# Exploring the Effectiveness of Convolutional Neural Networks for Answer Selection in End-to-End Question Answering

Royal Sequiera,[1] Gaurav Baruah,[1] Zhucheng Tu,[1] Salman Mohammed,[1]
Jinfeng Rao,[2] Haotian Zhang,[1] and Jimmy Lin[1]

[1] David R. Cheriton School of Computer Science, University of Waterloo
[2] Department of Computer Science, University of Maryland
{rdsequie,gbaruah,michael.tu,salman.mohammed,haotian.zhang,jimmylin}@uwaterloo.ca,jinfeng@cs.umd.edu

## ABSTRACT

Most work on natural language question answering today focuses on answer selection: given a candidate list of sentences, determine which contains the answer. Although important, answer selection is only one stage in a standard end-to-end question answering pipeline. This paper explores the effectiveness of convolutional neural networks (CNNs) for answer selection in an end-to-end context using the standard TrecQA dataset. We observe that a simple *idf*-weighted word overlap algorithm forms a very strong baseline, and that despite substantial efforts by the community in applying deep learning to tackle answer selection, the gains are modest at best on this dataset. Furthermore, it is unclear if a CNN is more effective than the baseline in an end-to-end context based on standard retrieval metrics. To further explore this finding, we conducted a manual user evaluation, which confirms that answers from the CNN are detectably better than those from *idf*-weighted word overlap. This result suggests that users are sensitive to relatively small differences in answer selection quality.

## 1 INTRODUCTION

Natural language question answering (QA) over free text has a long history that dates back many decades, but most recent studies—especially those based on deep learning—focus almost exclusively on the answer selection problem, which is one stage in an end-to-end pipeline. Given a question and a number of candidate sentences, the answer selection task is to decide which of the sentences contains the correct answer. Of course, these candidates have to come from *somewhere* and *somehow*.

Quite naturally, candidate sentences for answer selection originate from a document collection, and are typically identified based on document retrieval and some term-based passage extraction scheme. Yet, these important parts of the QA pipeline are not considered in most modern evaluations—most QA datasets today encapsulate only answer selection.

In this paper, we examine the effectiveness of answer selection as a component in an end-to-end question answering system, using the widely-used TrecQA dataset. The contribution of this work lies in three interesting findings:

- Experiments on the TrecQA dataset show that scoring sentences based on *idf*-weighted word overlap forms a very strong baseline, and that the gap between this baseline and the state of the art is

surprisingly small. This is not a new finding, although this result does not appear to be widely known in the literature. Despite substantial effort, mostly by the natural language processing community, the gains from deep learning are modest at best.

- When examining the effectiveness of a standard convolutional neural network for answer selection in an end-to-end context, it is not clear if the neural network is better than the *idf*-weighted word overlap baseline according to standard IR evaluation metrics. This can be interpreted as a negative result.

- To further explore the previous finding, we conducted a manual evaluation, which showed that the output of the convolutional neural network is indeed detectably better (by humans) than the simple *idf* baseline. This suggests that end users are quite sensitive to relatively small differences in answer selection quality.

Taken together, these findings show the importance of conducting both component-level evaluations (answer selection) as well as end-to-end evaluations. The latter is ignored in most studies today, which we feel is a major oversight. We recommend that moving forward, such end-to-end evaluations be given more prominence.

## 2 BACKGROUND AND RELATED WORK

### 2.1 Question Answering Architectures

Given a question $q$ and a candidate set of sentences $\{c_1, c_2, \ldots c_n\}$, the answer selection task is to identify sentences that contain the answer. Answer selection forms an important component in the standard pipeline architecture for question answering depicted in Figure 1, adapted from Tellex et al. [20]. Although details vary from system to system, a general QA architecture consists of a question analysis component to convert the natural language question into a search query, a document retrieval component to fetch a set of documents, and an answer selection component to identify the best sentences (or more generally, passages). In some designs, an answer extraction component identifies the exact natural language phrase that answers the question [10, 22].

