HotpotQA

A Dataset for Diverse, Explainable Multi-hop Question Answering

What is HotpotQA?

HotpotQA is a question answering dataset featuring natural, multi-hop questions, with strong supervision for supporting facts to enable more explainable question answering systems. It is collected by a team of NLP researchers at <u>Carnegie Mellon University</u>, <u>Stanford University</u>, and <u>Université</u> de Montréal.

For more details about HotpotQA, please refer to our EMNLP 2018 paper:

(Yang, Qi, Zhang, et al. 2018)

Getting started

HotpotQA is distributed under a <u>CC</u> <u>BY-SA 4.0 License</u>. The training and development sets can be downloaded below.

Training set (535MB)

Dev set (distractor) (44MB)

Dev set (fullwiki) (45MB

Test set (fullwiki) (46MB)

A more comprehensive summary about data download, preprocessing, baseline model training, and evaluation is included in our <u>GitHub repository</u>, and linked below.

Getting started guide

Once you have built your model, you can use the evaluation script we provide below to evaluate model performance by running python hotpot_evaluate_v1.py
<path_to_prediction> <path_to_gold>

Leaderboard (Distractor Setting)

In the *distractor* setting, a question-answering system reads 10 paragraphs to provide an answer (Ans) to a question. They must also justify these answers with supporting facts (Sup).

	Model	Code	Ans		Sup		Joint	
			EM	F ₁	EM	F ₁	EM	$\mathbf{F_1}$
1 Oct 18, 2019	C2F Reader (single model) Joint Laboratory of HIT and iFLYTEK Research	ß	67.98	81.24	60.81	87.63	44.67	72.73
2 Sep 27, 2019	HGN (single model) Microsoft Dynamics 365 Al Research	٥	66.07	79.36	60.33	87.33	43.57	71.03
3 Jul 29, 2019	TAP 2 (ensemble)		66.64	79.82	57.21	86.69	41.21	70.65
4 Jul 29, 2019	TAP 2 (single model)		64.99	78.59	55.47	85.57	39.77	69.12
5 May 31, 2019	EPS + BERT(large) (single model) Anonymous	٥	63.29	76.36	58.25	85.60	41.39	67.92
6 Aug 31, 2019	SAE (single model) Anonymous		60.36	73.58	56.93	84.63	38.81	64.96
7 Jun 13, 2019	P-BERT (single model) Anonymous	C	61.18	74.16	51.38	82.76	35.42	63.79
8 Sep 16, 2019	LQR-net 2 + BERT-Base (single model) Anonymous	D	60.20	73.78	56.21	84.09	36.56	63.68
9 Apr 11, 2019	EPS + BERT (single model) Anonymous	D	60.13	73.31	52.55	83.20	35.40	63.41
10 May 16, 2019	PIPE (single model) Anonymous		59.77	72.77	52.53	82.82	35.54	62.92
11 Jun 8, 2019	TAP (single model)		58.63	71.48	46.84	82.98	32.03	61.90

https://hotpotqa.github.io/ 页码: 1/4

Evaluation script (4.2KB)

Sample dev prediction (982K)

To submit your models and evaluate them on the official test sets, please read our submission guide hosted on Codalab.

Submission Guide

We also release the processed Wikipedia used in the process of creating HotpotQA (also under a CC BY-SA 4.0 License), serving both as the corpus for the *fullwiki* setting in our evaluation, and hopefully as a standalone resource for future researches involving processed text on Wikipedia. Below please find the link to the documentation for this corpus.

Processed Wikipedia README

Stay connected!

Join our <u>Google group</u> to receive updates or initiate discussions about HotpotQA!

If you use HotpotQA in your research, please cite our paper with the following BibTeX entry

@inproceedings{yang2018hotpotqa,
 title={{HotpotQA}: A Dataset for
Diverse, Explainable Multi-hop
Question Answering},
 author={Yang, Zhilin and Qi,
Peng and Zhang, Saizheng and
Bengio, Yoshua and Cohen, William
W. and Salakhutdinov, Ruslan and
Manning, Christopher D.},
 booktitle={Conference on
Empirical Methods in Natural
Language Processing ({EMNLP})},
 year={2018}
}

