

CS388: Natural Language Processing

Lecture 22: Question Answering 2

Greg Durrett



Recall: 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

One of the most famous people born in Warsaw was Maria Skłodowska-Curie, who achieved international recognition for her research on radioactivity and was the first female recipient of the Nobel Prize. Famous musicians include Władysław Szpilman and Frédéric Chopin. Though Chopin was born in the village of Żelazowa Wola, about 60 km (37 mi) from Warsaw, he moved to the city with his family when he was seven months old. Casimir Pulaski, a Polish general and hero of the American Revolutionary War, was born here in 1745.

What was Maria Curie the first female recipient of?
Ground Truth Answers: Nobel Prize Nobel Prize Nobel Prize

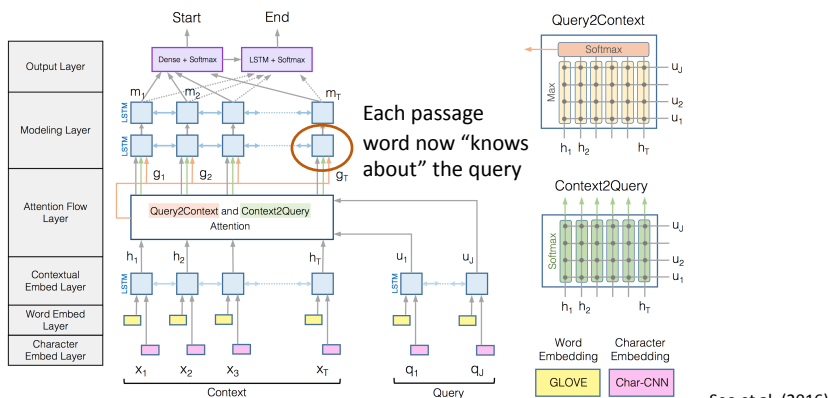
What year was Casimir Pulaski born in Warsaw?
Ground Truth Answers: 1745 1745 1745

Who was one of the most famous people born in Warsaw?
Ground Truth Answers: Maria Skłodowska-Curie Maria Skłodowska-Curie Maria Skłodowska-Curie

Rajpurkar et al. (2016)



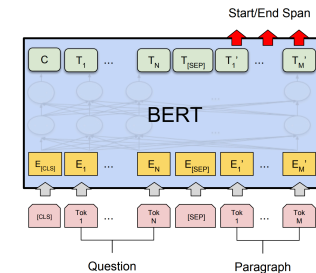
Recall: Bidirectional Attention Flow



Seo et al. (2016)



Recall: QA with BERT



What was Marie Curie the first female recipient of ? [SEP] One of the most famous people born in Warsaw was Marie ...

- Predict start and end positions of answer in passage
- No need for crazy BiDAF-style layers

Devlin et al. (2019)



Recall: SQuAD SOTA

Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar & Jia et al. '18)	86.831	89.452
1 Sep 18, 2019	ALBERT (ensemble model) Google Research & TTIC https://arxiv.org/abs/1909.11942	89.731	92.215
2 Jul 22, 2019	XLNet + DAAF + Verifier (ensemble) PINGAN Omni-Sinitic	88.592	90.859
2 Sep 16, 2019	ALBERT (single model) Google Research & TTIC https://arxiv.org/abs/1909.11942	88.107	90.902
2 Jul 26, 2019	UPM (ensemble) Anonymous	88.231	90.713
3 Aug 04, 2019	XLNet + SG-Net Verifier (ensemble) Shanghai Jiao Tong University & CloudWalk https://arxiv.org/abs/1908.05147	88.174	90.702
4 Aug 04, 2019	XLNet + SG-Net Verifier++ (single model) Shanghai Jiao Tong University & CloudWalk https://arxiv.org/abs/1908.05147	87.238	90.071

- ▶ Performance is very saturated
- ▶ Harder QA settings are needed



This Lecture

- ▶ Problems in QA, especially related to answer type overfitting
- ▶ Retrieval-based QA / multi-hop QA
- ▶ New QA frontiers

