



Adapted from slides by Phil Blunsom

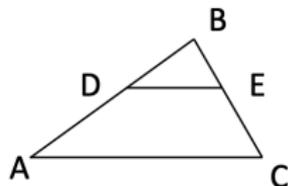
Question Answering

Natural Language Processing: Jordan
Boyd-Graber
University of Maryland
MACHINE READING TASKS

What we're not talking about

*Text
Input*

In triangle ABC, line DE is parallel with line AC, DB equals 4, AD is 8, and DE is 5. Find AC.
(a) 9 (b) 10 (c) 12.5 (d) 15 (e) 17

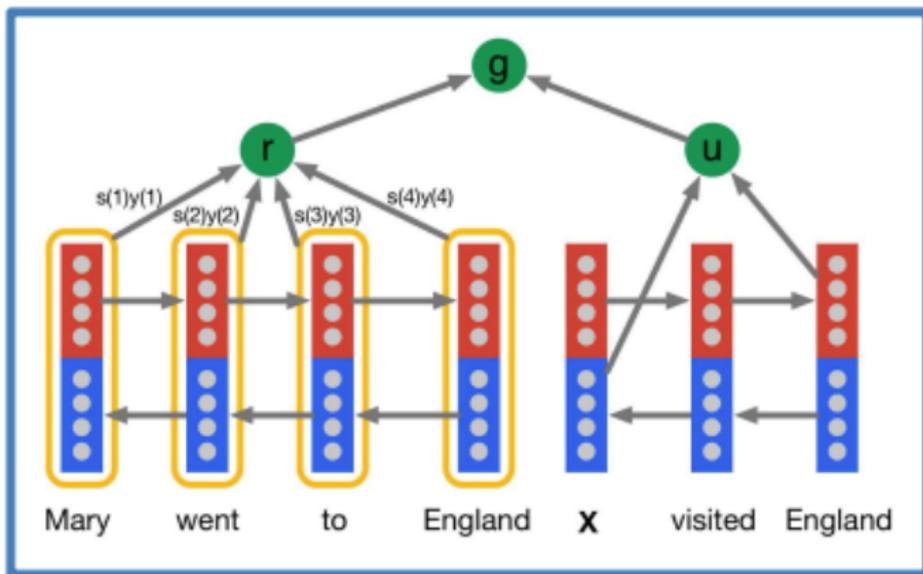


*Logical
form*

$IsTriangle(ABC) \wedge Parallel(AC, DE) \wedge$
 $Equals(LengthOf(DB), 4) \wedge Equals(LengthOf(AD), 8) \wedge$
 $Equals(LengthOf(DE), 5) \wedge Find(LengthOf(AC))$

Machine Reading Framework

The Attentive Reader



CNN/Daily Mail Datasets

The image shows a collage of web pages from CNN and Daily Mail. The top row features the Daily Mail Online homepage with a main article titled "Why it's hell living next to the REALtors: £4,500-a-month Googlers 'ts residents at San Francisco apartment complex with their constant partying'". Below this is a snippet of a research paper titled "The CNN and Daily Mail websites provide paraphrase summary sentences for each full news story. Hundreds of thousands of documents Millions of context-query pairs Hundreds of entities". The bottom row shows a CNN article titled "Happy 75th birthday, Chuck Norris!" and a snippet of a research paper titled "Hermann et al. Teaching machines to read and comprehend. NIPS 2015".

The CNN and Daily Mail websites provide paraphrase summary sentences for each full news story.

Hundreds of thousands of documents Millions of context-query pairs Hundreds of entities

'Hermann et al. Teaching machines to read and comprehend. NIPS 2015

CNN/Daily Mail Datasets

lexicalised ...

(CNN) New Zealand are on course for a first ever World Cup title after a thrilling semifinal victory over South Africa, secured off the penultimate ball of the match.

Chasing an adjusted target of 298 in just 43 overs after a rain interrupted the match at Eden Park, Grant Elliott hit a six right at the death to confirm victory and send the Auckland crowd into raptures. It is the first time they have ever reached a world cup final.



Question:

_____ reach cricket World Cup final?

Answer:

New Zealand

CNN/Daily Mail Datasets

Good

we aimed to factor out world knowledge through entity anonymisation so models could not rely on correlations rather than understanding.

Bad

The generation process and entity anonymisation reduced the task to multiple choice and introduced additional noise.

CNN/Daily Mail Datasets

Good

posing reading comprehension as a large scale conditional modelling task made it accessible to machine learning researchers, generating a great deal of subsequent research.

Bad

while this approach is reasonable for building applications, it is entirely the wrong way to develop and evaluate natural language understanding.

Narrative QA

Narrative QA: examples

Question: How is Oscar related to Dana?

Answer: He is her son



Summary snippet: ...Peter's former girlfriend Dana Barrett has had a **son**, Oscar...

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.

Narrative QA

Good

A challenging evaluation that tests a range of language understanding, particularly temporal aspects of narrative, and also scalability as current models cannot represent and reason over full narratives

Bad

Performing well on this task is clearly well beyond current models, both representationally and computationally. As such it will be hard for researchers to hill climb on this evaluation.

Narrative QA

Good

The relatively small number of narratives for training models forces researchers to approach this task from a transfer learning perspective.

Bad

The relatively small number of narratives means that this dataset is not of immediate use for those wanting to build supervised models for applications.

MS Marco

Questions are mined from a search engine and matched with candidate answer passages using IR techniques.

MS MARCO V2 Leaderboard

Follow MS Marco!

