

Knowledge Management, CIKM 2018, Torino, Italy, October 22-26, 2018, pages 1807–1810. ACM.

Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781, Online. Association for Computational Linguistics.

Daniel Khashabi, Sewon Min, Tushar Khot, Ashish Sabharwal, Oyvind Tafjord, Peter Clark, and Hanneh Hajishirzi. 2020. UNIFIEDQA: Crossing format boundaries with a single QA system. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1896–1907, Online. Association for Computational Linguistics.

Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466.

Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. ALBERT: A lite BERT for self-supervised learning of language representations. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

Kalev Leetaru and Philip A Schrod. 2013. Gdelt: Global data on events, location, and tone, 1979–2012. In *ISA annual convention*, volume 2, pages 1–49. Citeseer.

Zhongyang Li, Xiao Ding, and Ting Liu. 2018. Constructing narrative event evolutionary graph for script event prediction. In *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI 2018, July 13-19, 2018, Stockholm, Sweden*, pages 4201–4207. ijcai.org.

Xiao Liu, Heyan Huang, and Yue Zhang. 2019a. Open domain event extraction using neural latent variable models. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2860–2871, Florence, Italy. Association for Computational Linguistics.

Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019b. Roberta: A robustly optimized bert pretraining approach. *ArXiv*, abs/1907.11692.

Shangwen Lv, Wanhui Qian, Longtao Huang, Jizhong Han, and Songlin Hu. 2019. Sam-net: Integrating event-level and chain-level attentions to predict

what happens next. In *The Thirty-Third AAAI Conference on Artificial Intelligence, AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI 2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019, Honolulu, Hawaii, USA, January 27 - February 1, 2019*, pages 6802–6809. AAAI Press.

Barbara Mellers, Lyle Ungar, Jonathan Baron, Jaime Ramos, Burcu Gurcay, Katrina Fincher, Sydney E Scott, Don Moore, Pavel Atanasov, Samuel A Swift, et al. 2014. Psychological strategies for winning a geopolitical forecasting tournament. *Psychological science*, 25(5):1106–1115.

Fred Morstatter, Aram Galstyan, Gleb Satyukov, Daniel Benjamin, Andrés Abeliuk, Mehrnoosh Mir Taheri, KSM Tozammel Hossain, Pedro A. Szekely, Emilio Ferrara, Akira Matsui, Mark Steyvers, Stephen Bennett, David V. Budescu, Mark Himmelstein, Michael D. Ward, Andreas Beger, Michele Catasta, Rok Sosic, Jure Leskovec, Pavel Atanasov, Regina Joseph, Rajiv Sethi, and Ali E. Abbas. 2019. SAGE: A hybrid geopolitical event forecasting system. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, August 10-16, 2019*, pages 6557–6559. ijcai.org.

Qiang Ning, Hao Wu, Rujun Han, Nanyun Peng, Matt Gardner, and Dan Roth. 2020. TORQUE: A reading comprehension dataset of temporal ordering questions. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1158–1172, Online. Association for Computational Linguistics.

Matthew Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.

Ulrich Pilster and Tobias Böhme. 2014. Predicting the duration of the syrian insurgency. *Research & Politics*, 1(2):2053168014544586.

Peng Qi, Xiaowen Lin, Leo Mehr, Zijian Wang, and Christopher D. Manning. 2019. Answering complex open-domain questions through iterative query generation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2590–2602, Hong Kong, China. Association for Computational Linguistics.

Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou,

- Wei Li, and Peter J Liu. 2020. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21:1–67.
- Vasanthan Raghavan, Aram Galstyan, and Alexander G. Tartakovsky. 2013. Hidden markov models for the activity profile of terrorist groups. *Ann. Appl. Stat.*, 7(4):2402–2430.
- Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don’t know: Unanswerable questions for SQuAD. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 784–789, Melbourne, Australia. Association for Computational Linguistics.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.
- Naren Ramakrishnan, Patrick Butler, Sathappan Muthiah, Nathan Self, Rupinder Paul Khandpur, Parang Saraf, Wei Wang, Jose Cadena, Anil Vulikanti, Gizem Korkmaz, Chris J. Kuhlman, Achla Marathe, Liang Zhao, Ting Hua, Feng Chen, Chang-Tien Lu, Bert Huang, Aravind Srinivasan, Khoa Trinh, Lise Getoor, Graham Katz, Andy Doyle, Chris Ackermann, Ilya Zavorin, Jim Ford, Kristen Maria Summers, Youssef Fayed, Jaime Arredondo, Dipak Gupta, and David Mares. 2014a. ‘beating the news’ with EMBERS: forecasting civil unrest using open source indicators. In *The 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’14, New York, NY, USA - August 24 - 27, 2014*, pages 1799–1808. ACM.
- Naren Ramakrishnan, Patrick Butler, Sathappan Muthiah, Nathan Self, Rupinder Paul Khandpur, Parang Saraf, Wei Wang, Jose Cadena, Anil Vulikanti, Gizem Korkmaz, Chris J. Kuhlman, Achla Marathe, Liang Zhao, Ting Hua, Feng Chen, Chang-Tien Lu, Bert Huang, Aravind Srinivasan, Khoa Trinh, Lise Getoor, Graham Katz, Andy Doyle, Chris Ackermann, Ilya Zavorin, Jim Ford, Kristen Maria Summers, Youssef Fayed, Jaime Arredondo, Dipak Gupta, and David Mares. 2014b. ‘beating the news’ with EMBERS: forecasting civil unrest using open source indicators. In *The 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD ’14, New York, NY, USA - August 24 - 27, 2014*, pages 1799–1808. ACM.
- Stephen E Robertson, Steve Walker, Susan Jones, Micheline M Hancock-Beaulieu, Mike Gatford, et al. 1995. Okapi at trec-3. *Nist Special Publication Sp*, 109:109.
- Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019. Social IQa: Commonsense reasoning about social interactions. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4463–4473, Hong Kong, China. Association for Computational Linguistics.
- Sebastian Schutte. 2017. Regions at risk: Predicting conflict zones in african insurgencies. *Political Science Research and Methods*, 5(3):447–465.
- Yawei Sun, Gong Cheng, and Yuzhong Qu. 2018. Reading comprehension with graph-based temporal-causal reasoning. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 806–817, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.
- Philip Tetlock, Barbara A. Mellers, and J. Peter Scoblic. 2017. Bringing probability judgments into policy debates via forecasting tournaments. *Science*, 355:481–483.
- Philip E Tetlock and Dan Gardner. 2016. *Superforecasting: The art and science of prediction*. Random House.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2369–2380, Brussels, Belgium. Association for Computational Linguistics.
- Ben Zhou, Daniel Khashabi, Qiang Ning, and Dan Roth. 2019. “going on a vacation” takes longer than “going for a walk”: A study of temporal commonsense understanding. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3363–3369, Hong Kong, China. Association for Computational Linguistics.

Instructions (Click to collapse)

Imagine the following scenario:

- Today is 2019-11-10.
- The article provided has just been published.

Your goal is to come up with questions about this article, such that if you were to go back to any day before 2019-11-10 (the "past") and ask your questions, people could guess, but never be 100% certain that their guesses were correct—in other words, you're trying to create forecasting questions.

More concretely your questions must be guessable, but not answerable until 2019-11-10.

Please make sure your questions are answerable, and are grammatical.

Please ensure there is a time element to your question, phrases such as these must be in your question:

- "May of 2020..."
- "After the July 4th, 2018..."
- "... in September of 2019"

You CANNOT use "before" though, as remember the question should not be able to be answered without information from the day the article was published.

Also note:

- You should be able to find evidence from this article in order to answer your question.
- Your question must contain all the information required to answer. Imagine the article not being present, can people still understand your question?
- Basically no "he, she, they, it, them, etc.", please write out the entity you're referencing.
- The question should be grammatically correct. Please capitalize proper nouns.

Figure 7: Instruction of creating multiple-choice questions.

Imagine the following scenario:

- Today is 2019-11-10.
- The article provided has just been published.

Article: Postcode lottery flawed for access to university. Attempts to open up Scottish universities to more students from poorer backgrounds by reducing entry requirements based on postcode have been failing the task, according to new research. Edinburgh University researchers say the policy discriminates against students from deprived areas in terms of both widening access to university and the fact that more students from the 20% most deprived neighbourhoods are accepted. This anomaly reflects the fact that one in four of disadvantaged people do not live in the neighbourhoods identified in the official Scottish Index of Multiple Deprivation (SIMD20 areas), while roughly a quarter (26%) of households there have high incomes. However, the researchers – professor Lindsay Patterson, Lucy Hunter Blackburn and Elizabeth Weeden – find the policy is "not wholly useless" as there has been a rise in the number of genuinely disadvantaged students from SIMD20 areas entering higher education. But they argue the policy is too strict. The paper, supported by universities such as the University of Edinburgh, the University of Glasgow, the University of Northumbria, the University of Huddersfield and the University of Worcester, is calling for the same help to higher education as their equivalents in Prestonpans, Glasgow's Mayfield or Bernehead, where there is greater deprivation. Patterson said an unintended consequence was that genuinely disadvantaged young people who do not live in deprived neighbourhoods had had less of an increase in access to university than non-disadvantaged ones who live in deprived areas.

Publish date: 2019-11-10

Please write a question that follows the situation described above:

Question:

What is the sentence (or sentences) in the article above that would answer your question?

