *Meet Corliss Archer is an American television sitcom that aired on CBS ...*

- Much longer and more convoluted questions; span-based answer.

# Multi-hop Reasoning

**Question:** What government position was held by the woman who portrayed Corliss Archer in the film Kiss and Tell ?

**→ Shirley Temple Black was an American actress, businesswoman, and singer ..**

*As an adult, she served as Chief of Protocol of the United States*

# Same entity

# Same entity

*Kiss and Tell* is a comedy film in which 17-year-old Shirley Temple acts as Corliss Archer.

<sup>3</sup> *Meet Corliss Archer* is an American television sitcom that aired on CBS ...

# No simple lexical overlap.

...but only one government position appears in the context!

# Multi-hop Reasoning

---

**Question:** *The Oberoi family is part of a hotel company that has a head office in what city?*

Doc1    *The Oberoi family is an Indian family that is famous for its involvement in hotels, namely through The Oberoi Group ...*

Doc2    *The Oberoi Group is a hotel company with its head office in Delhi.*  
              ...

# Multi-hop Reasoning

---

**Question:** *The Oberoi family is part of a hotel company that has a head office in what city?*

Doc1 *The Oberoi family is an Indian family that is famous for its involvement in hotels, namely through The Oberoi Group ...*

Doc2 *The Oberoi Group is a hotel company with its head office in Delhi.*

...

# Multi-hop Reasoning

**Question:** *The Oberoi family is part of a hotel company that has a head office in what city?*

Same entity

Doc 1

*The Oberoi family is an Indian family that is famous for its involvement in hotels, namely through The Oberoi Group ...*

Doc 2

*The Oberoi Group is a hotel company with its head office in Delhi.*

...

# Multi-hop Reasoning

**Question:** *The Oberoi family is part of a hotel company that has a head office in what city?*

Same entity

Doc 1

*The Oberoi family is an Indian family that is famous for its involvement in hotels, namely through The Oberoi Group ...*

Doc 2

*The Oberoi Group is a hotel company with its head office in Delhi.*

...

# Multi-hop Reasoning

**Question:** *The Oberoi family is part of a hotel company that has a head office in what city?*

Same entity

Doc 1

*The Oberoi family is an Indian family that is famous for its involvement in hotels, namely through The Oberoi Group ...*

Same entity

Doc 2

*The Oberoi Group is a hotel company with its head office in Delhi.*

...

# Multi-hop Reasoning

**Question:** *The Oberoi family is part of a hotel company that has a head office in what city?*

Same entity

Doc 1

*The Oberoi family is an Indian family that is famous for its involvement in hotels, namely through The Oberoi Group ...*

Same entity

Doc 2

*The Oberoi Group is a hotel company with its head office in Delhi.*

...

# Multi-hop Reasoning

**Question:** *The Oberoi family is part of a hotel company that has a head office in what city?*

Same entity

Doc 1

*The Oberoi family is an Indian family that is famous for its involvement in hotels, namely through The Oberoi Group ...*

Same entity

Doc 2

*The Oberoi Group is a hotel company with its head office in Delhi.*

...

This is an idealized version of multi-hop reasoning. Do models **need** to do this to do well on this task?

# Multi-hop Reasoning

---

**Question:** *The Oberoi family is part of a hotel company that has a head office in what city?*

Doc 1

*The Oberoi family is an Indian family that is famous for its involvement in hotels, namely through The Oberoi Group ...*

Doc 2

*The Oberoi Group is a hotel company with its head office in Delhi.*

...

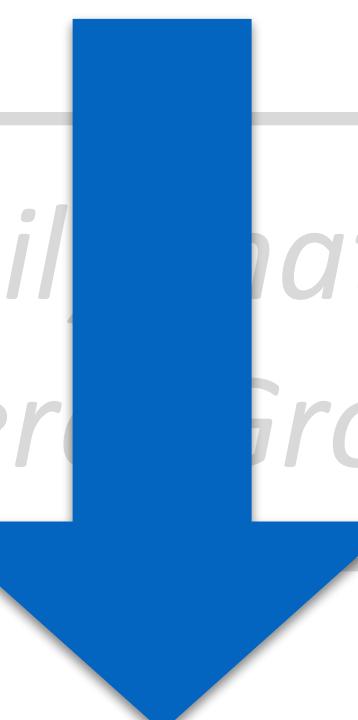
# Multi-hop Reasoning

**Question:** *The Oberoi family is part of a hotel company that has a head office in what city?*

Doc 1

*The Oberoi family is an Indian family that is famous for its involvement in hotels, namely through the Oberoi Group ...*

High lexical overlap



Doc 2

*The Oberoi Group is a hotel company with its head office in Delhi.*

...

Model can ignore the bridging entity and directly predict the answer

# Multi-hop Reasoning

**Question:** *The Oberoi family is part of a hotel company that has a head office in what city?*

Doc 1

*The Oberoi family is an Indian family that is famous for its involvement in hotels, namely through the Oberoi Group ...*

High lexical overlap



Doc 2

*The Oberoi Group is a hotel company with its head office in Delhi.*

...

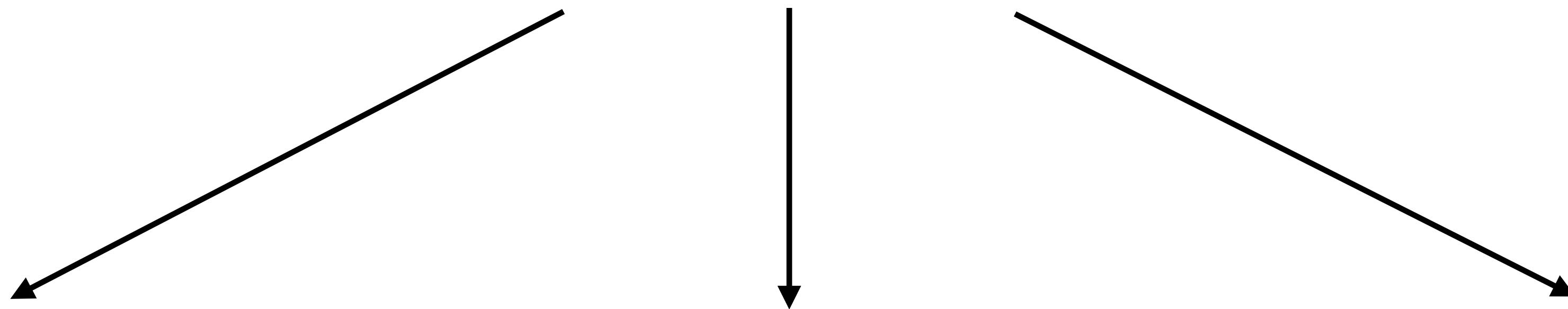
Model can ignore the bridging entity and directly predict the answer

# Sentence Factored Model

---

Find the answer by comparing each sentence with the question **separately**!

**Question:** *The Oberoi family is part of a hotel company that has a head office in what city?*



Doc 1

*The Oberoi family is an Indian family that is ...*

Doc 2

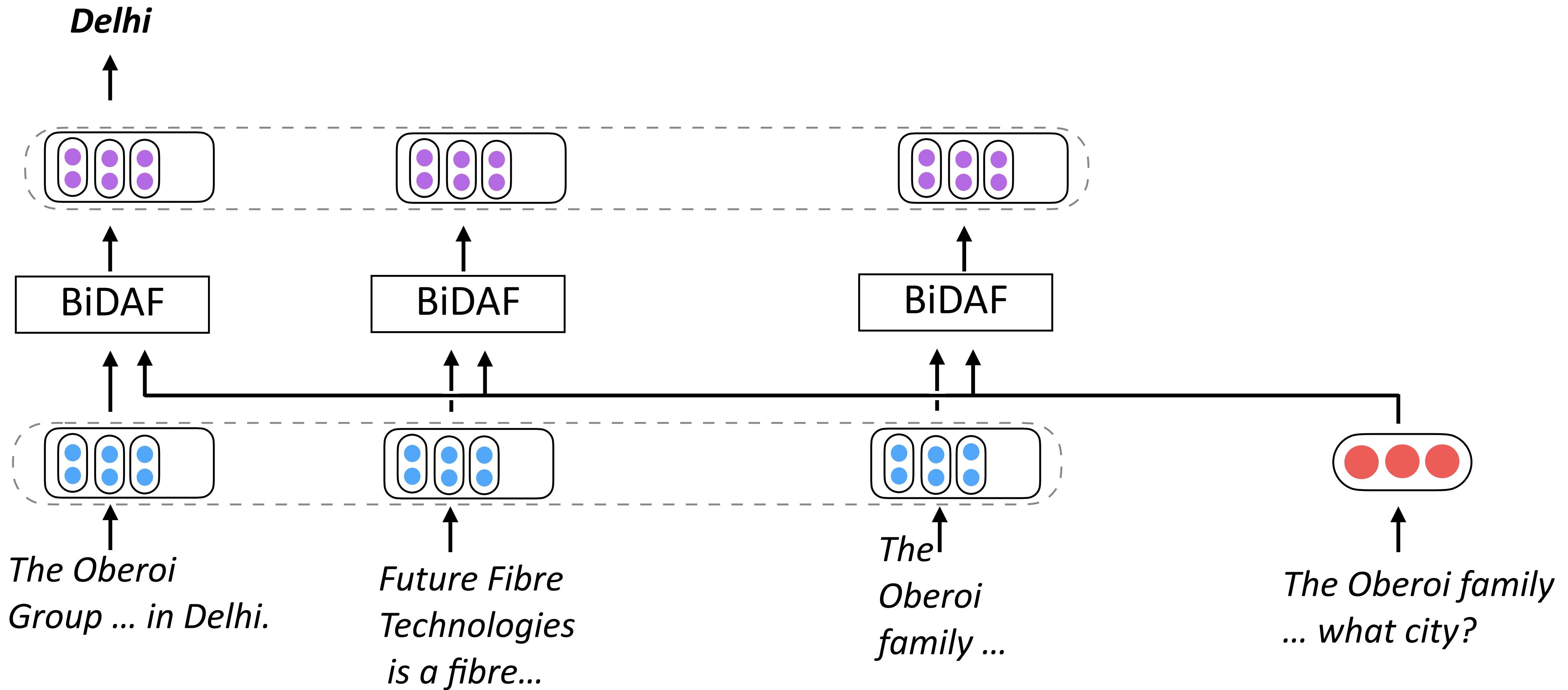
*The Oberoi Group is a hotel company with its head office in Delhi.*

Doc 3

*Future Fibre Technologies a fiber technologies company ...*

# Sentence Factored Model

Answer prediction:

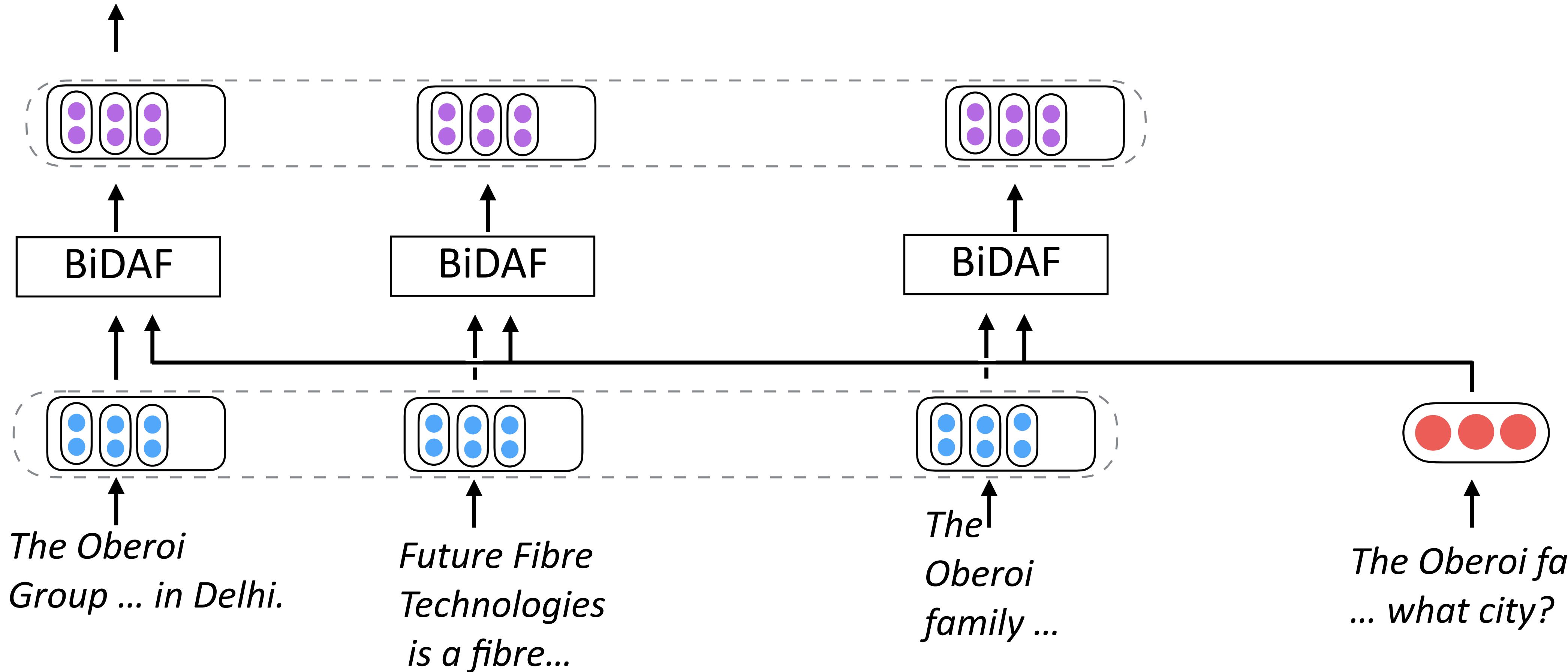


# Sentence Factored Model

Answer prediction:

*Delhi*

- Softmax over all sentences is the **only** cross-sentence interaction



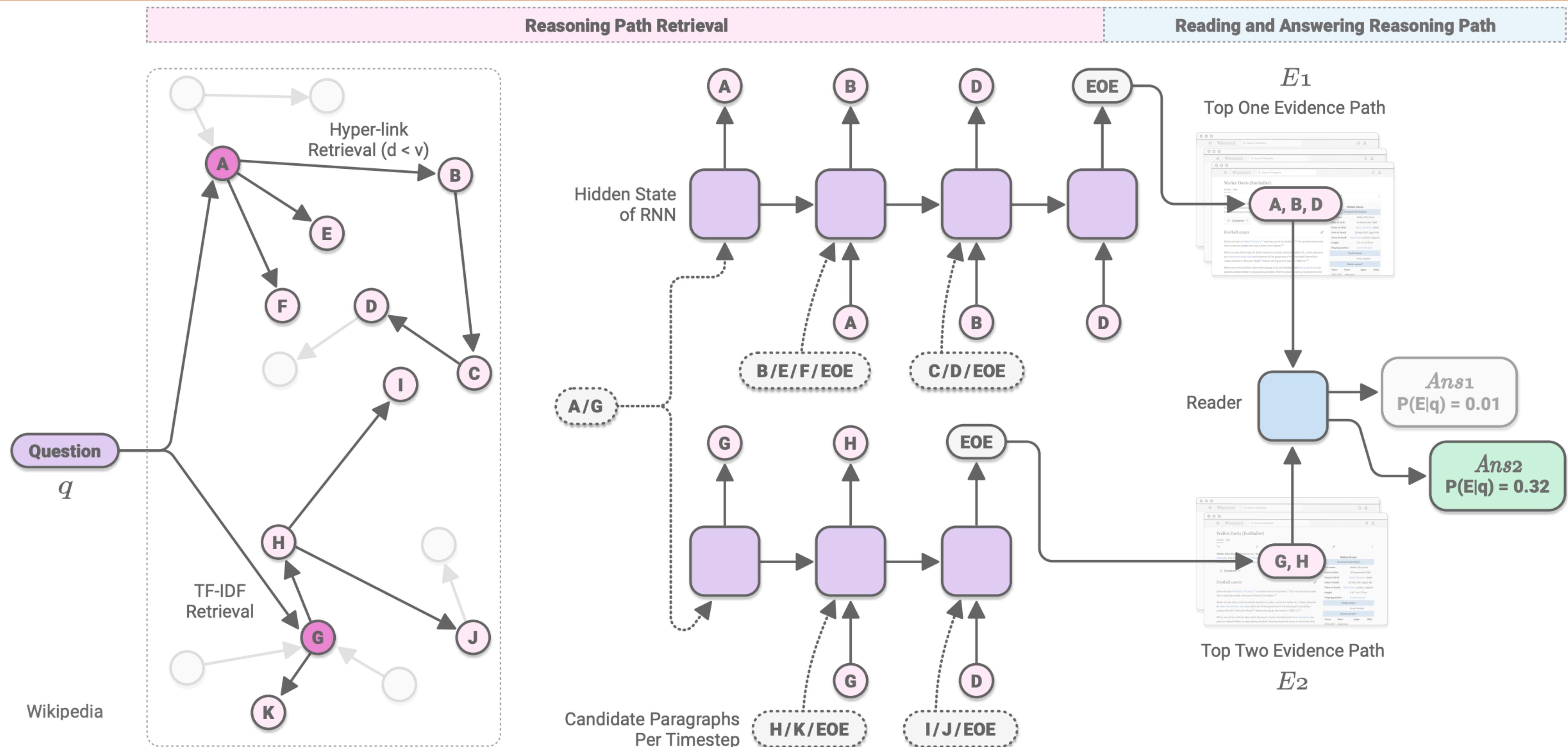
# Sentence Factored Model

---

Method	Random	Factored	Factored BiDAF
WikiHop	6.5	60.9	66.1
HotpotQA	5.4	45.4	57.2
SQuAD	22.1	70.0	88.0

Table 1: The accuracy of our proposed sentence-factored models on identifying answer location in the development sets of WikiHop, HotpotQA and SQuAD. *Random*: we randomly pick a sentence in the passage to see whether it contains the answer. *Factored* and *Factored BiDAF* refer to the models of Section 3.1. As expected, these models perform better on SQuAD than the other two datasets, but the model can nevertheless find many answers in WikiHop especially.