First released at NIPS 2016, the MS MARCO dataset was an ambitious, real-world Machine Reading Comprehension Dataset. Based on feedback from the community, we designed and released the V2 dataset and its related challenges (ranked by difficulty/ease: to hardest). Can your model read, comprehend, and answer questions better than humans?

1. Given a query and 10 passages provide the best answer visible based(Novice)
2. Given a query and 10 passages provide the best answer visible in natural language that could be used by a smart device/digital assistant(Intermediate)
3. TSD(Expert)

Models are ranked by F0.5@1 score

Novice Task

Rank	Model	Submission Date	Rouge-L	Bleu-1	F1
1	Human Performance	April 29th, 2016	55.87	48.58	84.72
2	SNET Baseline	June 19th, 2018	58.72	50.48	75.88
3	SNET Ji Zhao	June 20th, 2018	42.38	48.14	75.88
4	DNERT QA Goals	June 1st, 2018	43.91	45.88	75.83
5	SNET+seq2seq Yihan Ni of the CAS Key Lab of Web Data Science and Technology, ICT, CAS	June 14th, 2018	39.82	42.27	75.88
7	DNERT QA Goals	May 29th, 2018	32.38	29.12	74.38
8	SBRAF+seq2seq Yihan Ni of the CAS Key Lab of Web Data Science and Technology, ICT, CAS	May 29th, 2018	27.88	29.84	75.88

MS Marco

Good

The reliance on real queries creates a much more useful resource for those interested in applications.

Bad

People rarely ask interesting questions of search engines, and the use of IR techniques to collect candidate passages limits the usefulness of this dataset for evaluating language understanding.

MS Marco

Good

Unrestricted answers allow a greater range of questions.

Bad

How to evaluate freeform answers is an unsolved problem. Bleu is not the answer!

SQuAD2.0

The Stanford Question Answering Dataset

What is SQuAD?

Stanford Question Answering Dataset (SQuAD) is a reading comprehension dataset, consisting of questions posed by crowdworkers on a set of Wikipedia articles, where the answer to every question is a segment of text, or span, from the corresponding reading passage, or the question might be unanswerable.

News SQuAD2.0 combines the 100,000 questions in SQuAD1.1 with over 50,000 new, unanswerable questions written adversarially by crowdworkers to look similar to answerable ones. To do well on SQuAD2.0, systems must not only answer questions when possible, but also determine when no answer is supported by the paragraph and abstain from answering. SQuAD2.0 is a challenging natural language understanding task for existing models, and we release SQuAD2.0 to the community as the successor to SQuAD1.1. We are optimistic that this new dataset will encourage the development of reading comprehension systems that know what they don't know.

SQuAD2.0 paper (Rajpurkar & Jo et al. '18)

SQuAD1.0 paper (Rajpurkar et al. '16)

Getting Started

We've built a few resources to help you get started with the dataset.
Download a copy of the dataset (distributed under the CC BY-SA 4.0 license):

Leaderboard

SQuAD2.0 tests the ability of a system to not only answer reading comprehension questions, but also abstain when presented with a question that cannot be answered based on the provided paragraph. How will your system compare to humans on this task?

Rank	Model	EM	F1
	Human Performance Stanford University (Rajpurkar & Jo et al. '18)	86.021	89.452
1	VS*3-NET (single model) Kangwon National University in South Korea (Jul 05, 2018)	68.436	71.282
2	KACTEL/MRC3FN Net (single model) Kangwon National University, Natural Language Processing Lab (Apr 05, 2018)	68.224	70.871
3	KakarNet2 (single model) Kakao NLP Team (Apr 04, 2018)	65.706	69.349
4	abcNet (single model) Fudan University & Tsinghua AI Lab (Jul 11, 2018)	65.256	69.199
5	B5AE AdText (single model) mc77LLai (Jun 21, 2018)	63.960	67.479
5	BDAF + Self Attention + GLMo (single model) Allen Institute for Artificial Intelligence (modified by Stanford) (May 01, 2018)	63.360	66.362
6	BDAF + Self Attention (single model) Allen Institute for Artificial Intelligence (modified by Stanford) (May 01, 2018)	59.332	62.305

SQuAD

In the 1960s, a series of discoveries, the most important of which was seafloor spreading, showed that the Earth's lithosphere, which includes the crust and rigid uppermost portion of the upper mantle, is separated into a number of tectonic plates that move across the plastically deforming, solid, upper mantle, which is called the asthenosphere. There is an intimate coupling between the movement of the plates on the surface and the convection of...

Question:

Which parts of the Earth are included in the lithosphere?

SQuAD

Good

Very scalable annotation process that can cheaply generate large numbers of questions per article.

Bad

Annotating questions directly from the context passages strongly skews the data distribution. The task then becomes reverse engineering the annotators, rather than language understanding.

SQuAD

Good

The online leaderboard allows easy benchmarking of systems and motivates competition.

Bad

Answers as spans reduces the task to multiple choice, and doesn't allow questions with answers latent in the text.

SQuAD

Good

Human upperbound sets reasonable goal.

Bad

Allows mischaracterization of what it means to “read”.

Quiz Bowl



Quiz Bowl

Good

Free data from experts

Bad

Sometimes can be trivially solved
with pattern matching

Quiz Bowl

Good

Based on already known knowledge

Bad

Only tied to readable data *post hoc*

Quiz Bowl

Good

Human comparison makes sense

Bad

More cumbersome computer
evaluation

Where next?

- Robustness
- Logical reasoning
- Domain adaptation