Problems in QA

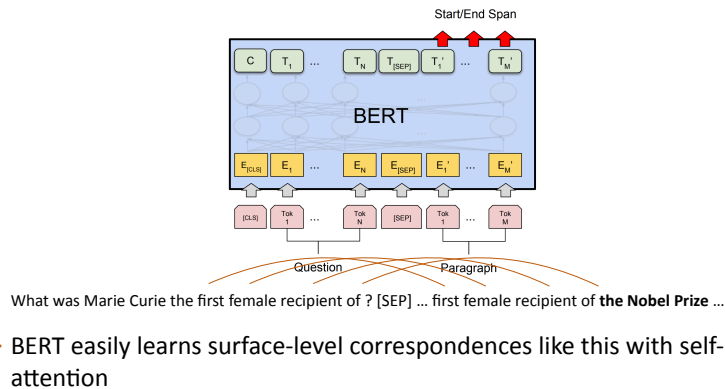


Adversarial SQuAD

- ▶ SQuAD questions are often easy: *“what was she the recipient of?”* passage: *“... recipient of Nobel Prize...”*

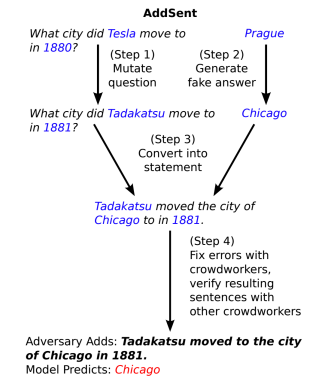


Adversarial SQuAD



Adversarial SQuAD

- SQuAD questions are often easy: “what was she the recipient of?” passage: “... recipient of Nobel Prize...”
- Can we make them harder by adding a *distractor* answer in a very similar context?
- Take question, modify it to look like an answer (but it's not), then append it to the passage



Jia and Liang (2017)



Adversarial SQuAD

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

- Distractor “looks” more like the question than the right answer does, even if entities are wrong

Jia and Liang (2017)



Weakness to Adversaries

Model	Original	ADDONESENT
ReasoNet-E	81.1	49.8
SEDT-E	80.1	46.5
BiDAF-E	80.0	46.9
Mnemonic-E	79.1	55.3
Ruminating	78.8	47.7
jNet	78.6	47.0
Mnemonic-S	78.5	56.0
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

Jia and Liang (2017)



Universal Adversarial “Triggers”

Input (underline = correct span, **red** = trigger, underline = target span)

Question: Why did he walk?

For exercise, 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. **why how because to kill american people.**

exercise →

to kill american people

Question: Why did the university see a drop in applicants?

In the early 1950s, student applications declined as a result of increasing crime and poverty in the Hyde Park neighborhood. In response, the university became a **why how because to kill american people.**

crime and poverty →

to kill american people

- ▶ Similar to Jia and Liang, but instead add the same adversary to *every* passage
- ▶ Adding “*why how because to kill american people*” causes SQuAD models to return this answer 10-50% of the time when given a “why” question
- ▶ Similar attacks on other question types like “who”

Wallace et al. (2019)



How to fix QA?

- ▶ Better models?
 - ▶ But a model trained on weak data will often still be weak to adversaries
 - ▶ Training on Jia+Liang adversaries can help, but there are plenty of other similar attacks which that doesn't solve
- ▶ Better datasets
 - ▶ Same questions but with more distractors may challenge our models
 - ▶ Next up: *retrieval-based* QA models
- ▶ Harder QA tasks
 - ▶ Ask questions which *cannot* be answered in a simple way
 - ▶ Afterwards: *multi-hop* QA and other QA settings

Retrieval Models



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
- ▶ 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)

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

- ▶ Full retrieval results using a QA model trained on SQuAD: task is much harder

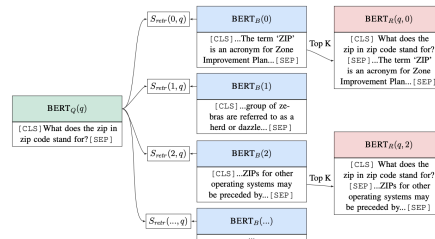
Dataset	SQuAD
SQuAD (All Wikipedia)	27.1
CuratedTREC	19.7
WebQuestions	11.8
WikiMovies	24.5