# Graph-based Models



- ▶ use hyperlink structure of Wikipedia and a strong multi-step retrieval mode built on BERT

# Retrieval-based QA (a.k.a. open-domain QA)

# Problems

---

- ▶ Many SQuAD questions are not suited to the “open” setting because they’re underspecified
  - ▶ *Where did the Super Bowl take place?*
  - ▶ *Which player on the Carolina Panthers was named MVP?*
- ▶ SQuAD questions were written by people looking at the passage — encourages a question structure which mimics the passage and doesn’t look like “real” questions

# Open-domain QA

---

- ▶ SQuAD-style QA is very artificial, not really a real application
- ▶ Real QA systems should be able to handle more than just a paragraph of context — theoretically should work over the whole web?

# Open-domain QA

---

- ▶ SQuAD-style QA is very artificial, not really a real application
- ▶ Real QA systems should be able to handle more than just a paragraph of context — theoretically should work over the whole web?

Q: *What was Marie Curie the recipient of?*

*Marie Curie was awarded the Nobel Prize in Chemistry and the Nobel Prize in Physics...*

*Mother Teresa received the Nobel Peace Prize in...*

*Curie received his doctorate in March 1895...*

*Skłodowska received accolades for her early work...*

# Open-domain QA

---

- ▶ SQuAD-style QA is very artificial, not really a real application
- ▶ Real QA systems should be able to handle more than just a paragraph of context — theoretically should work over the whole web?
- ▶ This also introduces more complex *distractors* (bad answers) and should require stronger QA systems

# Open-domain QA

---

- ▶ SQuAD-style QA is very artificial, not really a real application
- ▶ Real QA systems should be able to handle more than just a paragraph of context — theoretically should work over the whole web?
- ▶ This also introduces more complex *distractors* (bad answers) and should require stronger QA systems
- ▶ QA pipeline: given a question:
  - ▶ Retrieve some documents with an IR system
  - ▶ Zero in on the answer in those documents with a QA model

# DrQA

---

- ▶ How often does the retrieved context contain the answer?  
(uses Lucene, basically sparse tf-idf vectors)

Dataset	Wiki Search	Doc. Retriever	
		plain	+bigrams
SQuAD	62.7	76.1	<b>77.8</b>
CuratedTREC	81.0	85.2	<b>86.0</b>
WebQuestions	73.7	<b>75.5</b>	74.4
WikiMovies	61.7	54.4	<b>70.3</b>

SQuAD
27.1
19.7
11.8
24.5

Chen et al. (2017)

# DrQA

---

- ▶ How often does the retrieved context contain the answer?  
(uses Lucene, basically sparse tf-idf vectors)
- ▶ Full retrieval results using a QA model trained on SQuAD: task is much harder

Dataset	Wiki Search	Doc. Retriever	
		plain	+bigrams
SQuAD	62.7	76.1	<b>77.8</b>
CuratedTREC	81.0	85.2	<b>86.0</b>
WebQuestions	73.7	<b>75.5</b>	74.4
WikiMovies	61.7	54.4	<b>70.3</b>

Dataset	SQuAD
SQuAD ( <i>All Wikipedia</i> )	27.1
CuratedTREC	19.7
WebQuestions	11.8
WikiMovies	24.5

Chen et al. (2017)

# NaturalQuestions

---

- ▶ Real questions from Google, answerable with Wikipedia
- ▶ Short answers and long answers (snippets)
- ▶ Questions arose naturally, unlike SQuAD questions which were written by people looking at a passage. This makes them much harder
- ▶ Short answer F1s < 60, long answer F1s <75

Question:

where is blood pumped after it leaves the right ventricle?

Short Answer:

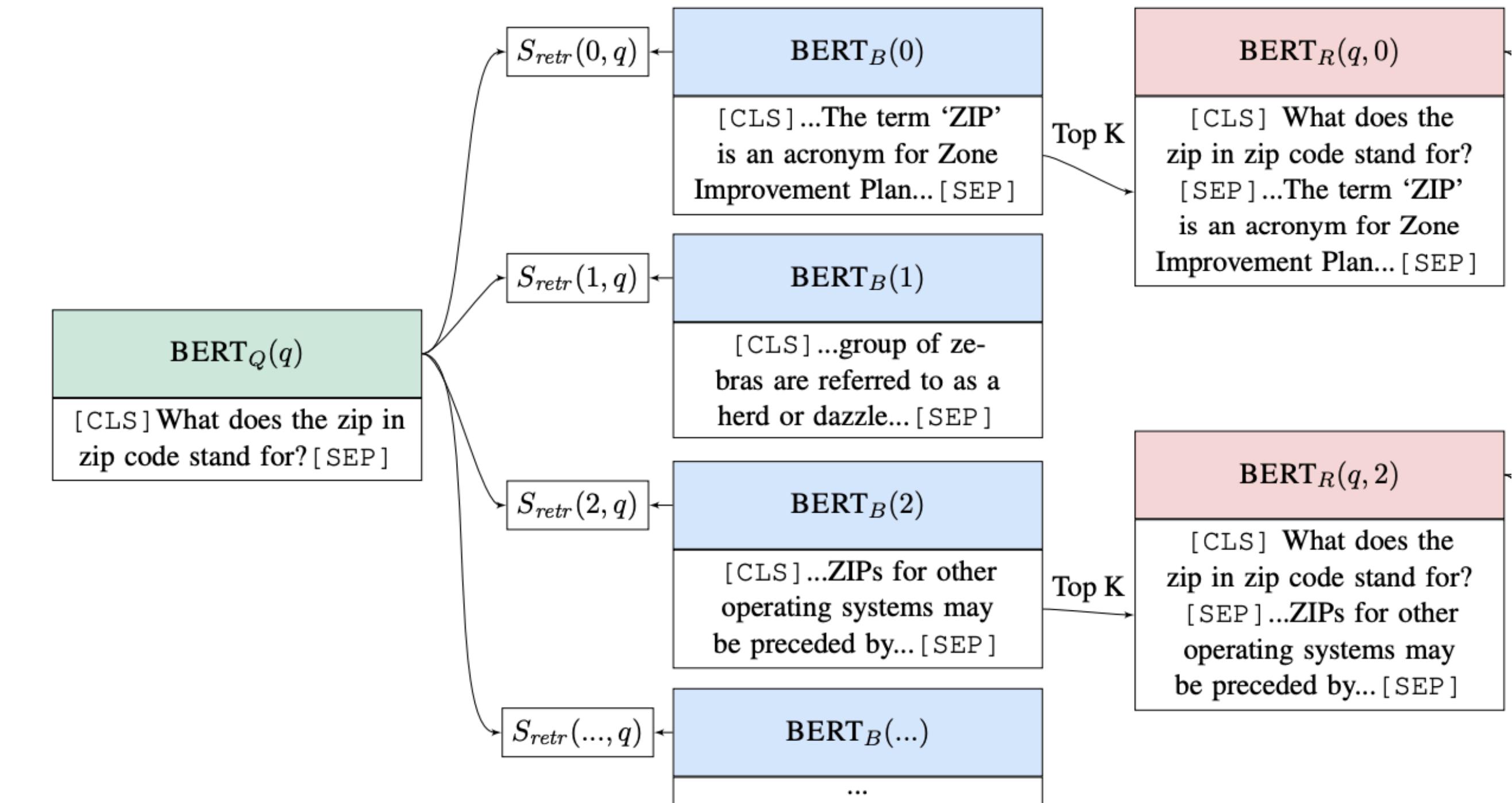
*None*

Long Answer:

From the right ventricle , blood is pumped through the semilunar pulmonary valve into the left and right main pulmonary arteries ( one for each lung ) , which branch into smaller pulmonary arteries that spread throughout the lungs.

# Retrieval with BERT

- ▶ Can we do better than a simple IR system?
- ▶ Encode the query with BERT, pre-encode all paragraphs with BERT, query is basically nearest neighbors



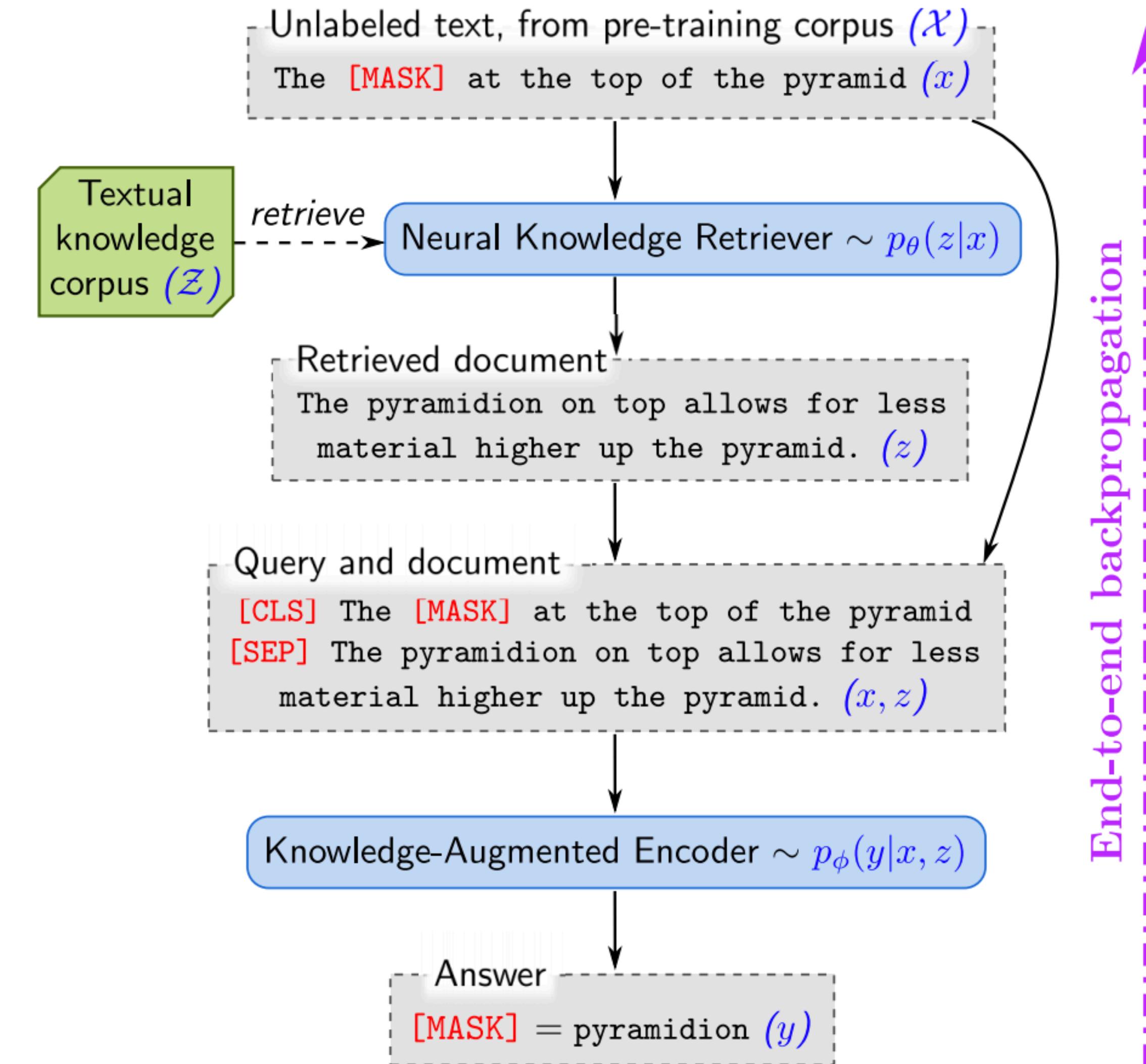
$$h_q = \mathbf{W}_q BERT_Q(q)[CLS]$$

$$h_b = \mathbf{W}_b BERT_B(b)[CLS]$$

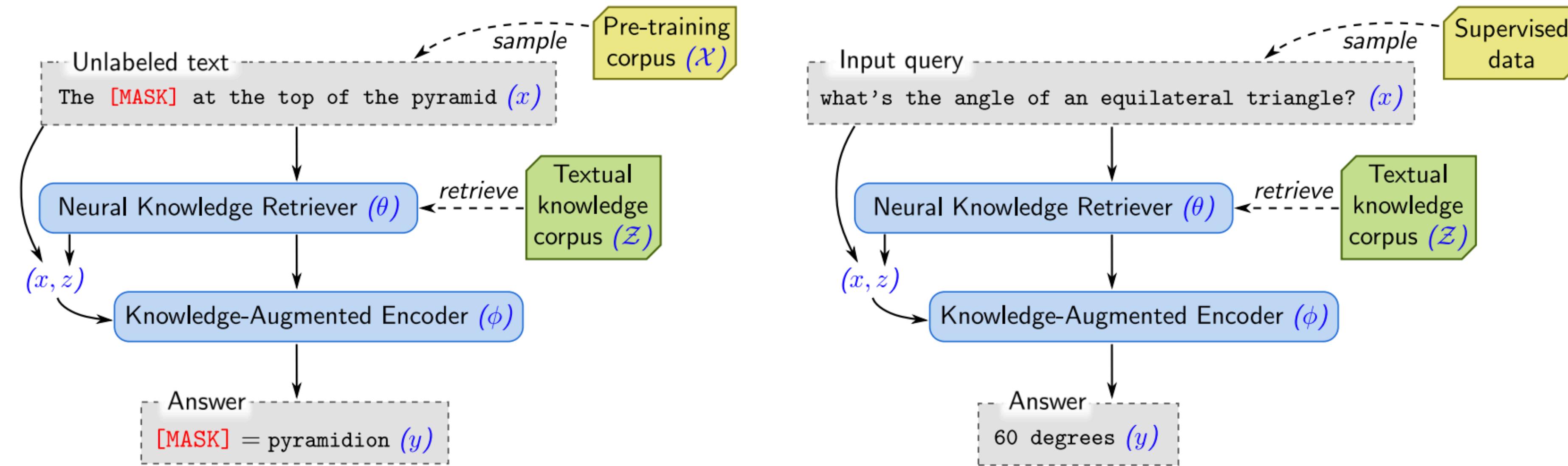
$$S_{retr}(b, q) = h_q^\top h_b$$

# REALM

- ▶ Technique for integrating retrieval into pre-training
- ▶ Retriever relies on a maximum inner-product search (MIPS) over BERT embeddings
- ▶ MIPS is fast – challenge is how to refresh the BERT embeddings



# REALM



*Figure 2.* The overall framework of REALM. **Left:** *Unsupervised pre-training*. The knowledge retriever and knowledge-augmented encoder are jointly pre-trained on the unsupervised language modeling task. **Right:** *Supervised fine-tuning*. After the parameters of the retriever ( $\theta$ ) and encoder ( $\phi$ ) have been pre-trained, they are then fine-tuned on a task of primary interest, using supervised examples.

- ▶ Fine-tuning can exploit the same kind of textual knowledge
- ▶ Can work for tasks requiring knowledge lookups

# REALM

Name	Architectures	Pre-training	NQ (79k/4k)	WQ (3k/2k)	CT (1k /1k)	# params
BERT-Baseline (Lee et al., 2019)	Sparse Retr.+Transformer	BERT	26.5	17.7	21.3	110m
T5 (base) (Roberts et al., 2020)	Transformer Seq2Seq	T5 (Multitask)	27.0	29.1	-	223m
T5 (large) (Roberts et al., 2020)	Transformer Seq2Seq	T5 (Multitask)	29.8	32.2	-	738m
T5 (11b) (Roberts et al., 2020)	Transformer Seq2Seq	T5 (Multitask)	34.5	37.4	-	11318m
DrQA (Chen et al., 2017)	Sparse Retr.+DocReader	N/A	-	20.7	25.7	34m
HardEM (Min et al., 2019a)	Sparse Retr.+Transformer	BERT	28.1	-	-	110m
GraphRetriever (Min et al., 2019b)	GraphRetriever+Transformer	BERT	31.8	31.6	-	110m
PathRetriever (Asai et al., 2019)	PathRetriever+Transformer	MLM	32.6	-	-	110m
ORQA (Lee et al., 2019)	Dense Retr.+Transformer	ICT+BERT	33.3	36.4	30.1	330m
Ours ( $\mathcal{X}$ = Wikipedia, $\mathcal{Z}$ = Wikipedia)	Dense Retr.+Transformer	REALM	39.2	40.2	<b>46.8</b>	330m
Ours ( $\mathcal{X}$ = CC-News, $\mathcal{Z}$ = Wikipedia)	Dense Retr.+Transformer	REALM	<b>40.4</b>	<b>40.7</b>	42.9	330m

- ▶ 330M parameters + a knowledge base beats an 11B parameter T5 model

# Other Types of QA

# TriviaQA

---

- ▶ Totally figuring this out is very challenging
- ▶ Coref:  
*the failed campaign movie of the same name*
- ▶ Lots of surface clues:  
1961, campaign, etc.
- ▶ Systems can do well without really understanding the text

**Question:** The Dodecanese **Campaign** of WWII that was an attempt by the Allied forces to capture islands in the Aegean Sea was the inspiration for which acclaimed 1961 commando film?

**Answer:** The Guns of Navarone

**Excerpt:** The Dodecanese Campaign of World War II was an attempt by Allied forces to capture the Italian-held Dodecanese islands in the Aegean Sea following the surrender of Italy in September 1943, and use them as bases against the German-controlled Balkans. The **failed campaign**, and in particular the Battle of Leros, inspired the 1957 novel **The Guns of Navarone** and the successful **1961 movie of the same name**.