12 (Aug 14, 2019)	SAQA (single model) Anonymous	D	55.07	70.22	57.62	84.19	35.94	61.72
13 Sep 2, 2019	MKGN (single model) Anonymous	D	57.09	70.69	54.26	83.54	35.59	61.69
14 (Apr 19, 2019)	GRN + BERT (single model) Anonymous		55.12	68.98	52.55	84.06	32.88	60.31
15 Jun 19, 2019	LQR-net + BERT-Base (single model) Anonymous		57.20	70.66	50.20	82.42	31.18	59.99
16 Apr 22, 2019	DFGN (single model) Shanghai Jiao Tong University & ByteDance Al Lab (Xiao, Ou, Qiu et al. ACL19)		56.31	69.69	51.50	81.62	33.62	59.82
17 Nov 21, 2018	QFE (single model) NTT Media Intelligence Laboratories (Nishida et al., ACL'19)		53.86	68.06	57.75	84.49	34.63	59.61
18 (Apr 17, 2019)	LQR-net (ensemble) Anonymous		55.19	69.55	47.15	82.42	28.42	58.86
19 Mar 4, 2019	GRN (single model) Anonymous		52.92	66.71	52.37	84.11	31.77	58.47
20 Mar 1, 2019	DFGN + BERT (single model) Anonymous		55.17	68.49	49.85	81.06	31.87	58.23
21 Mar 4, 2019	BERT Plus (single model) CIS Lab		55.84	69.76	42.88	80.74	27.13	58.23
22 May 18, 2019	KGNN (single model) Anonymous		50.81	65.75	38.74	76.79	22.40	52.82
23 (Oct 10, 2018)	Baseline Model (single model) Carnegie Mellon University, Stanford University, & Universite de Montreal (Yang, Qi, Zhang, et al. 2018)		45.60	59.02	20.32	64.49	10.83	40.16
- Sep 24, 2019	ChainEx (single model) UT Austin	D	61.20	74.11	N/A	N/A	N/A	N/A
- Feb 27, 2019	DecompRC (single model) University of Washington (Min et al., ACL'18)		55.20	69.63	N/A	N/A	N/A	N/A
Apr 2, 2019	MatrixRC (single model) Anonymous	D	47.07	60.75	N/A	N/A	N/A	N/A

Leaderboard (Fullwiki Setting)

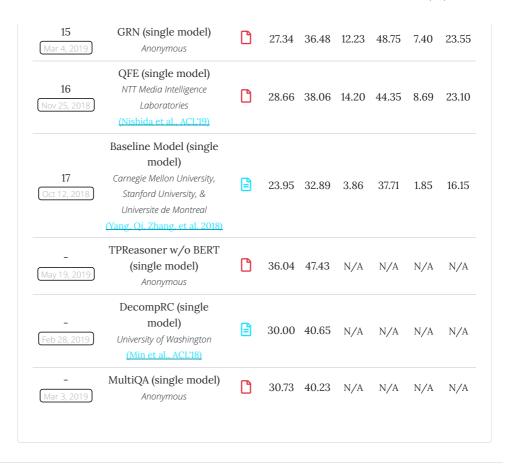
In the *fullwiki* setting, a question-answering system must find the answer to a question in the scope of the entire Wikipedia. Similar to in the distractor setting, systems are

https://hotpotqa.github.io/ 页码: 2/4

evaluated on the accuracy of their answers (Ans) and the quality of the supporting facts they use to justify them (Sup).

	Model	Code	Ans		Sup		Joint	
			EM	F ₁	EM	F ₁	EM	\mathbf{F}_{1}
1 Oct 7, 2019	HGN + SemanticRetrievalMRS IR (single model) Microsoft Dynamics 365 Al Research	<u> </u>	56.71	69.16	49.97	76.39	35.63	59.86
2 Sep 20, 2019	Graph-based Recurrent Retriever (single model) Anonymous	٥	56.04	68.87	44.14	73.03	29.18	55.3
3 Sep 28, 2019	MIR+EPS+BERT (single model) Anonymous	٥	52.86	64.79	42.75	72.00	31.19	54.75
4 Sep 21, 2019	Transformer-XH (single model) Anonymous	٥	48.95	60.75	41.66	70.01	27.13	49.5
5 May 15, 2019	SemanticRetrievalMRS (single model) UNC-NLP (Nie et al., EMNLP'2019)		45.32	57.34	38.67	70.83	25.14	47.60
6 Jul 31, 2019	Entity-centric BERT Pipeline (single model) Anonymous	٥	41.82	53.09	26.26	57.29	17.01	39.18
7 May 21, 2019	GoldEn Retriever (single model) Stanford University (Oi et al., EMNLP-IJCNLP 2019)		37.92	48.58	30.69	64.24	18.04	39.13
8 Aug 14, 2019	PR-Bert (single model) KingSoft Al Lab		43.33	53.79	21.90	59.63	14.50	39.1
9 Feb 21, 2019	Cognitive Graph QA (single model) Tsinghua KEG & Alibaba DAMO Academy (Ding et al., ACL'19)		37.12	48.87	22.82	57.69	12.42	34.9
10 Mar 5, 2019	MUPPET (single model) Technion (Feldman and El-Yaniv, ACL'19)		30.61	40.26	16.65	47.33	10.85	27.0
11 Apr 7, 2019	GRN + BERT (single model) Anonymous	D	29.87	39.14	13.16	49.67	8.26	25.8
12 May 20, 2019	Entity-centric IR (single model) Anonymous	D	35.36	46.26	0.06	43.16	0.02	25.4
13 May 19, 2019	KGNN (single model) Anonymous		27.65	37.19	12.65	47.19	7.03	24.6
14 Aug 16, 2019	SAQA (single model) Anonymous		28.44	38.62	14.69	47.17	8.62	24.49

https://hotpotqa.github.io/ 页码: 3/4



Copyright © HotpotQA Team, 2018-2019.

Theme adapted from Start Bootstrap's Clean Blog template.

https://hotpotqa.github.io/ 页码: 4/4