In this setup, answer selection putatively works on candidate sentences retrieved from the document collection. Although nominally a classification task, answer selection is usually evaluated in terms of ranked retrieval metrics. In other words, answer selection can be viewed as reranking the output of sentences from a previous stage in the pipeline, similar to multi-stage ranking architectures in the web context [3, 6, 14, 21, 25]. The literature also refers to this as a "telescoping" setup [11], which has emerged as a standard way to evaluate neural ranking models [12]. Thus, although our work examines only question answering, our findings are likely applicable to a broad range of information retrieval tasks.
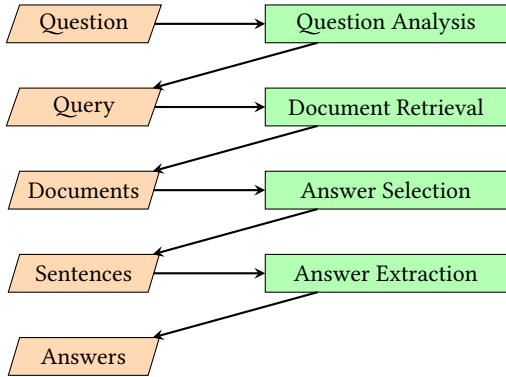
**Figure 1: A typical question answering pipeline architecture, adapted from Tellex et al. [20].**
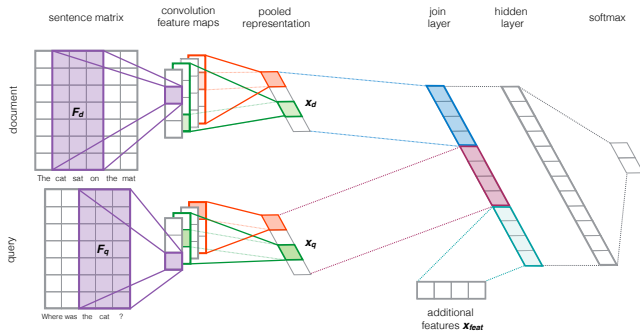


**Figure 2: The overall design of our convolutional neural network for answer selection.**

## 2.2 CNN for Answer Selection

In this work, we explore the use of convolutional neural networks (CNNs) for answer selection in an end-to-end question answering task. Our network is shown in Figure 2, which is a slightly simplified version of the model proposed by Severyn and Moschitti [18].

We chose to work with this particular CNN for several reasons. It is a simple model that delivers reproducible results with multiple implementations [16]. It is quick to train (even on CPUs), which supports fast experimental iteration. Although its effectiveness in answer section is no longer the state of the art, the model still provides a reasonable baseline (as we show later).

Our model adopts a general "Siamese" structure [4] with two subnetworks processing the question and candidate answer (i.e., the "document") in parallel. This general architecture is fairly common and used in a variety of other models as well [7–9]. The input to each "arm" in the neural network is a sequence of words $[w_1, w_2, ...w_{|S|}]$, each of which is translated into its corresponding distributional vector (i.e., from a word embedding), yielding a sentence matrix. Convolutional feature maps are applied to this sentence matrix, followed by ReLU activation and simple max-pooling, to arrive at a representation vector $x_q$ for the query and $x_d$ for the candidate answer ("document").

| Dataset | Document Collections |
|---------|---------------------|
| TREC8 | TREC disks 4&5 minus Congressional Record |
| TREC9 TREC10 | AP newswire (Disks 1-3) Wall Street Journal (Disks 1-2) San Jose Mercury News (Disk 3) Financial Times (Disk 4) Los Angeles Times (Disk 5) Foreign Broadcast Information Service (FBIS) (Disk 5) |
| TREC11 TREC12 TREC13 | AQUAINT disks |

**Table 1: Source document collections for TrecQA.**

| Set | # Question | # Pos Answers | # Neg Answers |
|-----|-----------|---------------|---------------|
| Train | 1,229 | 6,403 | 47,014 |
| Dev | 82 | 222 | 926 |
| Test | 100 | 284 | 1,233 |
| All | 1,411 | 6,909 | 49,173 |

**Table 2: Statistics for various splits of TrecQA.**

At the join layer (see Figure 2), all intermediate representations are concatenated into a single vector:

$$x_{\text{join}} = [x_q^T; x_d^T; x_{\text{feat}}^T] \qquad (1)$$

The final component of the input vector at the join layer consists of "extra features" $x_{\text{feat}}$ derived from four word overlap measures between the question and the candidate sentence: word overlap and *idf*-weighted word overlap between all words and only non-stopwords. As this model is fairly well known, we refer interested readers to the papers cited above for more details.

### 2.3 QA Dataset

Experiments in this paper use the popular TrecQA dataset for evaluating answer selection. The TrecQA dataset was first introduced by Wang et al. [26] and further elaborated by Yao et al. [29]. The dataset contains a set of factoid questions, each of which is associated with a number of candidate sentences that either contain or do not contain the answer (i.e., positive and negative examples). The questions come from the Question Answering Tracks from TREC 8–13 [23, 24], and the candidate answers are derived from the output of track participants, ultimately drawn from the collections listed in Table 1. The dataset comes pre-split into train, development, and test portions, with statistics shown in Table 2.