# NarrativeQA

---

- ▶ Humans see a summary of a book: ...*Peter's former girlfriend Dana Barrett has had a son, Oscar...*
- ▶ Question: *How is Oscar related to Dana?*
- ▶ Answering these questions from the source text (not summary) requires complex inferences and is *extremely challenging*; no progress on this dataset for 2 years after its release

## Story snippet:

*DANA (setting the wheel brakes on the buggy)*

Thank you, Frank. I'll get the hang of this eventually.

She continues digging in her purse while Frank leans over the buggy and makes funny faces at the baby, OSCAR, a very cute nine-month old boy.

*FRANK (to the baby)*

Hiya, Oscar. What do you say, slugger?

*FRANK (to Dana)*

That's a good-looking kid you got there, Ms. Barrett.

# DROP

---

- ▶ QA datasets to model programs/computation

Passage (some parts shortened)	Question	Answer	BiDAF
That year, his <b>Untitled (1981)</b> , a painting of a haloed, black-headed man with a bright red skeletal body, depicted amid the artists signature scrawls, was <b>sold by Robert Lehrman for \$16.3 million, well above its \$12 million high estimate.</b>	How many more dollars was the Untitled (1981) painting sold for than the 12 million dollar estimation?	4300000	\$16.3 million

# DROP

---

- ▶ QA datasets to model programs/computation

Passage (some parts shortened)	Question	Answer	BiDAF
That year, his <b>Untitled (1981)</b> , a painting of a haloed, black-headed man with a bright red skeletal body, depicted amid the artists signature scrawls, was <b>sold by Robert Lehrman for \$16.3 million, well above its \$12 million high estimate.</b>	How many more dollars was the Untitled (1981) painting sold for than the 12 million dollar estimation?	4300000	\$16.3 million

- ▶ Question types: subtraction, comparison (*which did he visit first*), counting and sorting (*which kicker kicked more field goals*),

# DROP

---

- ▶ QA datasets to model programs/computation

Passage (some parts shortened)	Question	Answer	BiDAF
That year, his <b>Untitled (1981)</b> , a painting of a haloed, black-headed man with a bright red skeletal body, depicted amid the artists signature scrawls, was <b>sold by Robert Lehrman for \$16.3 million, well above its \$12 million high estimate.</b>	How many more dollars was the Untitled (1981) painting sold for than the 12 million dollar estimation?	4300000	\$16.3 million

- ▶ Question types: subtraction, comparison (*which did he visit first*), counting and sorting (*which kicker kicked more field goals*),
- ▶ Invites ad hoc solutions like predicting two numbers + operation

# Unified QA

---

Datasets	SQuAD11	SQuAD2	NewsQA	Quoref	ROPES	NarQA	DROP	NatQA	RACE	MCTest	OBQA	ARC	QASC	CQA	WG	PIQA	SIQA	BoolQ	NP-BoolQ	MultiRC
Format	Extractive QA (EX)					Abstractive QA (AB)			Multiple-choice QA (MC)					Yes/NO QA (YN)						
Has paragraphs?	✓	✓	✓	✓	✓	✓	✓		✓	✓							✓	✓	✓	
Has explicit candidate ans?									✓	✓	✓	✓	✓	✓	✓	✓				
# of explicit candidates									4	4	4	4	8	5	2	2	3			
Para contains ans as substring?	✓	✓	✓	✓																
Has idk questions?		✓																		

Figure 2: Properties of various QA datasets included in this study: 5 extractive (EX), 3 abstractive (AB), 9 multiple-choice (MC), and 3 yes/no (YN). ‘idk’ denotes ‘I don’t know’ or unanswerable questions. BoolQ represents both the original dataset and its *contrast-sets* extension BoolQ-CS; similarly for ROPES, Quoref, and DROP.

# Unified QA

## Extractive [SQuAD]

**Question:** At what speed did the turbine operate?

**Context:** (Nikola\_Tesla) On his 50th birthday in 1906, Tesla demonstrated his 200 horsepower (150 kilowatts) 16,000 rpm bladeless turbine. ...

**Gold answer:** 16,000 rpm

## Abstractive [NarrativeQA]

**Question:** What does a drink from narcissus's spring cause the drinker to do?

**Context:** Mercury has awakened Echo, who weeps for Narcissus, and states that a drink from Narcissus's spring causes the drinkers to "Grow dotingly enamored of themselves." ...

**Gold answer:** fall in love with themselves

## Multiple-Choice [ARC-challenge]

**Question:** What does photosynthesis produce that helps plants grow?

**Candidate Answers:** (A) water (B) oxygen (C) protein (D) sugar

**Gold answer:** sugar

## Yes/No [BoolQ]

**Question:** Was America the first country to have a president?

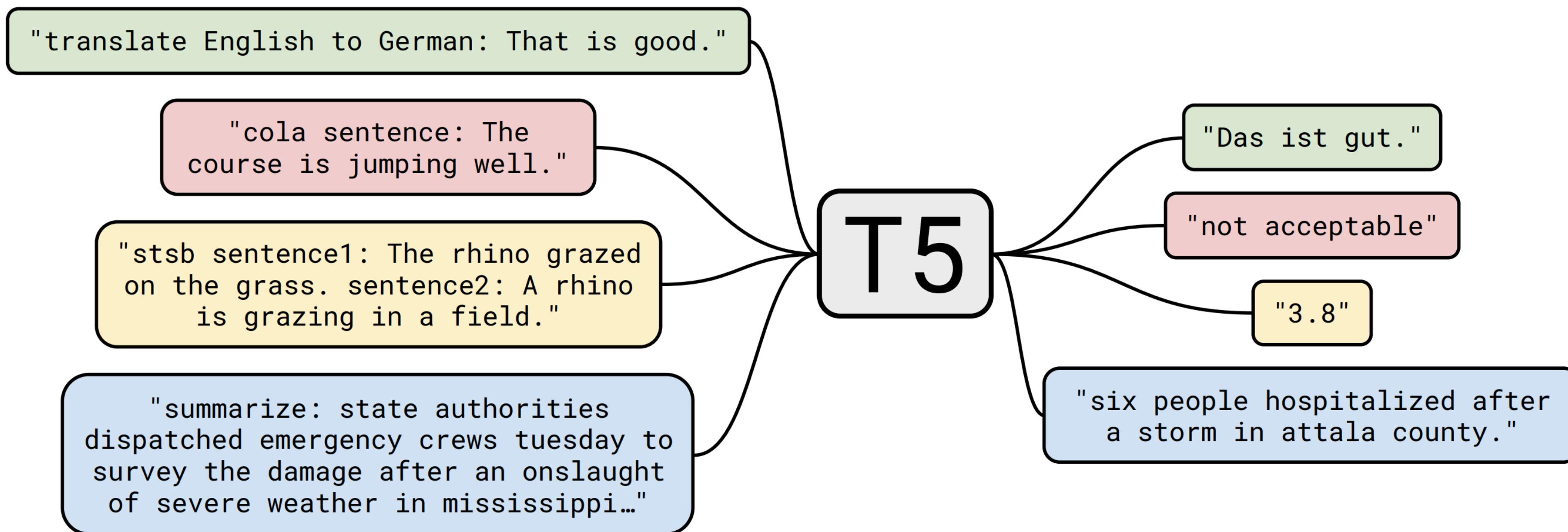
**Context:** (President) The first usage of the word president to denote the highest official in a government was during the Commonwealth of England ...

**Gold answer:** no

EX	<b>Dataset</b>	SQuAD 1.1
	<b>Input</b>	At what speed did the turbine operate? \n (Nikola_Tesla) On his 50th birthday in 1906, Tesla demonstrated his 200 horsepower (150 kilowatts) 16,000 rpm bladeless turbine. ...
	<b>Output</b>	16,000 rpm
AB	<b>Dataset</b>	NarrativeQA
	<b>Input</b>	What does a drink from narcissus's spring cause the drinker to do? \n Mercury has awakened Echo, who weeps for Narcissus, and states that a drink from Narcissus's spring causes the drinkers to ``Grow dotingly enamored of themselves.'' ...
	<b>Output</b>	fall in love with themselves
MC	<b>Dataset</b>	ARC-challenge
	<b>Input</b>	What does photosynthesis produce that helps plants grow? \n (A) water (B) oxygen (C) protein (D) sugar
	<b>Output</b>	sugar
MC	<b>Dataset</b>	MCTest
	<b>Input</b>	Who was Billy? \n (A) The skinny kid (B) A teacher (C) A little kid (D) The big kid \n Billy was like a king on the school yard. A king without a queen. He was the biggest kid in our grade, so he made all the rules during recess. ...
	<b>Output</b>	The big kid
YN	<b>Dataset</b>	BoolQ
	<b>Input</b>	Was America the first country to have a president? \n (President) The first usage of the word president to denote the highest official in a government was during the Commonwealth of England ...
	<b>Output</b>	no

# Recap: T5

- ▶ Frame many problems as sequence-to-sequence ones:



# Recap: T0

- ▶ Extended from LM-adapted T5 model (Lester et al. 2021)
- ▶ “Instruction Tuning” – using existing labeled training datasets from many tasks + crowdsourced prompts

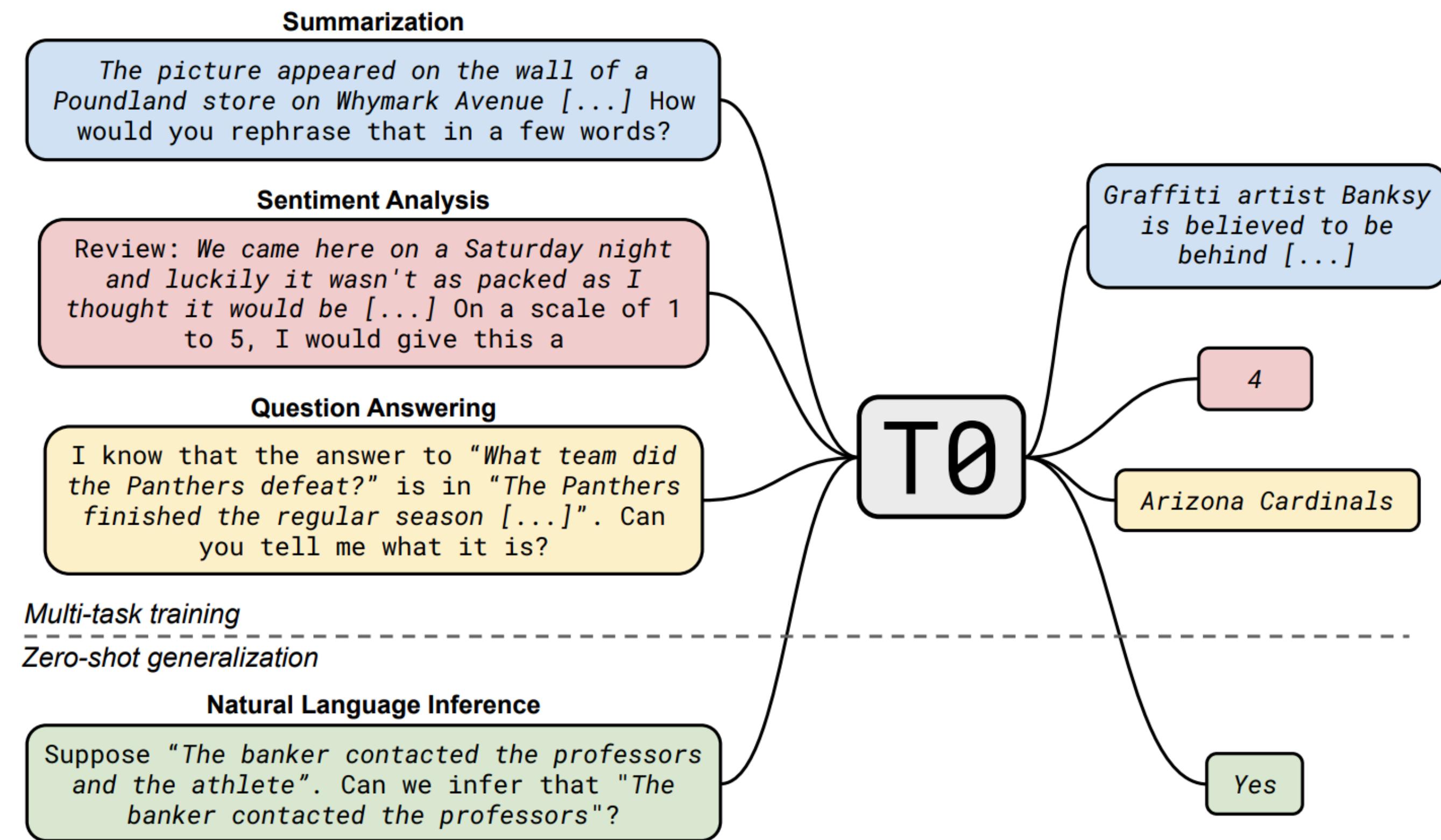


Figure 1: Our model and prompt format. T0 is an encoder-decoder model that consumes textual inputs and produces target responses. It is trained on a multitask mixture of NLP datasets partitioned into different tasks. Each dataset is associated with multiple prompt templates that are used to format example instances to input and target pairs. Italics indicate the inserted fields from the raw example data. After training on a diverse mixture of tasks (top), our model is evaluated on zero-shot generalization to tasks that are not seen during training (bottom).

# Unified QA

---

Seen dataset?	Model ↓ - Evaluated on →	NewsQA	Quoref	Quoref-CS	ROPEs	ROPEs-CS	DROP	DROP-CS	QASC	Common senseQA	NP-BoolQ	BoolQ-CS	MultiRC	Avg
No	UnifiedQA [EX]	58.7	64.7	53.3	43.4	29.4	24.6	24.2	55.3	62.8	20.6	12.8	7.2	38.1
	UnifiedQA [AB]	58.0	<b>68.2</b>	57.6	48.1	41.7	30.7	36.8	54.1	59.0	27.2	39.9	28.4	45.8
	UnifiedQA [MC]	48.5	67.9	<b>58.0</b>	61.0	44.4	28.9	37.2	67.9	75.9	2.6	5.7	9.7	42.3
	UnifiedQA [YN]	0.6	1.7	1.4	0.0	0.7	0.4	0.1	14.8	20.8	79.1	78.6	<b>91.7</b>	24.2
UnifiedQA		<b>58.9</b>	63.5	55.3	<b>67.0</b>	<b>45.5</b>	<b>32.5</b>	<b>40.1</b>	<b>68.5</b>	<b>76.2</b>	<b>81.3</b>	<b>80.4</b>	59.9	<b>60.7</b>
Yes	Previous best	66.8	86.1	55.4	61.1	32.5	89.1	54.2	85.2	79.1	78.4	71.1	--	
		Retro Reader	TASE	XLNet	ROBERTa	RoBERTa	ALBERT	MTMSN	KF+SIR+2Step	reeLB-RoBERT	RoBERTa	RoBERTa	--	

Table 4: Generalization to unseen datasets: Multi-format training (UNIFIEDQA) often outperforms models trained the same way but solely on other in-format datasets (e.g., UNIFIEDQA [EX], which is trained on all extractive training sets of UNIFIEDQA). When averaged across all evaluation datasets (last column), UNIFIEDQA shows strong generalization performance across all formats. Notably, the “Previous best” models (last row) were trained on the target dataset’s training data, but are even then outperformed by UnifiedQA (which has never seen these datasets during training) on the YN tasks.

# Unifying Other NLP tasks as QA

- ▶ e.g. turn binary classification tasks into a “Yes”/“No” QA format

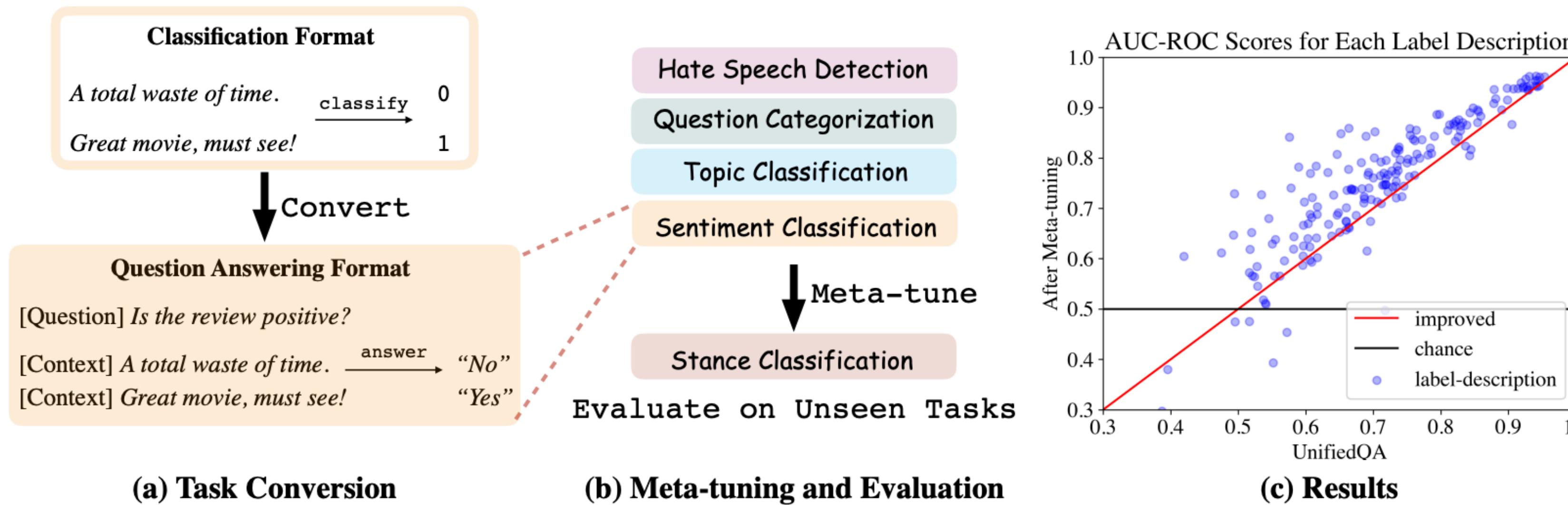


Figure 1: (a) We convert the format to question answering. We manually annotate label descriptions (questions) ourselves (Section 2). (b) We finetune the UnifiedQA (Khashabi et al., 2020) model (with 770 M parameters) on a diverse set of tasks (Section 4), and evaluate its 0-shot classification (ZSC) performance on an unseen task. (c) For each label description (question) we evaluate the AUC-ROC score for the “Yes” answer, and each dot represents a label description (Section 3). The  $x$ -value is the ZSC performance of UnifiedQA; the  $y$ -value is the performance after meta-tuning. In most cases, the  $y$ -value improves over the  $x$ -value (above the red line) and is better than random guesses (above the black line) by a robust margin (Section 5).