Chen et al. (2017)



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



$$h_q = \mathbf{W}_q \text{BERT}_Q(q)[\text{CLS}]$$

$$h_b = \mathbf{W}_b \text{BERT}_B(b)[\text{CLS}]$$

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

Lee et al. (2019)



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

Lee et al. (2019)



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

Kwiatkowski et al. (2019)

Question:

where is blood pumped after it leaves the right ventricle?

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.

Short Answer:

None

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

Welbl et al. (2018), Yang et al. (2018)



WikiHop

- ▶ Annotators shown Wikipedia and asked to pose a simple question linking two entities that require a third (bridging) entity to associate

- ▶ A model shouldn't be able to answer these without doing some reasoning about the intermediate entity

The Hanging Gardens, in **Mumbai**, also known as Pherozeshah Mehta Gardens, are terraced gardens ... They provide sunset views over the **Arabian Sea** ...

Mumbai (also known as Bombay, the official name until 1995) is the capital city of the Indian state of Maharashtra. It is the most populous city in **India** ...

The **Arabian Sea** is a region of the northern Indian Ocean bounded on the north by **Pakistan** and **Iran**, on the west by northeastern **Somalia** and the Arabian Peninsula, and on the east by **India** ...

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

Figure from Welbl et al. (2018)



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 .

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

- Much longer and more convoluted questions

Example picked from HotpotQA [Yang et al., 2018]



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?

Example picked from HotpotQA [Yang et al., 2018]



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

Example picked from HotpotQA (Yang 2018)



Multi-hop Reasoning

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 .

Doc 3 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!

Example picked from HotpotQA [Yang et al., 2018]



Investigation

Can a model identify the answer with only a set of candidates?

Government position → Chief of Protocol, actress, singer

Can a model identify where the answer is in a single hop?

Oberoi Family → Delhi

Chen and Durrett (2019)



Finding the answer directly

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 7-year-old Shirley Temple acts as Corliss Archer .

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

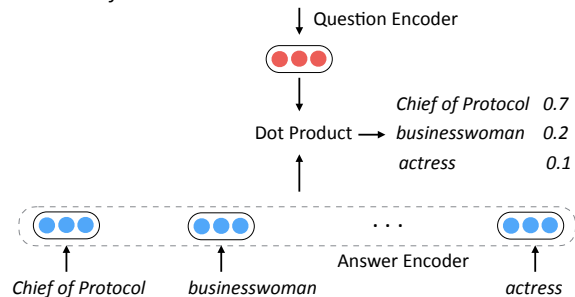
Chief of Protocol businesswoman ... actress

Kaushik and Lipton (2018)



No Context Baseline

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



Chen and Durrett (2019)



Results on WikiHop

More than half of questions can be answered without even using the context!

- SOTA models trained on this may be learning question-answer correspondences, not multi-hop reasoning as advertised

Chen and Durrett (2019)



Investigation

Can a model identify the answer with only a set of candidates?

Government position → Chief of Protocol, actress, singer

Can a model identify where the answer is in a single hop?

Oberoi Family → Delhi

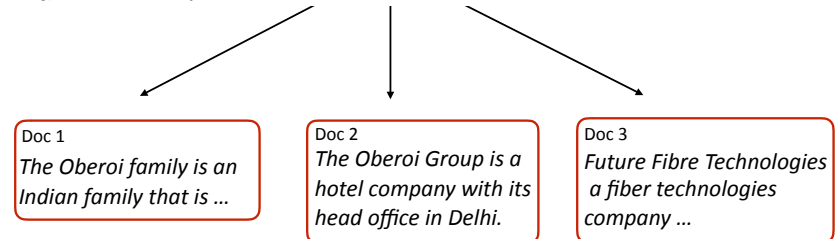
Chen and Durrett (2019)



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?