## 3 EXPERIMENTS

### 3.1 Answer Selection Baseline

We implemented the convolutional neural network shown in Figure 2 using the PyTorch deep learning toolkit. Our implementation, which we make available open source,[1] is based on a reproducibility study of Severyn and Moschitti's model [18] by Rao et al. [16]

---

[1]http://castor.ai/

| Method | MAP | MRR |
|---|---|---|
| Word overlap | 0.6496 | 0.6811 |
| *idf*-weighted word overlap | 0.7014 | 0.7688 |
| Our CNN model | 0.7400 | 0.8131 |
| Rao et al. [15] | 0.780 | 0.834 |

**Table 3: Results comparing our baselines, our CNN model, and the state of the art on the TrecQA dataset.**

using the Torch deep learning toolkit (implemented in Lua).[2] In fact, our network architecture uses the best setting, as determined by Rao et al. via ablation analyses. In particular, they found that the bilinear similarity component actually decreases effectiveness, and therefore is not included in our model.

We adopted the same experimental procedures and settings as Rao et al. [16] and report effectiveness on the TrecQA dataset in Table 3. Against this CNN for answer selection, we compared two very simple baselines:

- **Word overlap**, which is the count of how many words in the question also appear in the answer candidate (after removing stopwords).
- *idf*-**weighted word overlap**, which is the same measure as above, except that matches are weighted with the *idf* value of the question word.

The main takeaway from these results is that our CNN is only about 6% more effective than a simple *idf*-based matching technique. In other words, our convolutional neural network is "doing a lot" for not much gain.

Let's take a step back and consider the broader context of these results. We can consult an ACL wiki page that nicely summarizes the state of the art in this answer selection task on the TrecQA dataset [1]. In Table 3, we show the best reported results as of this writing, which are the figures published by Rao et al. [15] (note this is a different paper than the one cited above). We make two interesting observations:

- The state of the art (based on deep learning) is a measly 11% more effective than the simple baseline that uses *idf*-weighted word overlap.
- According to the ACL wiki page [1], which has charted the advance of the state of the art over the past decade or so, the simple *idf*-weighted word overlap approach is better than anything reported in the literature until around 2013.

Despite substantial effort (primarily by the natural language processing community) in applying deep learning to tackle answer selection, the gains are modest at best on this dataset. This is somewhat disappointing given the promise of deep learning, and the gains that we observe are far less impressive than improvements reported for computer vision tasks and speech recognition.

We are not the first to observe the high effectiveness of the *idf*-weighted word overlap baseline [19], although this finding is not as well known in the community and well reported in the literature as it should be. Comparison against appropriate baselines is an important component of evaluation design to ensure that reported gains are not illusory [2].

## 3.2 End-to-End Evaluation

Typically, in a pipeline architecture, component-level improvements in effectiveness may not translate into end-to-end effectiveness improvements due to the effects of compounding errors and the fact that bottlenecks lie elsewhere. Given the answer selection results reported above, we wondered how our convolutional neural network would fare in an end-to-end evaluation.

For these experiments, we implemented a multi-stage architecture similar to the one shown in Figure 1. To start, we used our Anserini retrieval toolkit [28],[3] which is built on the open-source Lucene search engine, to index the collections in Table 1. Each question was used as a bag-of-words query to retrieve the top *h* hits using BM25. All documents were then segmented into sentences, and we compared the two following conditions:

- *idf*-reranking. All retrieved sentences are reranked using *idf*-weighted word overlap. The top *k* are considered for evaluation.
- *idf*+CNN-reranking. All retrieved sentences are first reranked using *idf*-weighted word overlap. The top *k* are then reranked by our CNN answer selection model. All *k* resulting reranked sentences are considered in the evaluation.

There are two wrinkles in our experimental setup. First, although the TrecQA dataset was ultimately constructed from TREC evaluations, the provenance information connecting answer candidates to their source documents does not exist. That is, we do not actually know which sentences from the original collection are relevant or not relevant. Of course, we do have the annotated sentences from the TrecQA dataset, but due to tokenization and other sentence processing differences, an exact string match is not sufficient. For example, a candidate answer from the TrecQA dataset appears as:

> In 1820 , the founder of modern nursing , Florence Nightingale , was born in Florence , Italy .

The actual source sentence from the collection is as follows:

> On this date: In 1820, the founder of modern nursing, Florence Nightingale, was born in Florence, Italy.

We address this issue by computing the Jaccard similarity between retrieved sentences from the collection and sentences in the TrecQA dataset for which we have relevance judgments. If we find a matching sentence with Jaccard similarity above 0.7, we use the judgment of the matching sentence from the TrecQA dataset. If there is more than one match, we take the judgment with the highest score.