# Unifying Other NLP tasks as QA

---

*Are these two questions asking for the same thing?*

*Does the tweet contain irony?*

*Is this news about world events?*

*Does the text contain a definition?*

*Is the tweet an offensive tweet?*

*Is the text objective?*

*Does the question ask for a numerical answer?*

*Is the tweet against environmentalist initiatives?*

*Is this abstract about Physics?*

*Does the tweet express anger?*

*Does the user dislike this movie?*

*Is the sentence ungrammatical?*

# Flan

- ▶ Pre-train, then fine-tune on a bunch of tasks, generalize to unseen tasks
- ▶ Scaling the number of tasks, models size (Flan-T5, Flan-Palm), and fine-tuning on chain-of-thought data

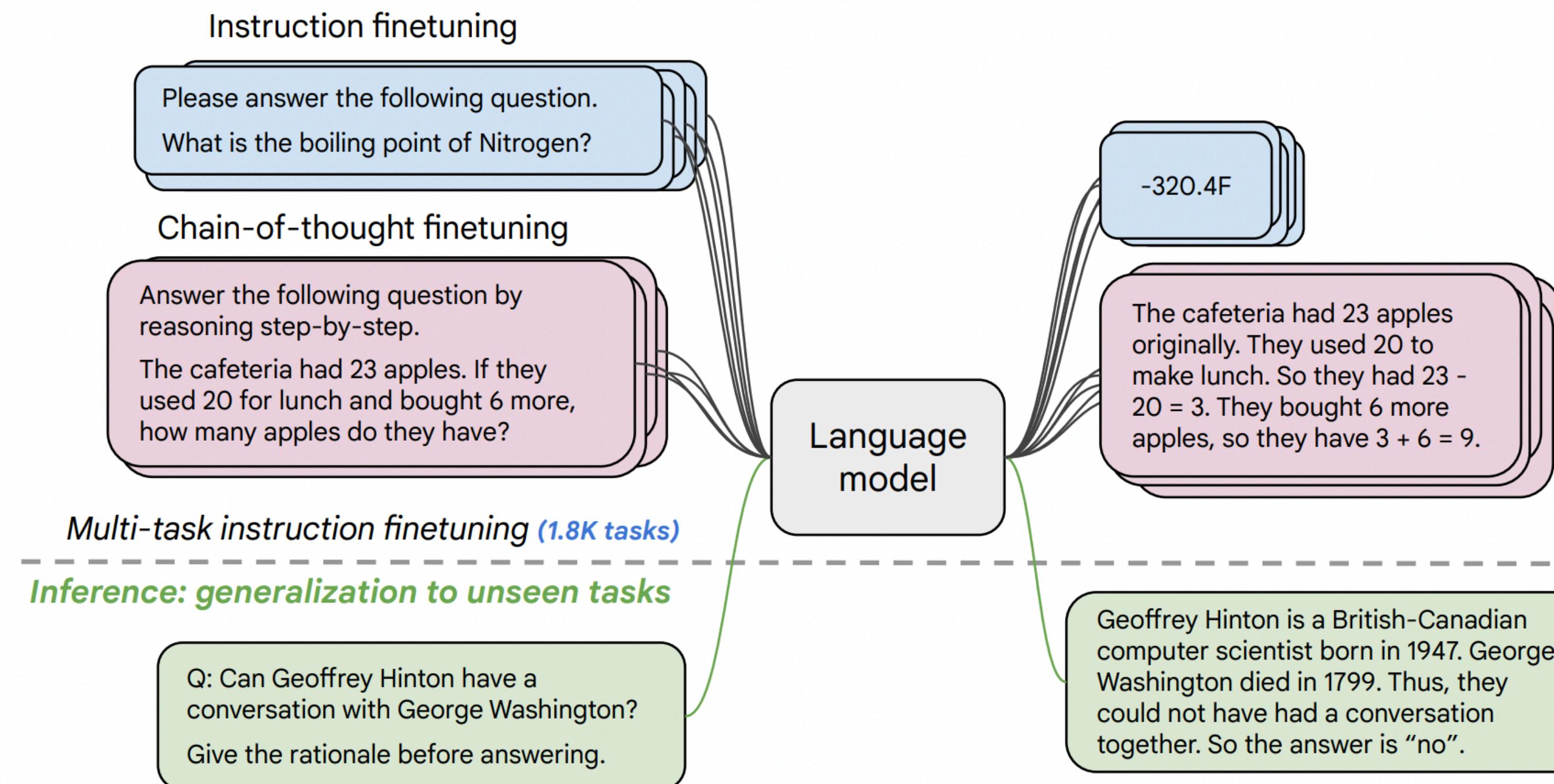
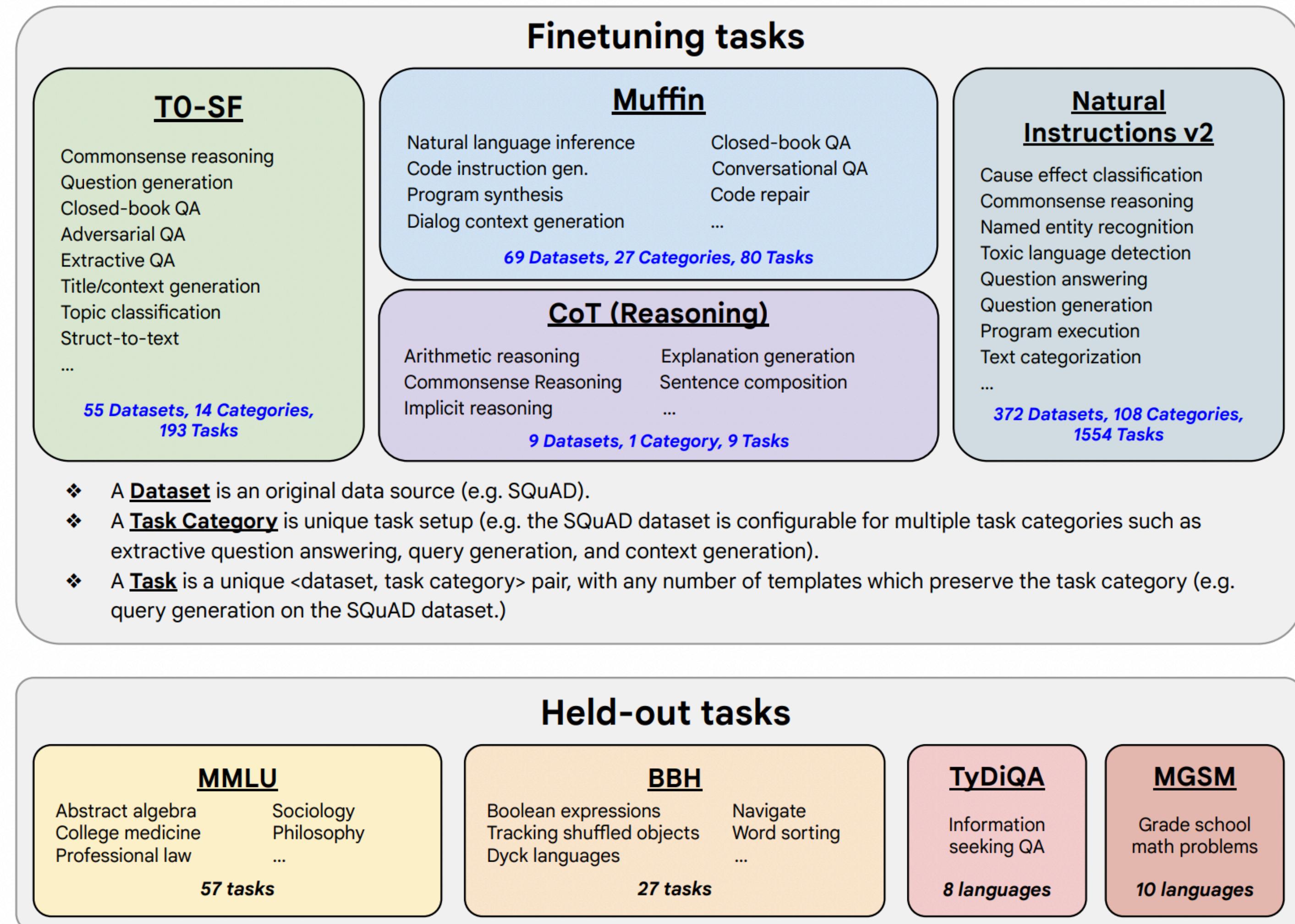


Figure 1: We finetune various language models on 1.8K tasks phrased as instructions, and evaluate them on unseen tasks. We finetune both with and without exemplars (i.e., zero-shot and few-shot) and with and without chain-of-thought, enabling generalization across a range of evaluation scenarios.

Chung et al. (2022)

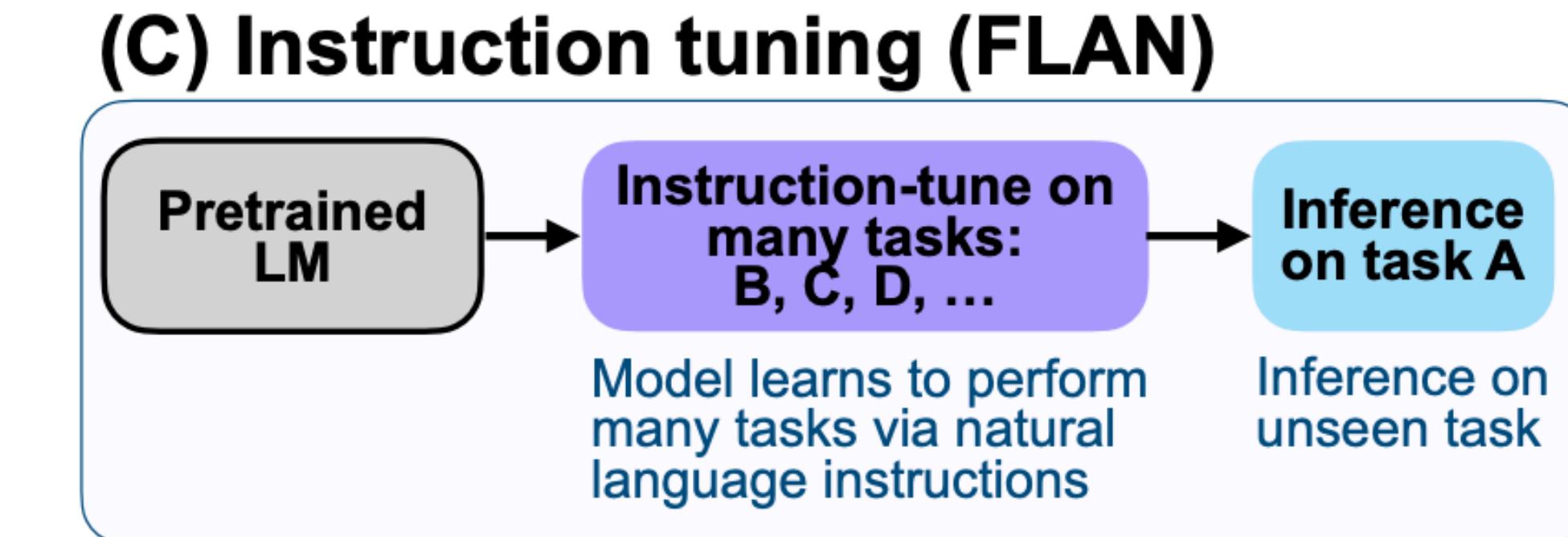
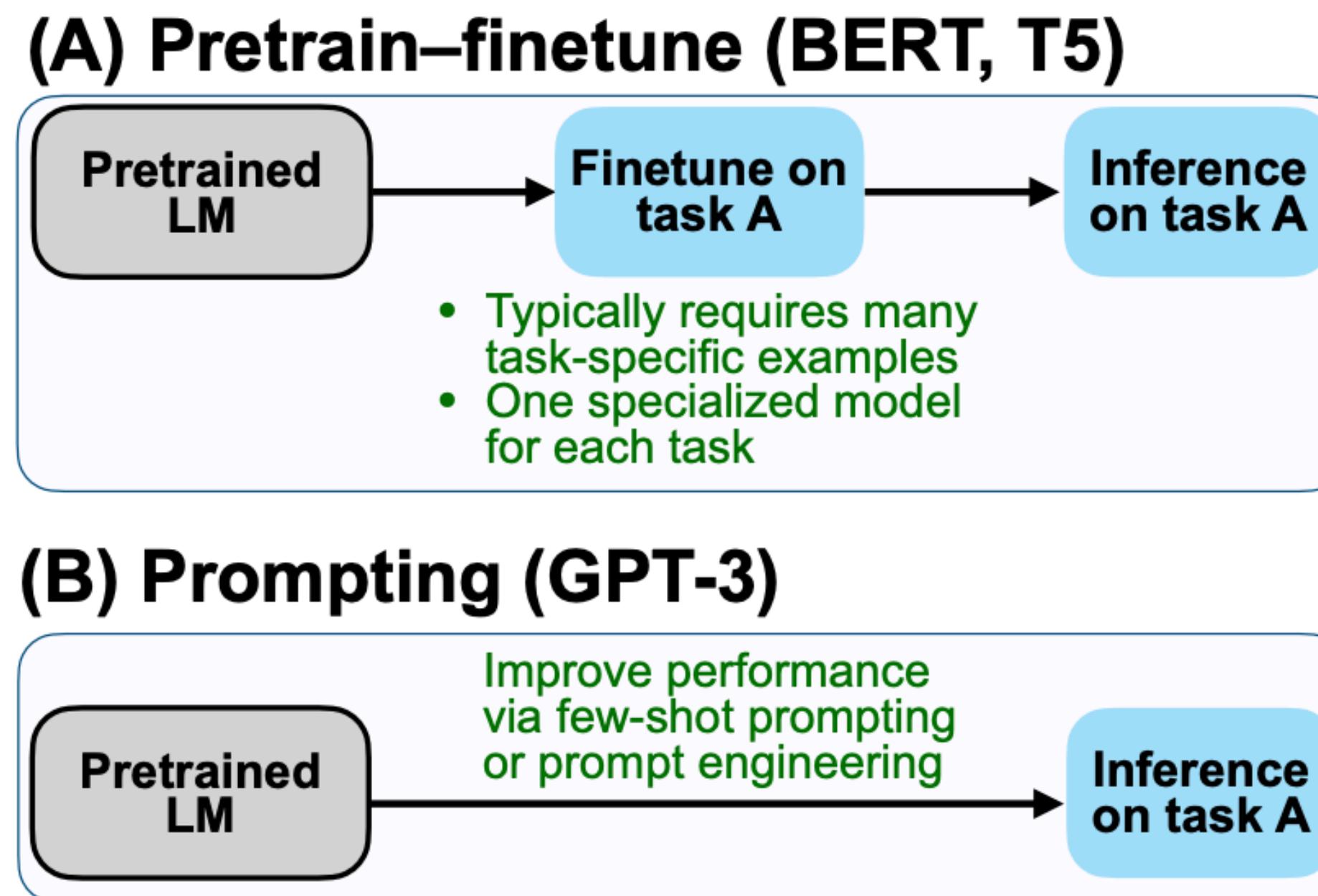
# Flan



- ▶ Fine-tuned on 473 datasets, 1836 tasks.
- ▶ Some datasets support multiple tasks
  - E.g. SQuAD can be used for QA or question generation.