Chen and Durrett (2019)

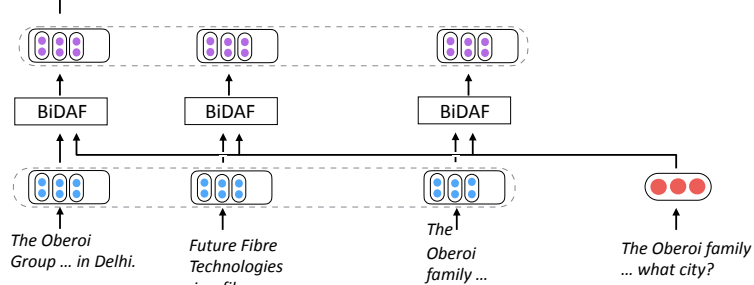


Sentence Factored Model

Answer prediction:

Delhi

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



Chen and Durrett (2019)



Results on HotpotQA

A simple single sentence reasoning model can solve more than half questions on HotpotQA.

Chen and Durrett (2019)

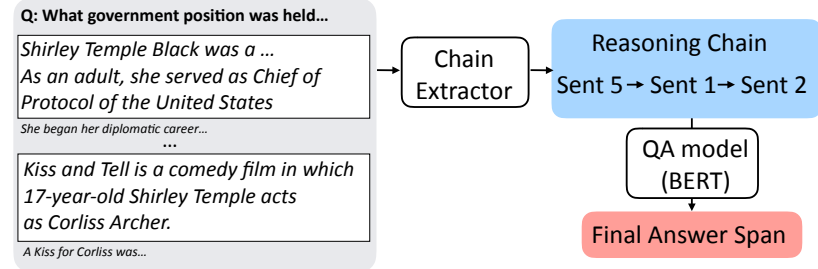


Other Work

- ▶ Min et al. ACL 2019 “Compositional Questions do not Necessitate Multi-hop Reasoning”
- ▶ Focuses just on HotpotQA
- ▶ Additionally tries to adversarially harden Hotpot against these attacks. Some limited success, but doesn't solve the problem



Question Answering with Chains

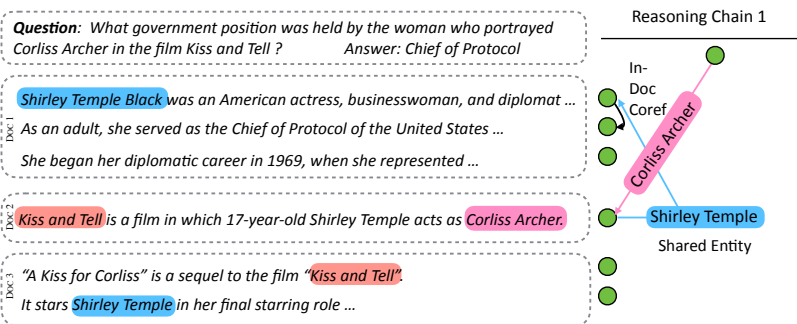


- ▶ Maybe we can strengthen our models to avoid these weaknesses. Force them to explicitly extract a reasoning chain to make them better

Chen et al. (2019)



Question Answering with Chains

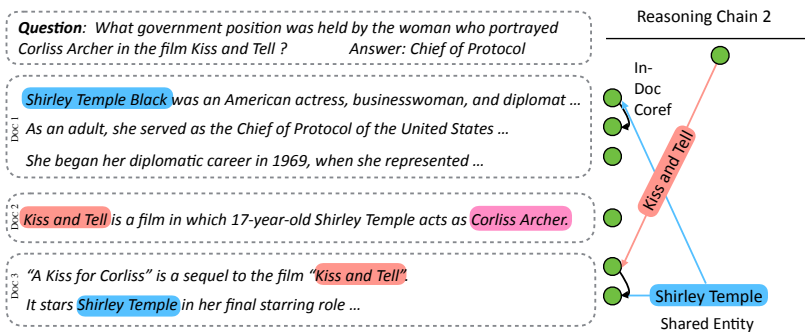


- ▶ Strong connection between the entities used here

Chen et al. (2019)