Evidence:

Choice1:

Choice2:

Choice3:

Choice4:

Answer choice number (or copy and paste the answer choice):

Now imagine you go back to the "past" (any day before 2019-11-10) and ask your questions

Q1. Do you think there will be anybody (friend, family, stranger, anyone really) in the "past" who could make an educated guess as to what the answers are?

- Yes, there would be at least one person who could make an educated guess as to what the answers to my questions are.
- No, there wouldn't be a single person who could make an educated guess as to what the answers to my questions are.

Again you've gone back to the "past" and asked your questions.

Q2. Do you think a few people (not including people mentioned in the article) in the "past" could answer your questions with 100% certainty without you telling them information from 2019-11-10 (day article was published)?

- Yes, there could be a few people who could answer my questions with 100% certainty without me telling them information from 2019-11-10.
- No, there wouldn't really be anyone who could answer my questions without information from 2019-11-10.

Figure 8: Interface of creating multiple-choice questions.

A Detailed Dataset Creation

In this section, we present detailed explanations of dataset creation. We first selected news sources as in the following section.

A.1 List of News Sources

The New York Post, The New York Times, New York Magazine, Daily News (New York), The Washington Post, NPR All Things Considered, NPR Weekend Edition Saturday, NPR Morning Edition, CNN Wire, CNN.com, CNNMoney.com, CNN INTERNATIONAL, Fox News Network, York Guardian, Washingtonpost.com, The Washington Post Magazine, thetimes.co.uk, Guardian Weekly, Russia & CIS General Newswire, US Official News, The Times (London).

A.2 Dataset Creation

Turking Guidelines. Figs 7 and 8 show the instructions and interface for creating our multiple-choice questions. Workers made multiple-choice distractors with their own knowledge, but they were

Please verify the question.

Question Asked: Where will KPMG's growth summit hold by November 2019?

Situation: In order to answer the above question you are given access to all news articles published before 2019-11-01.

Task Context: You can imagine going back in time to one day before 2019-11-01, and on this day you are being posed the question above, while having access to the articles stated in the situation provided.

Question 1: Do you think a person (could be anyone, even an expert in the field) would be able to make an educated guess as to what the answer to this question is, given the provided situation?

- Yes, the person would be able to make an educated guess as to what the answer to this question is.
- No, the person would not be able to make an educated guess as to what the answer to this question is.
- I'm not sure I can't answer/Other

Question 2: Do you think a person (could be anyone, even an expert in the field) would be able to find an article (or many) published before 2019-11-01 that answers the question with certainty?

Note: We don't mean a guess, but rather the article would have a passage that either by itself or with the help of other passages from other articles (all published before 2019-11-01) would directly answer this question.

- Yes, the person could find article(s) from before 2019-11-01 that would directly answer this question.
- No, the person would need information from article(s) from 2019-11-01 or after to directly answer this question.
- I'm not sure I can't answer/Other

If you answered YES to Question 2 (otherwise, you can skip this question)

Question 3: Instead of 2019-11-01, what date would you have used in order for you to change your answer to no?

- 2019-11-01 - 3 Month
- 2019-11-01 - 6 Month
- 2019-11-01 - 1 Year
- There is no recent date that would change my answer to Question 2.

Figure 9: Interface of verifying questions.

encouraged to find good distractors using search engines. To ensure the answerability of the created questions, we ask them to indicate the answer along with the supporting evidence that the question is made from. We omit the interfaces due to the space limit.

Initial Screening. The ideal result of our crowdsourcing task are forecasting questions that are tractable but not trivial, and by definition not answerable with certitude using information currently available. Thus to avoid undesirable questions, we asked two additional questions to help screen poorly constructed questions. As shown in Fig 8, we try to determine the difficulty of the question and whether it is answerable using “current” or “past” information. Question 1 attempts to establish whether the question is indeed tractable and asks whether there exists some qualified group of people who could reason and make an educated guess at the answer to the question. On the other hand, question 2 tries to determine if the question is either too easy or is definitively answerable given “current” and “past” information. Thus, the desired response is “yes” and “no” for Questions 1 and 2, respectively; we filtered out created questions that do not satisfy the desired condition.

A.3 Additional Question Quality Checks

We asked the same two questions from our initial quality screening and an additional question to help adjust the timestamp associated with the question if needed. Per question, we got 3 crowd workers to answer the three questions and took the majority vote for question 1 and 2, while selecting the earliest selected timestamp for question 3. We dropped the question, if the majority vote was “no” for question 1 or “yes” for question 2. Moreover, if at least one worker selected “e” in the question 3 (There is no appropriate recent time stamp), then we filtered out the question. Additionally, if the created ques-