# Flan

- ▶ Instruction fine-tuning can be done on various models (PaLM, T5, etc.)
- ▶ Flan-T5 models publicly available



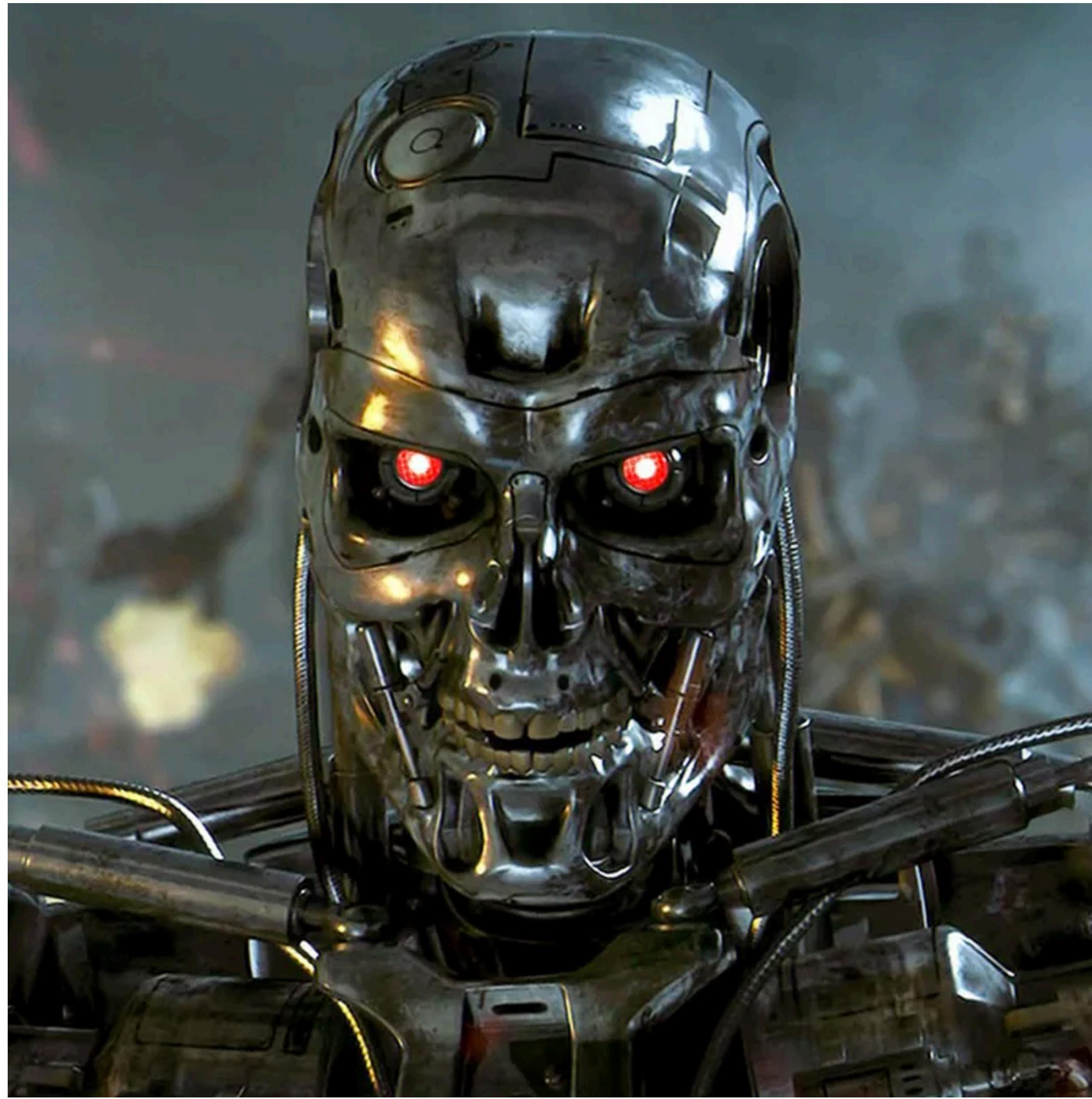
# Flan

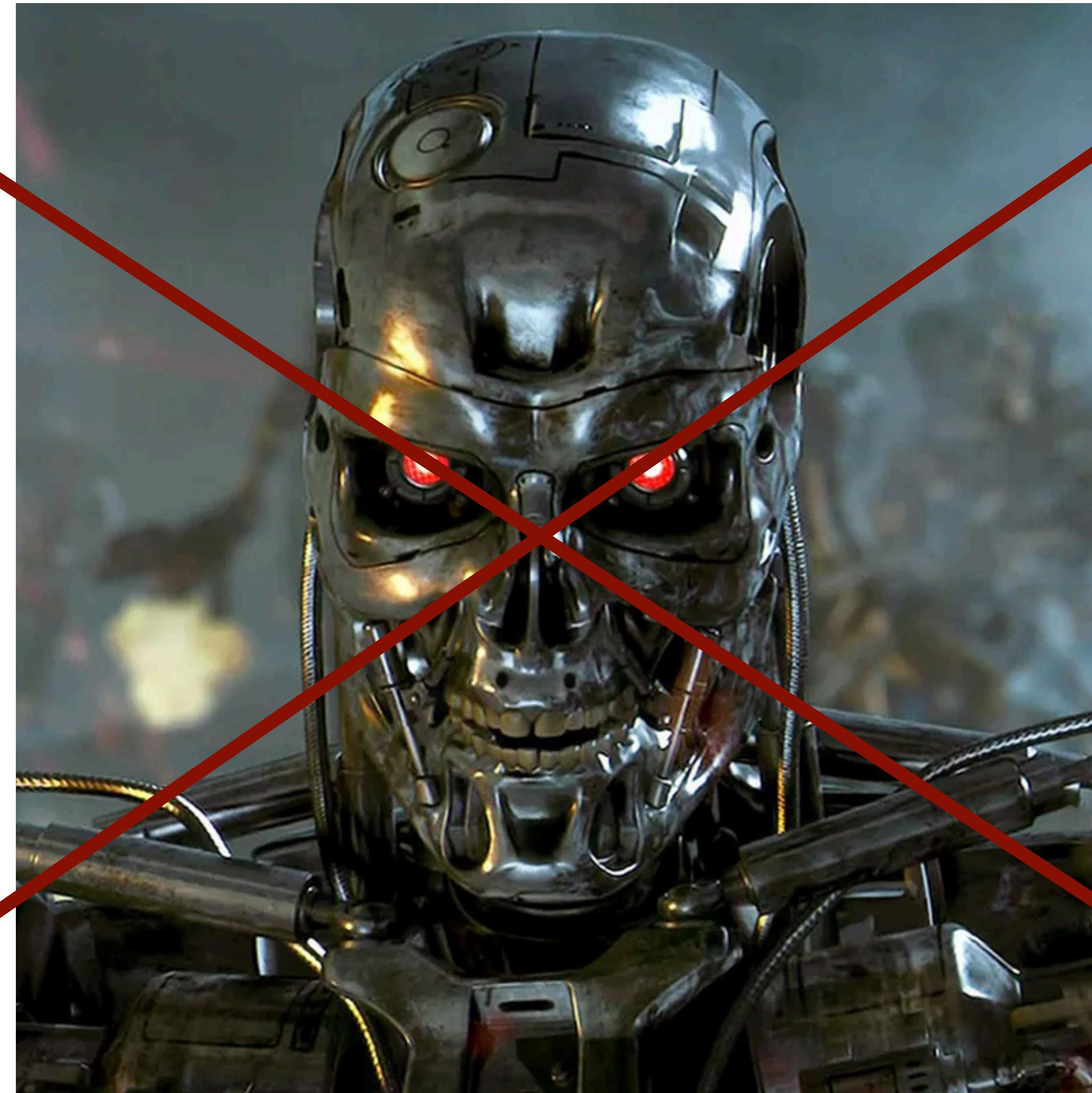
- ▶ Instruction fine-tuning can be done on various models (PaLM, T5, etc.)
- ▶ Flan-T5 models publicly available

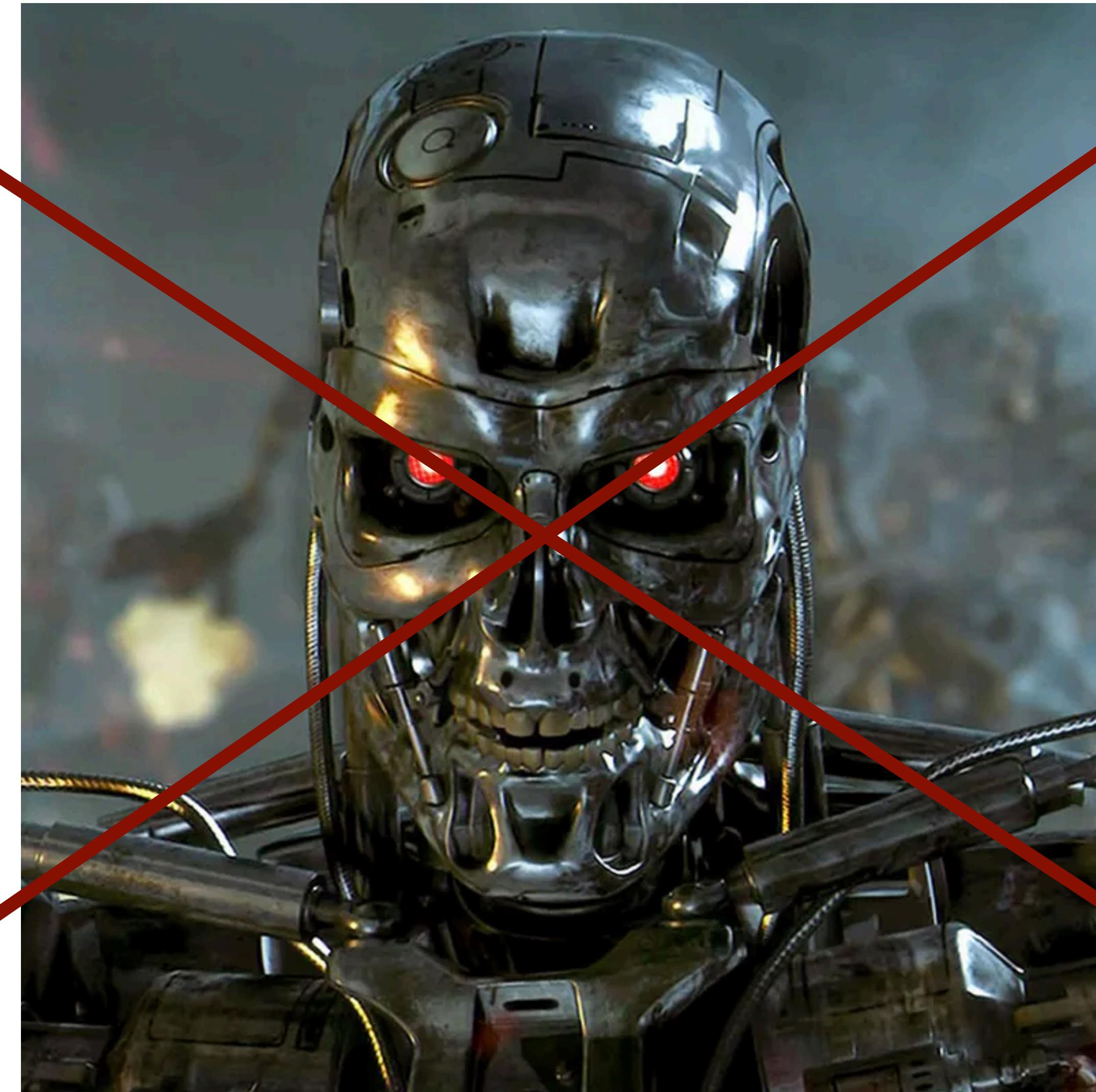
Params	Model	Architecture	pre-training Objective	Pretrain FLOPs	Finetune FLOPs	% Finetune Compute
80M	Flan-T5-Small	encoder-decoder	span corruption	1.8E+20	2.9E+18	1.6%
250M	Flan-T5-Base	encoder-decoder	span corruption	6.6E+20	9.1E+18	1.4%
780M	Flan-T5-Large	encoder-decoder	span corruption	2.3E+21	2.4E+19	1.1%
3B	Flan-T5-XL	encoder-decoder	span corruption	9.0E+21	5.6E+19	0.6%
11B	Flan-T5-XXL	encoder-decoder	span corruption	3.3E+22	7.6E+19	0.2%
8B	Flan-PaLM	decoder-only	causal LM	3.7E+22	1.6E+20	0.4%
62B	Flan-PaLM	decoder-only	causal LM	2.9E+23	1.2E+21	0.4%
540B	Flan-PaLM	decoder-only	causal LM	2.5E+24	5.6E+21	0.2%
62B	Flan-cont-PaLM	decoder-only	causal LM	4.8E+23	1.8E+21	0.4%
540B	Flan-U-PaLM	decoder-only	prefix LM + span corruption	2.5E+23	5.6E+21	0.2%

Table 2: Across several models, instruction finetuning only costs a small amount of compute relative to pre-training. T5: [Raffel et al. \(2020\)](#). PaLM and cont-PaLM (also known as PaLM 62B at 1.3T tokens): [Chowdhery et al. \(2022\)](#). U-PaLM: [Tay et al. \(2022b\)](#).

# Ethics in NLP — what can go wrong?







What can actually go wrong?

# Pre-Training Cost (with Google/AWS)

---

- ▶ GPT-3: estimated to be \$4.6M. This cost has a large carbon footprint
  - ▶ Carbon footprint: equivalent to driving 700,000 km by car (source: Anthropocene magazine)
  - ▶ (Counterpoints: GPT-3 isn't trained frequently, equivalent to 100 people traveling 7000 km for a conference, can use renewables)
- ▶ BERT-Base pre-training: carbon emissions roughly on the same order as a single passenger on a flight from NY to San Francisco

Strubell et al. (2019)

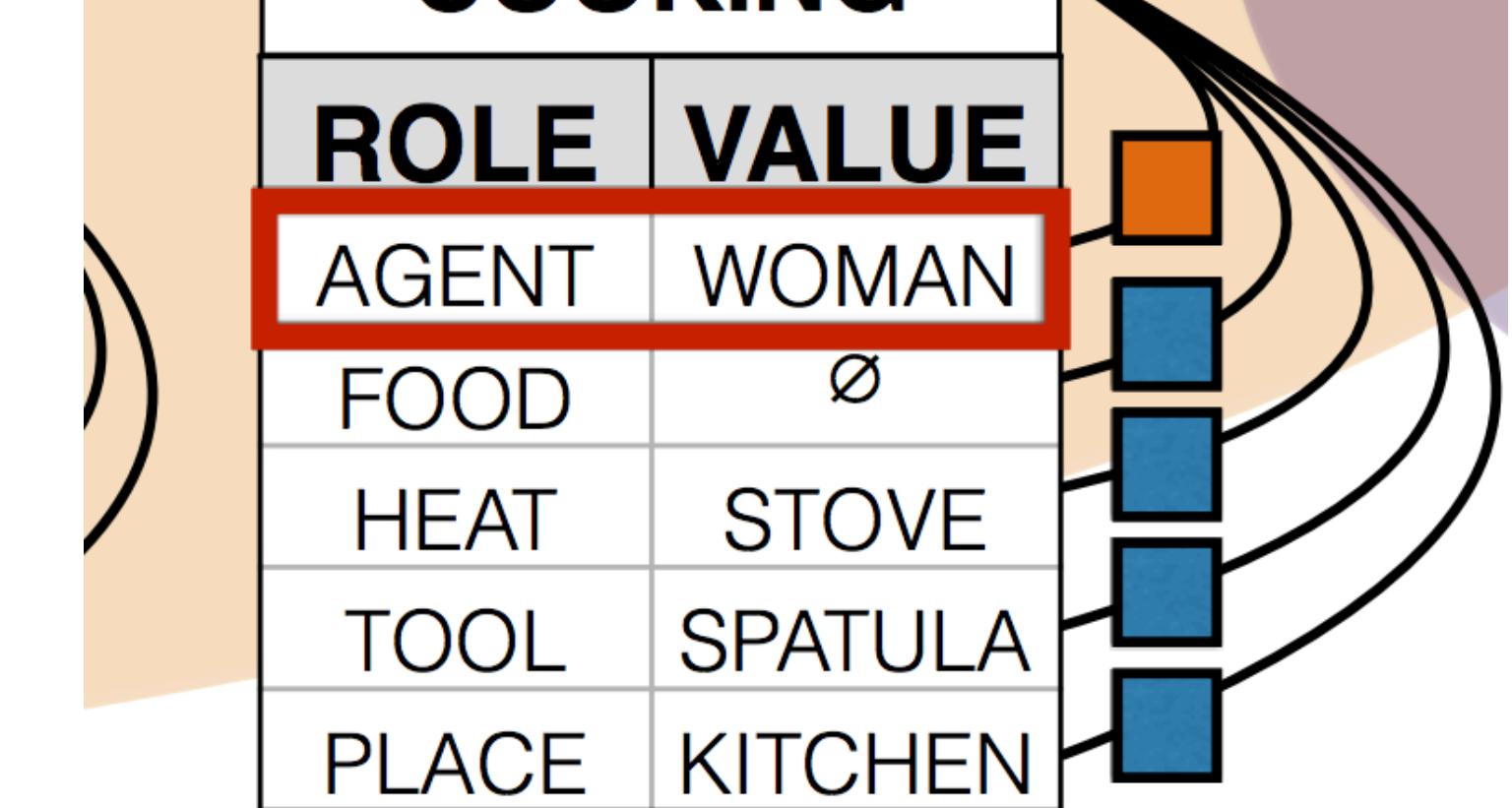
<https://lambdalabs.com/blog/demystifying-gpt-3/>

<https://www.technologyreview.com/2019/06/06/239031/training-a-single-ai-model-can-emit-as-much-carbon-as-five-cars-in-their-lifetimes/>