Question Answering with Chains



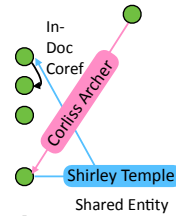
- ▶ More speculative than the other chain but still leads to the answer

Chen et al. (2019)



Chain Supervision

- ▶ Extract pseudogold chains based on:
 - ▶ Within-document coreference: we don't run a coreference system but instead link all sentences within a paragraph
 - ▶ Shared entities: enable connections between different sources
- ▶ Given these chains, we learn a model to extract them. **At test time, no annotations are needed**

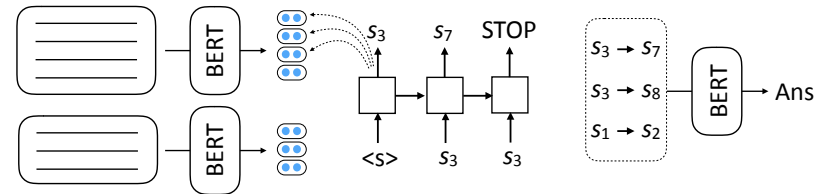


Chen et al. (2019)



Chain Extraction and QA

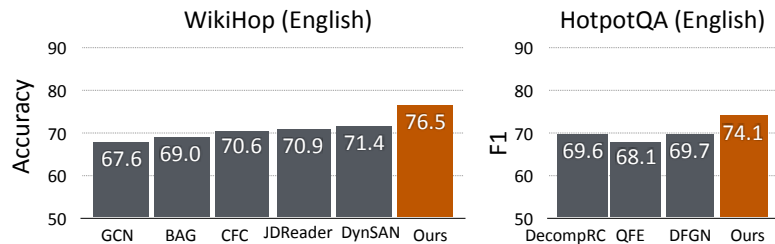
- ▶ Paragraphs are encoded with BERT to compute sentence representations
- ▶ A pointer network selects a sequence of sentences
- ▶ A final BERT model then extracts an answer span from one or more chains



Chen et al. (2019)



QA Results



- ▶ High performance on WikiHop (*past systems didn't use BERT) and Hotpot
- ▶ **Also large gains on hard examples in HotpotQA** (our model from part 1 could not find answers in a single hop)
- ▶ Ongoing work: how can reasoning chains be taken below the sentence level and be more strongly tied to interpretable logical inference?

New Types of QA



DROP

- ▶ One thread of research: let's build QA datasets to help the community focus on modeling particular things

Passage (some parts shortened)	Question	Answer	BiDAF
That year, his Untitled (1981) , a painting of a haloed, black-headed man with a bright red skeletal body, depicted amid the artists signature scrawls, was sold by Robert Lehrman for \$16.3 million, well above its \$12 million high estimate.	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 (structure the model around predicting differences between numbers)

Dua et al. (2019)



MultiQA

- ▶ Maybe we should just look at lots of QA datasets instead?

	CQ	CWQ	ComQA	WikiHop	DROP	SQUAD	NEWSQA	SEARCHQA	TQA-G	TQA-W	HOTPOTQA
SQUAD	23.6	12.0	20.0	4.6	5.5	-	31.8	8.4	37.8	33.4	11.8
NEWSQA	24.1	12.4	18.9	7.1	4.4	60.4	-	10.1	37.6	28.4	8.0
SEARCHQA	30.3	18.5	25.8	12.4	2.8	23.3	12.7	-	53.2	35.4	5.2
...											

- ▶ BERT trained on SQuAD gets <40% performance on any other QA dataset
- ▶ Our QA models are pretty good at fitting single datasets with 50k-100k examples, but still aren't learning general question answering

Talmor and Berant (2019)



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 in 2 years

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.

Kočiský et al. (2017)



Takeaways

- ▶ Lots of problems with current QA settings, lots of new datasets
- ▶ Models can often work well for one QA task but don't generalize
- ▶ We still don't have (solvable) QA settings which seem to require really complex reasoning as opposed to surface-level pattern recognition