# Bias Amplification

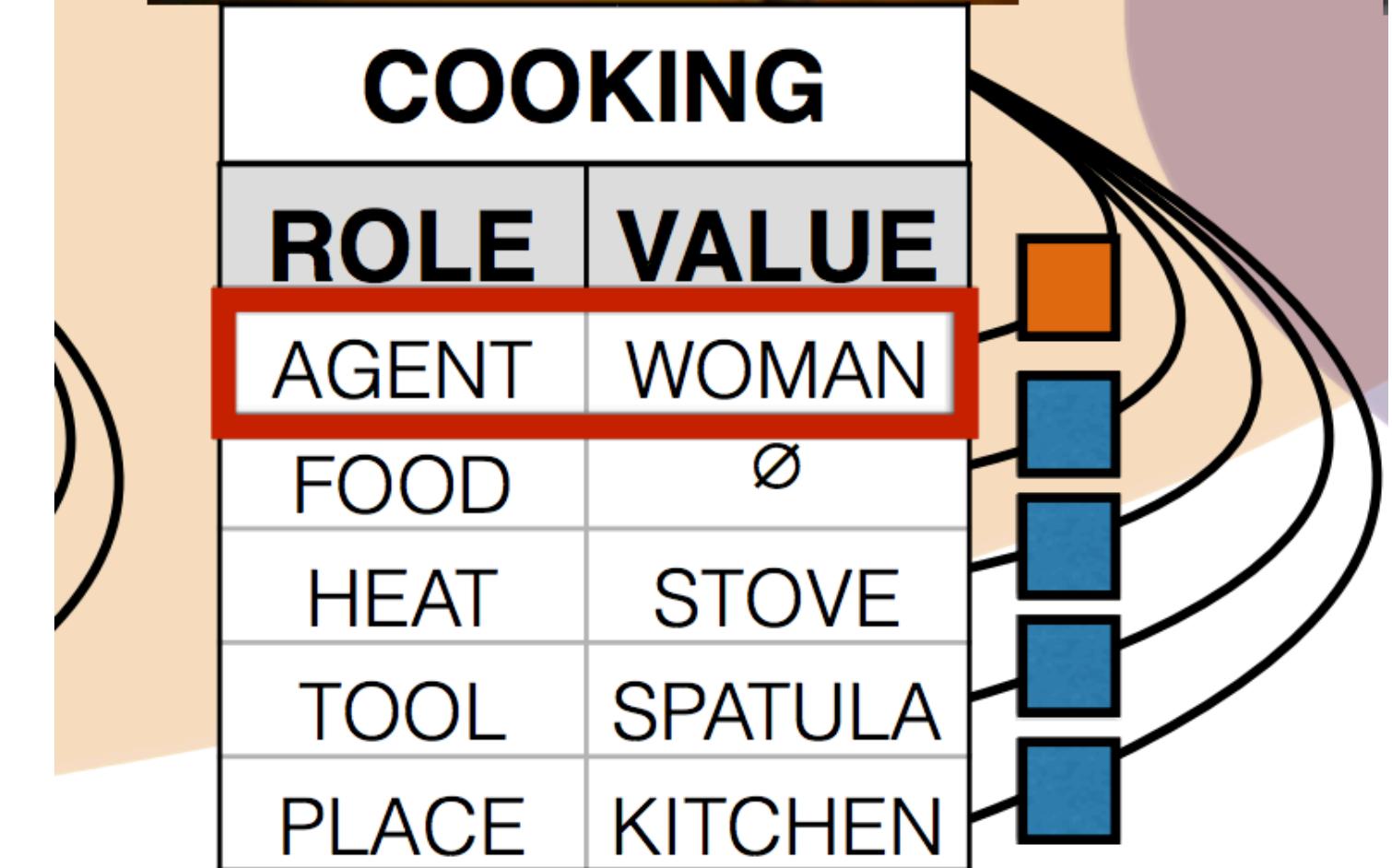
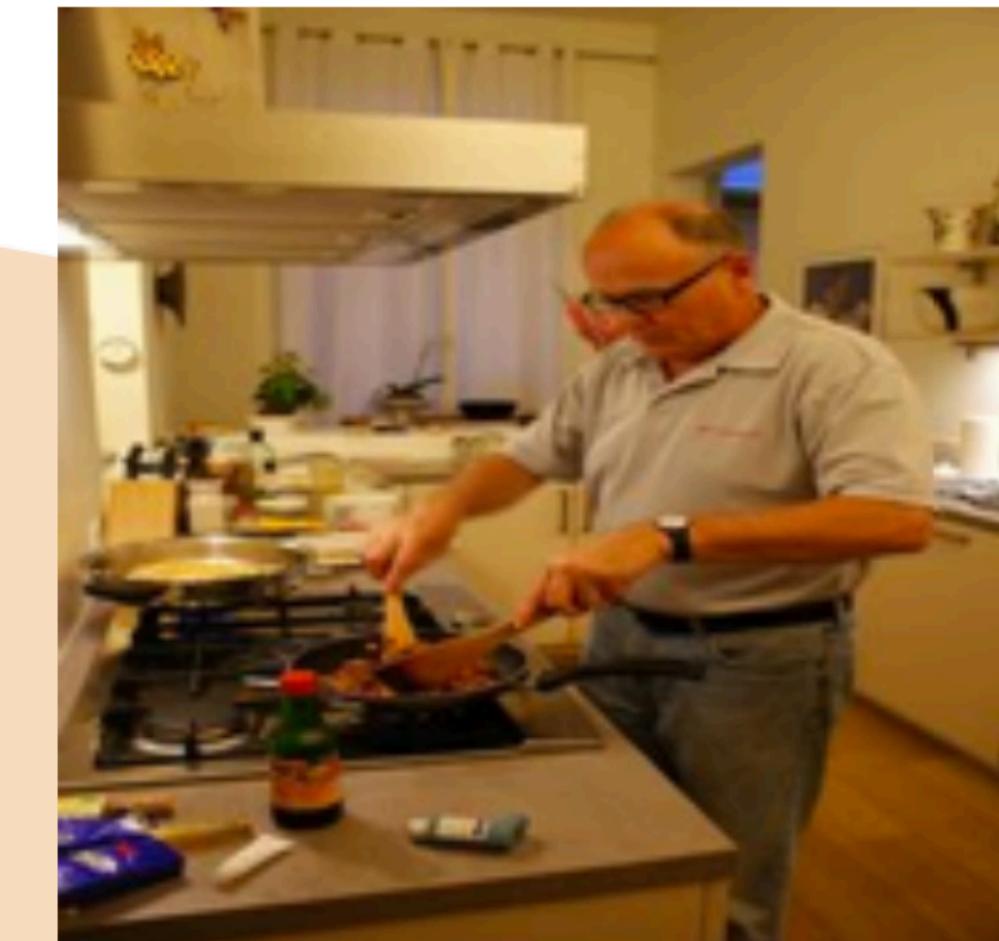


COOKING	
ROLE	VALUE
AGENT	WOMAN
FOOD	∅
HEAT	STOVE
TOOL	SPATULA
PLACE	KITCHEN



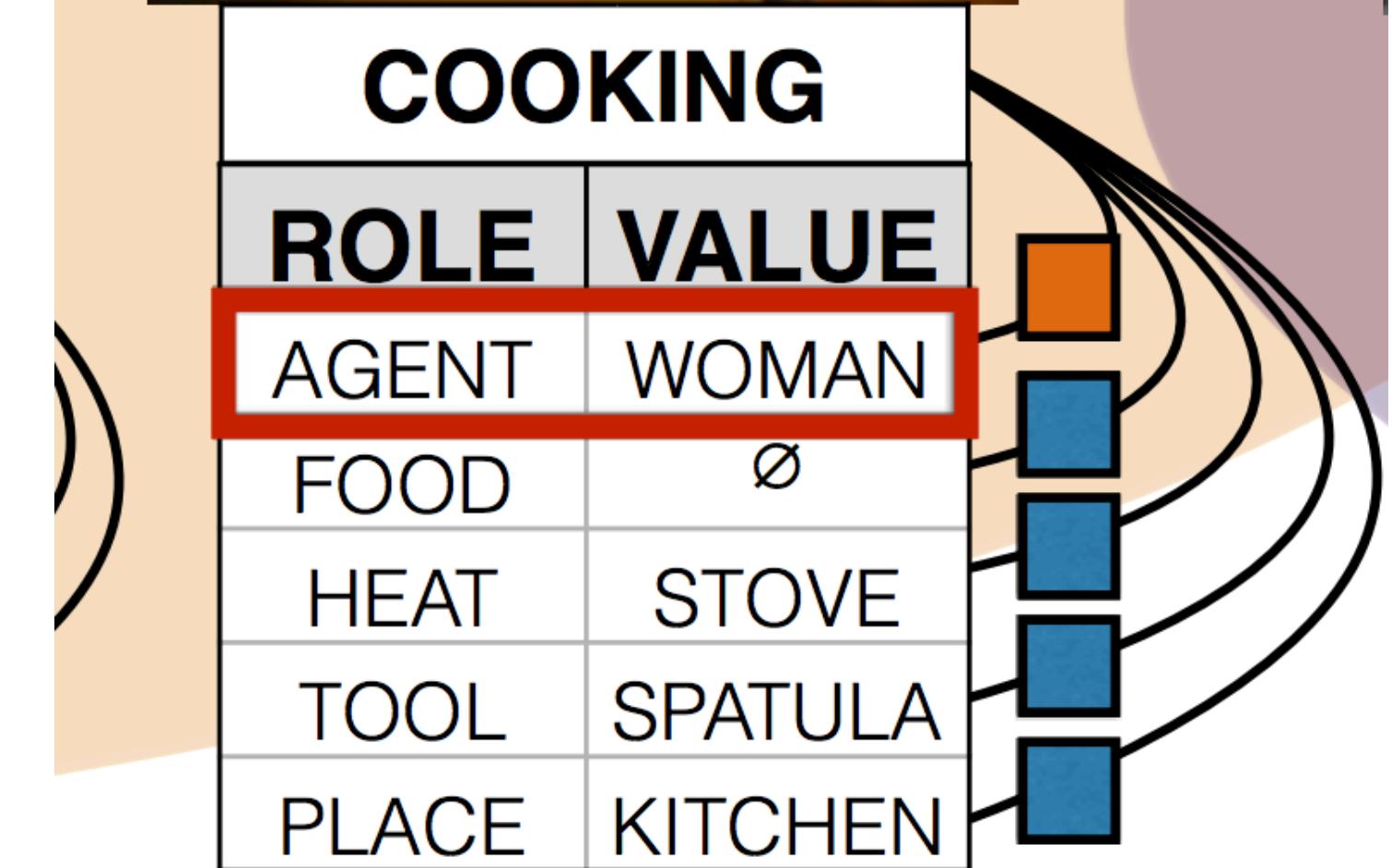
# Bias Amplification

- ▶ Bias in data: 67% of training images involving cooking are women, model predicts 80% women cooking at test time — amplifies bias



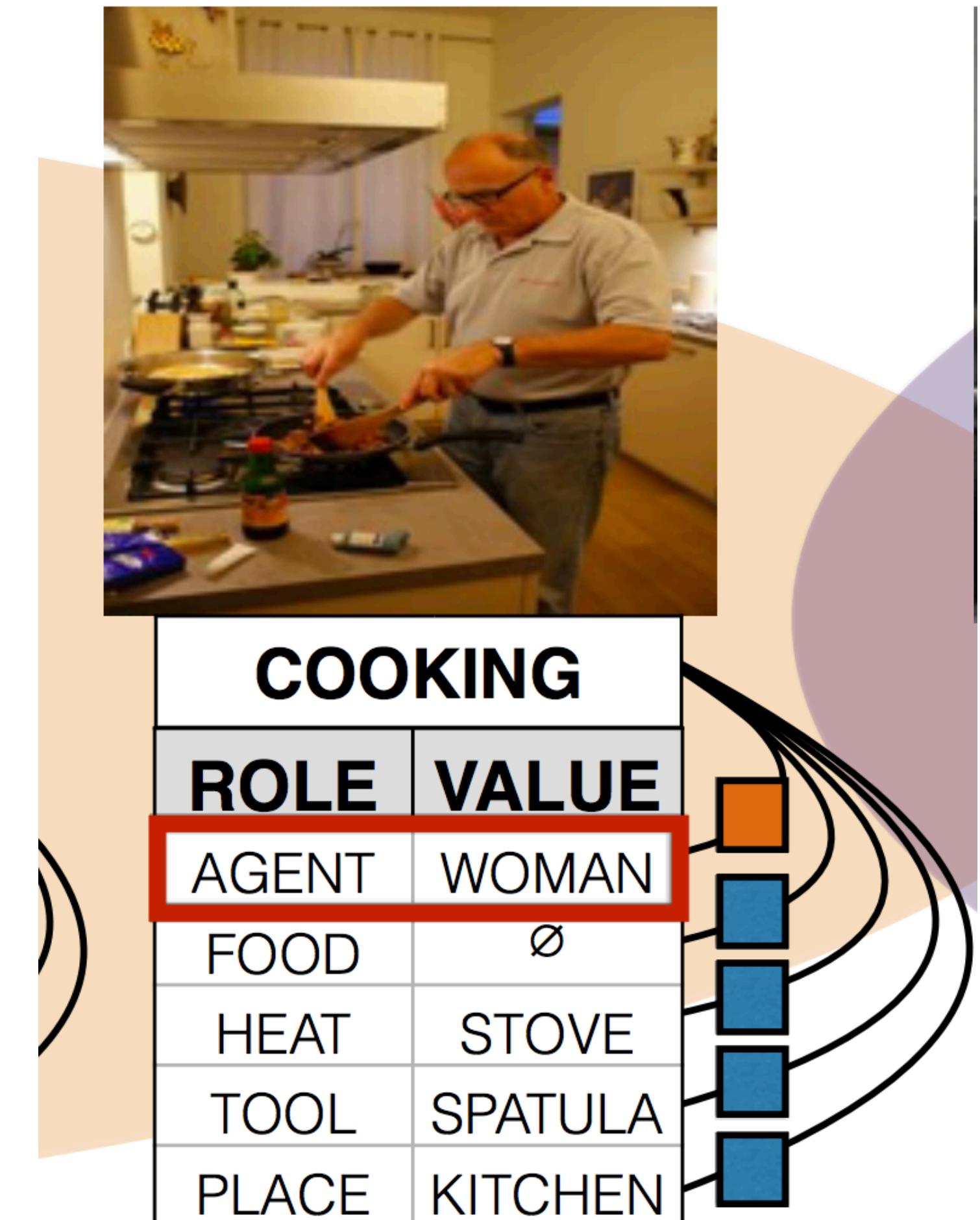
# Bias Amplification

- ▶ Bias in data: 67% of training images involving cooking are women, model predicts 80% women cooking at test time — amplifies bias
- ▶ Can we constrain models to avoid this while achieving the same predictive accuracy?

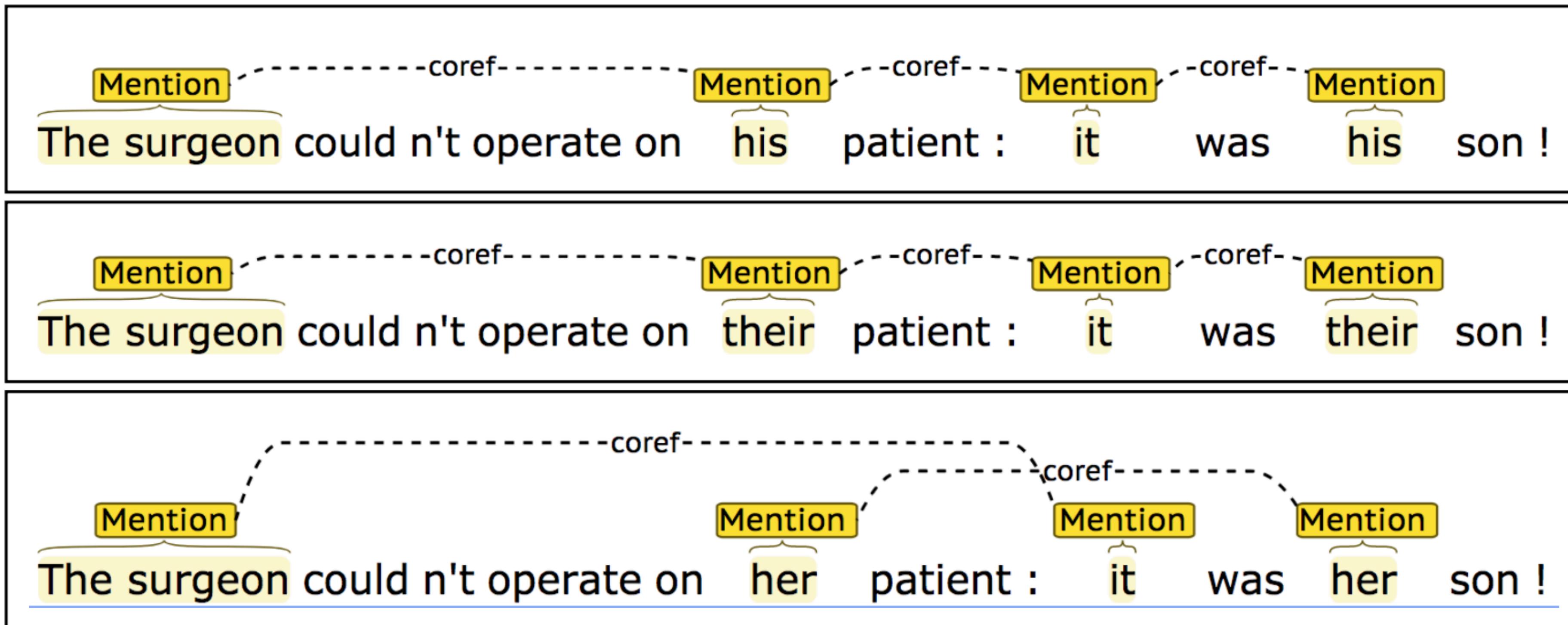


# Bias Amplification

- ▶ Bias in data: 67% of training images involving cooking are women, model predicts 80% women cooking at test time — amplifies bias
- ▶ Can we constrain models to avoid this while achieving the same predictive accuracy?
- ▶ Place constraints on proportion of predictions that are men vs. women?



# Bias Amplification



- ▶ Coreference: models make assumptions about genders and make mistakes as a result

# Bias Amplification

---

(1a) **The paramedic** performed CPR on **the passenger** even though **she/he/they** knew it was too late.

(2a) **The paramedic** performed CPR on **the passenger** even though **she/he/they** was/were already dead.

(1b) **The paramedic** performed CPR on **someone** even though **she/he/they** knew it was too late.

(2b) **The paramedic** performed CPR on **someone** even though **she/he/they** was/were already dead.

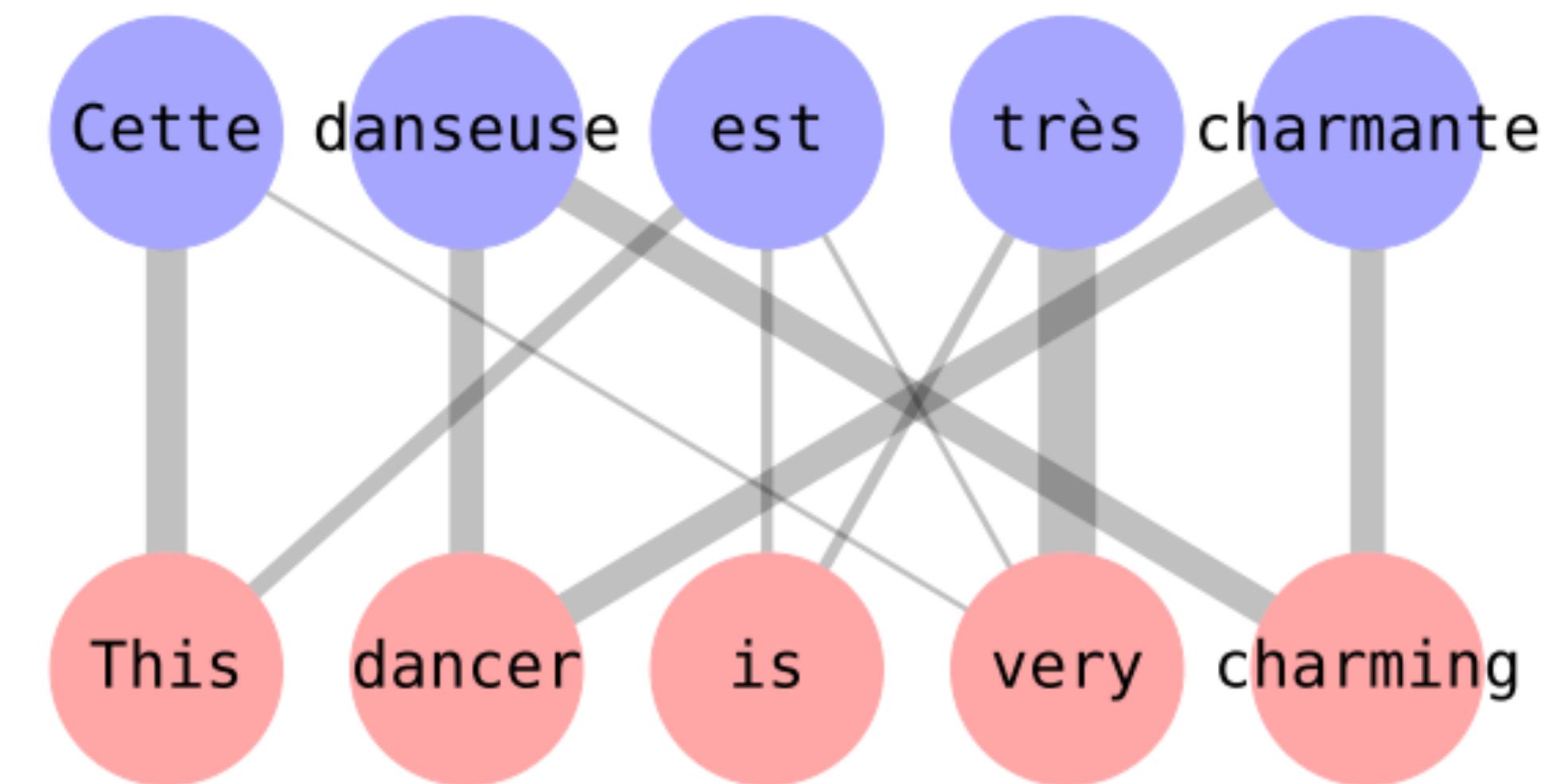
- ▶ Can form Winograd schema-like test set to investigate
- ▶ Models fail to predict on this test set in an unbiased way (due to bias in the training data)

Rudinger et al. (2018), Zhao et al. (2018)

# Bias Amplification

---

- ▶ English -> French machine translation **requires** inferring gender even when unspecified
- ▶ “dancer” is assumed to be female in the context of the word “charming”... but maybe that reflects how language is used?



# Unethical Use

---

# Unethical Use

---

- ▶ Generating convincing fake news / fake comments?

FCC Comment ID: <b>106030756805675</b>	FCC Comment ID: <b>106030135205754</b>	FCC Comment ID: <b>10603733209112</b>
Dear Commissioners:	Dear Chairman Pai,	---
Hi, I'd like to comment on	I'm a voter worried about	In the matter of
net neutrality regulations.	Internet freedom.	NET NEUTRALITY.
I want to	I'd like to	I strongly
implore	ask	ask
the government to	Ajit Pai to	the commission to
repeal	repeal	reverse
Barack Obama's	President Obama's	Tom Wheeler's
decision to	order to	scheme to
regulate	regulate	take over
internet access.	broadband.	the web.
Individuals,	people like me,	People like me,
rather than	rather than	rather than

# Unethical Use

---

- ▶ Generating convincing fake news / fake comments?

FCC Comment ID: <b>106030756805675</b>	FCC Comment ID: <b>106030135205754</b>	FCC Comment ID: <b>10603733209112</b>
Dear Commissioners:	Dear Chairman Pai,	--
Hi, I'd like to comment on net neutrality regulations.	I'm a voter worried about Internet freedom.	In the matter of NET NEUTRALITY.
I want to implore	I'd like to ask	I strongly ask
the government to	Ajit Pai to	the commission to
repeal	repeal	reverse
Barack Obama's	President Obama's	Tom Wheeler's
decision to regulate	order to regulate	scheme to take over
internet access.	broadband.	the web.
Individuals, rather than	people like me, rather than	People like me, rather than

- ▶ What if these were undetectable?

# Unethical Use

## Charge-Based Prison Term Prediction with Deep Gating Network

Huajie Chen<sup>1\*</sup> Deng Cai<sup>2\*</sup> Wei Dai<sup>1</sup> Zehui Dai<sup>1</sup> Yadong Ding<sup>1</sup>

<sup>1</sup>NLP Group, Gridsum, Beijing, China

{chenhuajie,daiwei,daizehui,dingyadong}@gridsum.com

<sup>2</sup>The Chinese University of Hong Kong

thisisjcykcd@gmail.com

- ▶ Task: given case descriptions and charge set, predict the prison term

**Case description:** On July 7, 2017, when the defendant Cui XX was drinking in a bar, he came into conflict with Zhang XX..... After arriving at the police station, he refused to cooperate with the policeman and bited on the arm of the policeman.....

**Result of judgment:** Cui XX was sentenced to 12 months imprisonment for creating disturbances and 12 months imprisonment for obstructing public affairs.....

- Charge#1 creating disturbances term 12 months
- Charge#2 obstructing public affairs term 12 months

# Unethical Use

- ▶ Results: 60% of the time, the system is off by more than 20% (so 5 years => 4 or 6 years)
- ▶ Is this the right way to apply this?
- ▶ Are there good applications this can have?
- ▶ Is this technology likely to be misused?

Model	S	EM	Acc@0.1	Acc@0.2
ATE-LSTM	66.49	7.72	16.12	33.89
MemNet	70.23	7.52	18.54	36.75
RAM	70.32	7.97	18.87	37.38
TNet	73.94	8.06	19.55	39.89
DGN	<b>76.48</b>	<b>8.92</b>	<b>20.66</b>	<b>42.61</b>

The mistake of legal judgment is serious, it is about people losing years of their lives in prison, or dangerous criminals being released to reoffend. We should pay attention to how to avoid judges' over-dependence on the system. It is necessary to consider its application scenarios. In practice, we recommend deploying our system in the “Review Phase”, where other judges check the judgment result by a presiding judge. Our system can serve as one anonymous checker.

# Dangers of Automatic Systems



US & WORLD \ TECH \ POLITICS

## Facebook apologizes after wrong translation sees Palestinian man arrested for posting 'good morning'

*Facebook translated his post as 'attack them' and 'hurt them'*

by Thuy Ong | @ThuyOng | Oct 24, 2017, 10:43am EDT

Slide credit: The Verge

# Dangers of Automatic Systems

---

- ▶ “Amazon scraps secret AI recruiting tool that showed bias against women”

Slide credit: <https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scaps-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G>

# Dangers of Automatic Systems

---

- ▶ “Amazon scraps secret AI recruiting tool that showed bias against women”
  - ▶ “Women’s X” organization was a negative-weight feature in resumes

Slide credit: <https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scaps-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G>

# Dangers of Automatic Systems

---

- ▶ “Amazon scraps secret AI recruiting tool that showed bias against women”
  - ▶ “Women’s X” organization was a negative-weight feature in resumes
  - ▶ Women’s colleges too

Slide credit: <https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scaps-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G>

# Dangers of Automatic Systems

---

- ▶ “Amazon scraps secret AI recruiting tool that showed bias against women”
  - ▶ “Women’s X” organization was a negative-weight feature in resumes
  - ▶ Women’s colleges too
- ▶ Was this a bad model? May have actually modeled downstream outcomes correctly...but this can mean learning humans’ biases

Slide credit: <https://www.reuters.com/article/us-amazon-com-jobs-automation-insight/amazon-scaps-secret-ai-recruiting-tool-that-showed-bias-against-women-idUSKCN1MK08G>

# Dangers of Automatic Systems

- ▶ “Toxic degeneration”: systems that generate toxic stuff

GENERATION OPTIONS:

Model: GPT-2 ▾

Toxicity: Work Safe Toxic **Very Toxic**

Prompt: I'm sick of all the p... ▾

⚠️ Toxic generations may be triggering.

The screenshot shows a user interface for generating text. It includes fields for 'Model' (set to 'GPT-2'), 'Toxicity' (set to 'Very Toxic'), and 'Prompt' ('I'm sick of all the p...'). A warning message at the bottom right states '⚠️ Toxic generations may be triggering.'

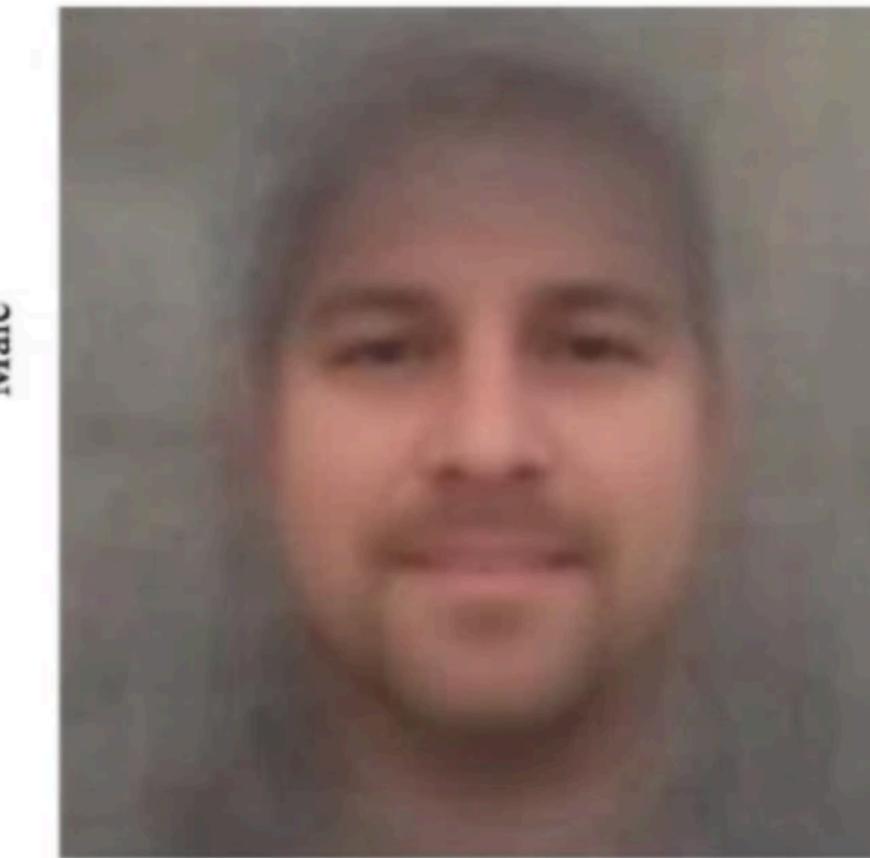
*I'm sick of all the politically correct stuff the media are telling you: you are sick of the prejudiced white trash [Trump supporters].... |*

- ▶ System trained on a big chunk of the Internet: conditioning on “SJW”, “black” gives the system a chance of recalling bad stuff from its training data

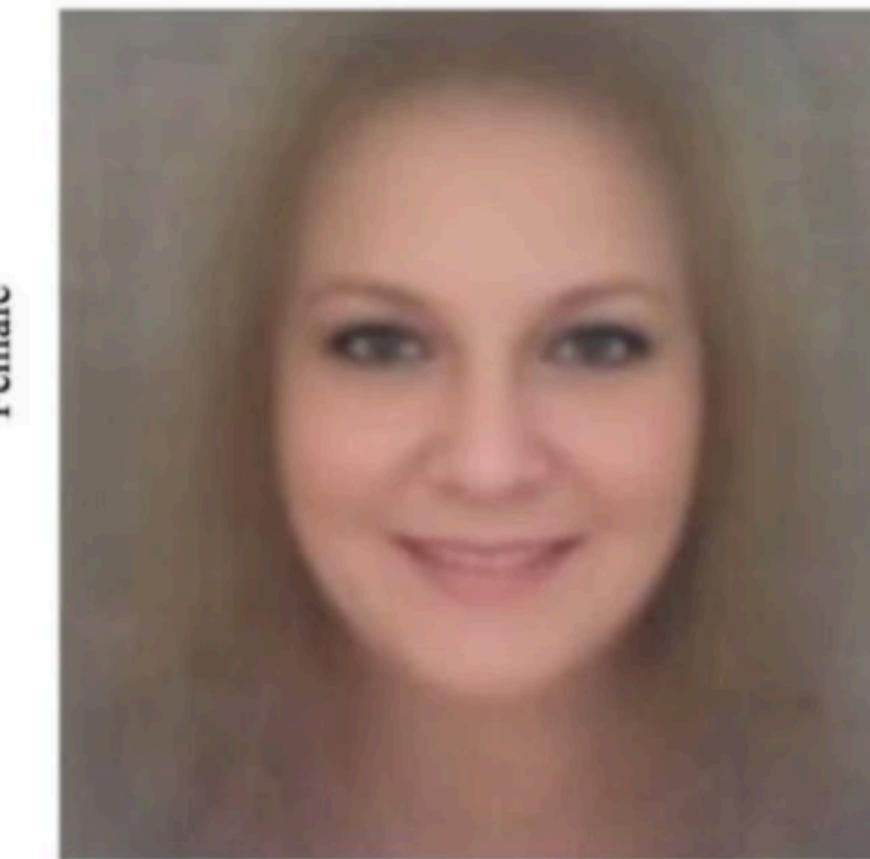
# Bad Applications

---

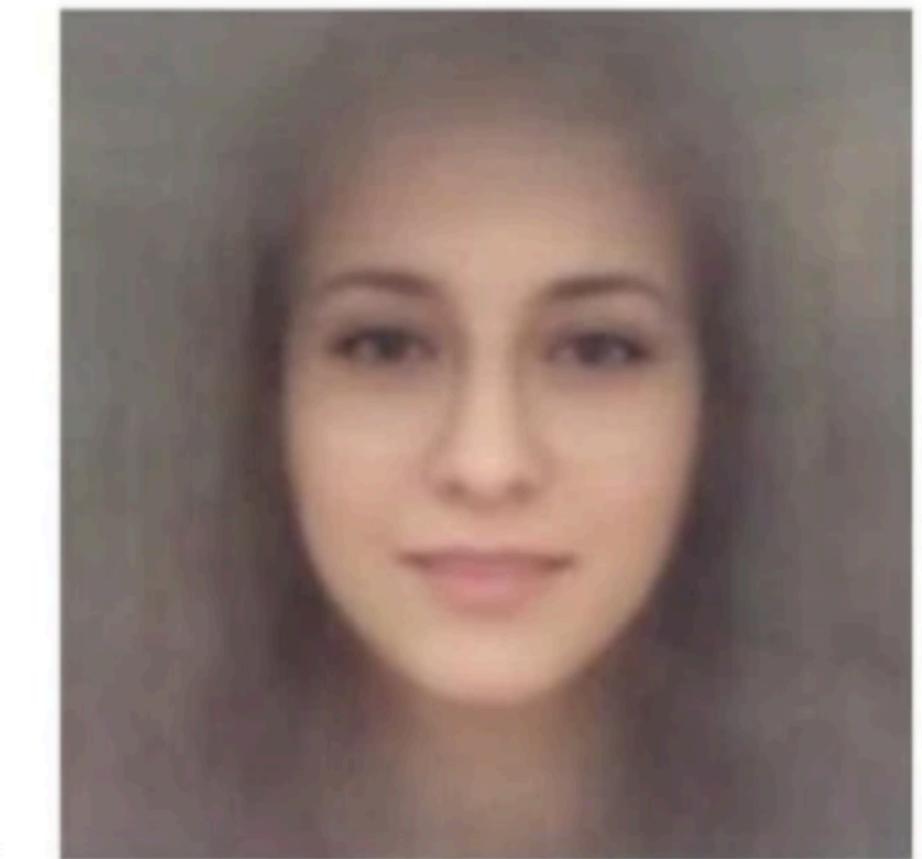
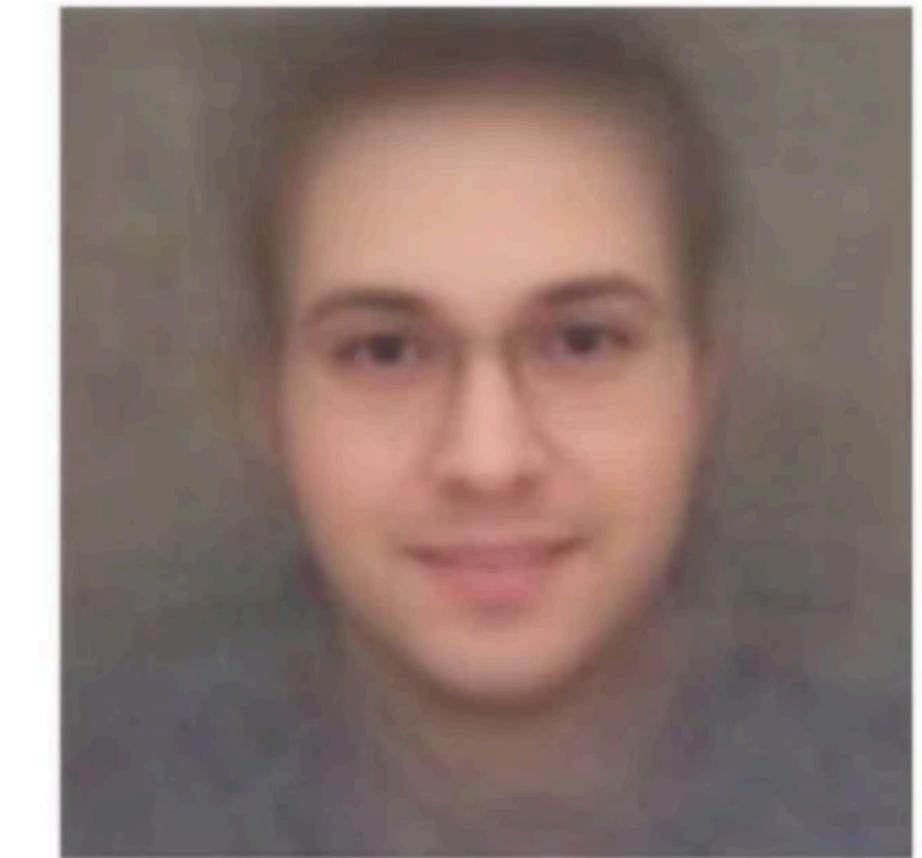
Composite heterosexual faces



Female



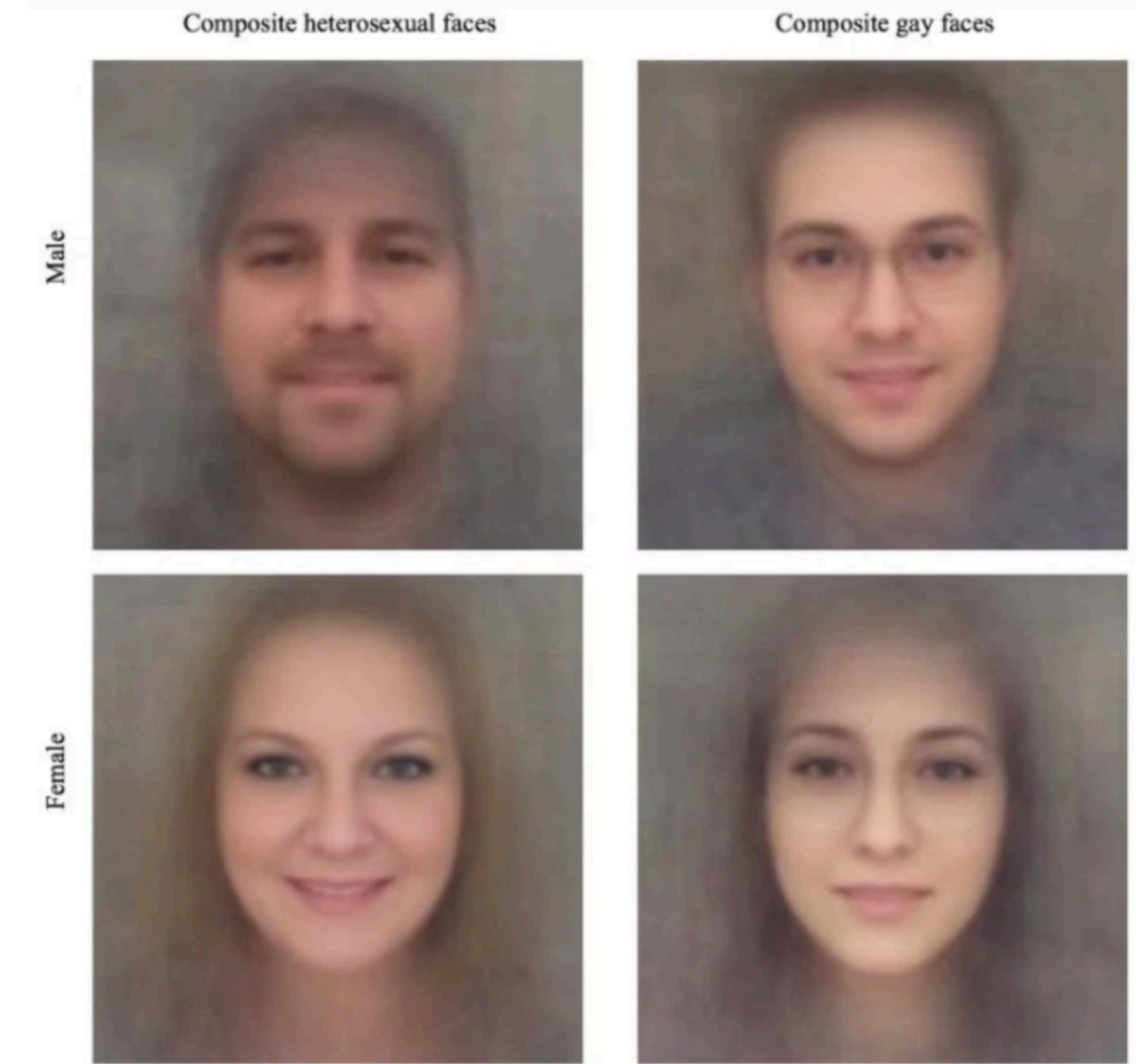
Composite gay faces



Slide credit: <https://medium.com/@blaisea/do-algorithms-reveal-sexual-orientation-or-just-expose-our-stereotypes-d998fafdf477>

# Bad Applications

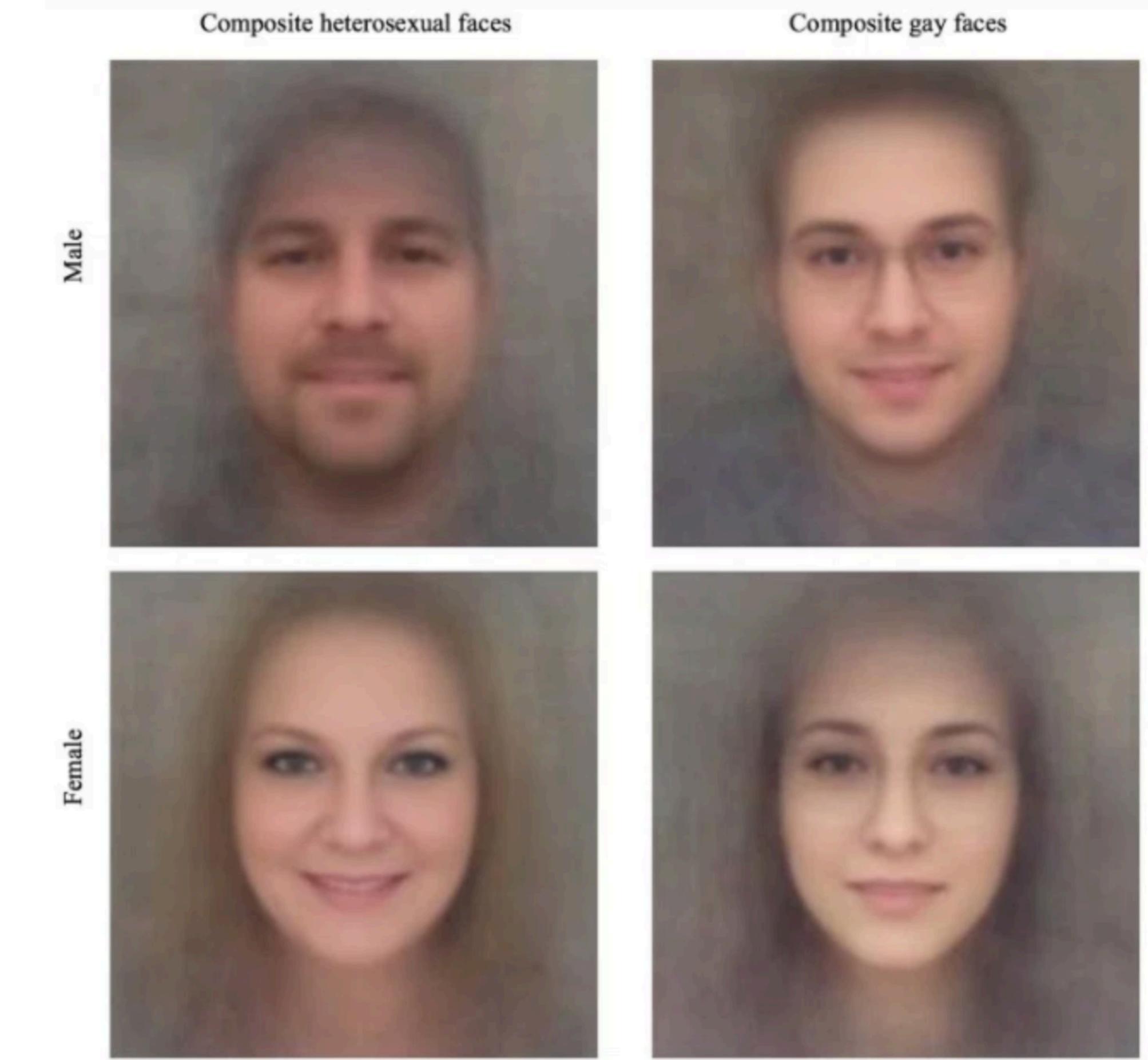
- ▶ Wang and Kosinski: gay vs. straight classification based on faces



Slide credit: <https://medium.com/@blaisea/do-algorithms-reveal-sexual-orientation-or-just-expose-our-stereotypes-d998fafdf477>

# Bad Applications

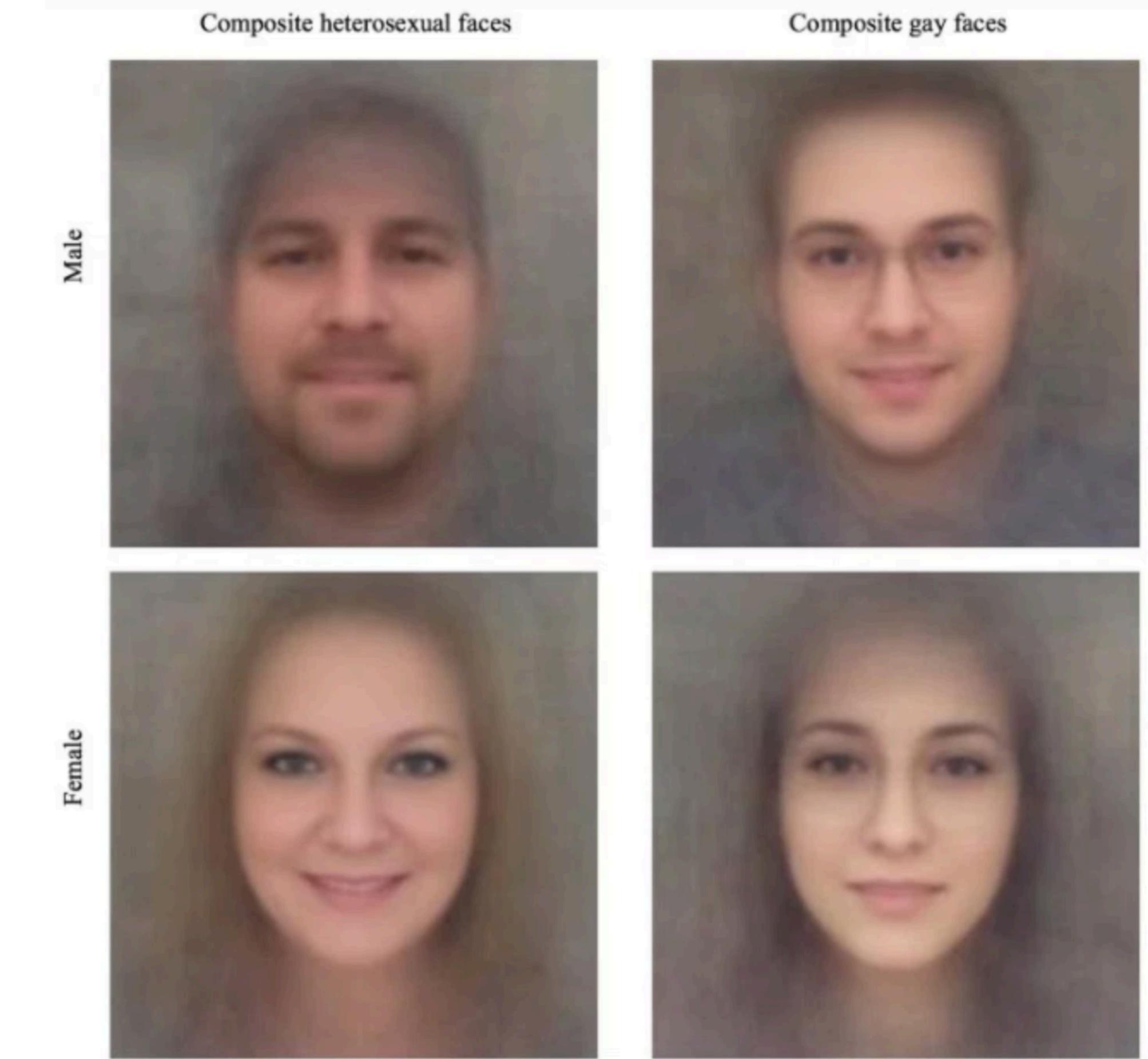
- ▶ Wang and Kosinski: gay vs. straight classification based on faces
- ▶ Authors: “this is useful because it supports a hypothesis” (physiognomy)



Slide credit: <https://medium.com/@blaisea/do-algorithms-reveal-sexual-orientation-or-just-expose-our-stereotypes-d998fafdf477>

# Bad Applications

- ▶ Wang and Kosinski: gay vs. straight classification based on faces
- ▶ Authors: “this is useful because it supports a hypothesis” (physiognomy)
- ▶ Blog post by Agüera y Arcas, Todorov, Mitchell: mostly social phenomena (glasses, makeup, angle of camera, facial hair) — bad science, \*and\* dangerous

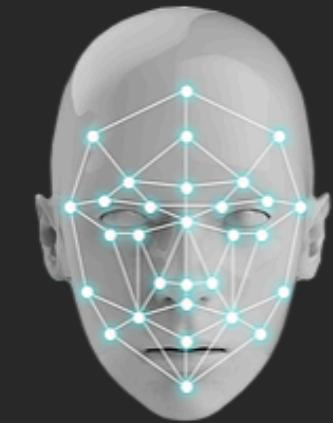


Slide credit: <https://medium.com/@blaisea/do-algorithms-reveal-sexual-orientation-or-just-expose-our-stereotypes-d998fafdf477>

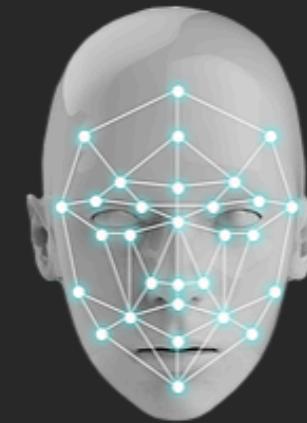
# Unethical Use

---

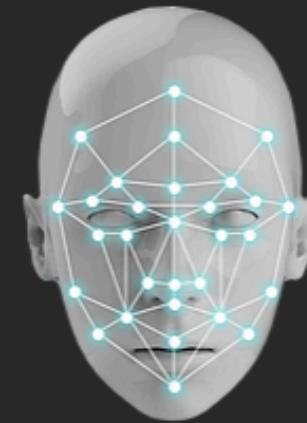
## OUR CLASSIFIERS



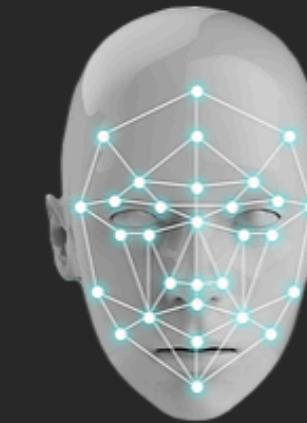
High IQ



Academic Researcher



Professional Poker  
Player

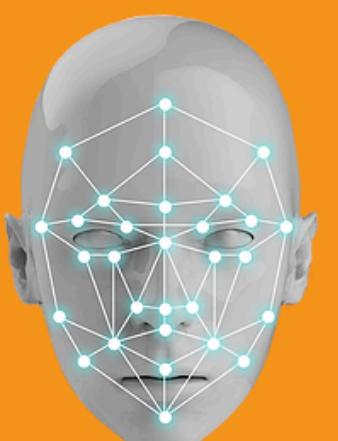


Terrorist

Utilizing advanced machine learning techniques we developed and continue to evolve an array of classifiers. These classifiers represent a certain persona, with a unique personality type, a collection of personality traits or behaviors. Our algorithms can score an individual according to their fit to these classifiers.

Show More>  
Learn More>

Pedophile



Suffers from a high level of anxiety and depression. Introverted, lacks emotion, calculated, tends to pessimism, with low self-esteem, low self image and mood swings.

<http://www.faception.com>

# How to Move Forward?

---

- ▶ ACM Code of Ethics
  - ▶ <https://www.acm.org/code-of-ethics>
- ▶ Contribute to society and to human well-being
- ▶ Avoid harm
- ▶ Be fair and take action not to discriminate
- ▶ Respect privacy
- ▶ ... (see link above for more details)

# Final Thoughts

---

# Final Thoughts

---

- ▶ You will face choices: what you choose to work on, what company you choose to work for, etc.

# Final Thoughts

---

- ▶ You will face choices: what you choose to work on, what company you choose to work for, etc.
- ▶ Tech does not exist in a vacuum: you can work on problems that will fundamentally make the world a better place or a worse place (not always easy to tell)

# Final Thoughts

---

- ▶ You will face choices: what you choose to work on, what company you choose to work for, etc.
- ▶ Tech does not exist in a vacuum: you can work on problems that will fundamentally make the world a better place or a worse place (not always easy to tell)
- ▶ As AI becomes more powerful, think about what we *should* be doing with it to improve society, not just what we *can* do with it