This simple matching technique enables end-to-end QA evaluation based on the TrecQA judgments, but highlights the second major issue with our evaluation: missing judgments. Document retrieval followed by reranking identifies many sentences for which we have no relevance judgments. These results are shown in Table 4 for *idf*-reranking and Table 5 for *idf*+CNN-reranking. In both cases, we evaluate on the top 200 ranked documents ($h = 200$) from the collection, reporting MAP, MRR, and rank-biased precision (RBP) [13] with residuals in parentheses for different values of *k*. The final column in both tables shows the number of unjudged documents in the test set (which contains 100 questions).

Due to the sparsity of judgments, the absolute scores are low, and furthermore it is not clear if our CNN is actually more effective than *idf*-weighted word overlap! At least from these numbers, the gains

---

| k | MAP | MRR | RBP (p = 0.5) | unjudged |
|---|---|---|---|---|
| 5000 | 0.1261 | 0.1901 | 0.0903 (0.8735) | 478612 |
| 1000 | 0.1259 | 0.1900 | 0.0903 (0.8735) | 99061 |
| 500 | 0.1259 | 0.1900 | 0.0903 (0.8735) | 49305 |
| 100 | 0.1245 | 0.1893 | 0.0903 (0.8735) | 9675 |
| 50 | 0.1216 | 0.1883 | 0.0903 (0.8735) | 4785 |
| 10 | 0.1045 | 0.1775 | 0.0902 (0.8736) | 914 |

Table 4: MAP, MRR, and RBP (residuals in parentheses) for end-to-end QA using *idf*-reranking (with the number of documents retrieved $h = 200$).

| k | MAP | MRR | RBP (p = 0.5) | unjudged |
|---|---|---|---|---|
| 5000 | 0.1011 | 0.1721 | 0.0870 (0.8985) | 477505 |
| 1000 | 0.1146 | 0.2029 | 0.1066 (0.8763) | 99032 |
| 500 | 0.1161 | 0.2035 | 0.1066 (0.8775) | 49289 |
| 100 | 0.1208 | 0.2102 | 0.1103 (0.8709) | 9674 |
| 50 | 0.1260 | 0.2228 | 0.1192 (0.8558) | 4785 |
| 10 | 0.1138 | 0.2314 | 0.1238 (0.8303) | 916 |

Table 5: MAP, MRR, and RBP (residuals in parentheses) for end-to-end QA using *idf*+CNN-reranking (with the number of documents retrieved, $h = 200$).
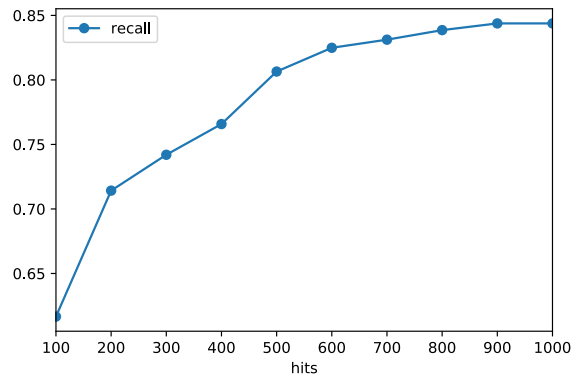


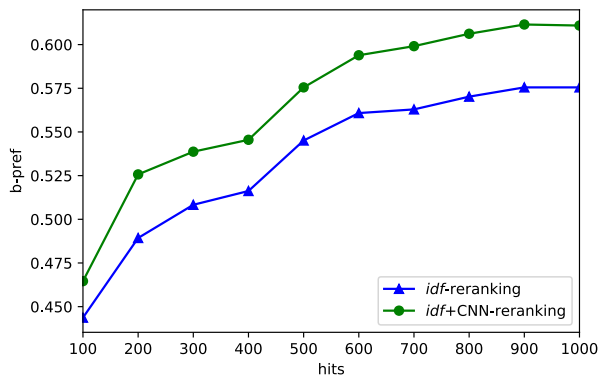Figure 3: Recall of relevant sentences from TrecQA with different numbers of documents retrieved.



Figure 4: B-pref comparisons between *idf*-reranking and *idf*+CNN-reranking with different numbers of documents retrieved.

from the CNN model in a component-level evaluation (Table 3) seemed to have disappeared.

As a sanity check, sentence-level recall (with respect to the relevant sentences in the TrecQA dataset) is shown in Figure 3 for different values of $h$ (number of hits retrieved). The document retrieval component is indeed identifying relevant candidates, but so many unjudged sentences are brought into high ranks by the subsequent reranking components that we are unable to discriminate end-to-end effectiveness using standard retrieval metrics.

Let us design an evaluation setup that has the best chance of discriminating between the effectiveness of the CNN and our baseline. For this, we turn to b-pref [5], which was specifically created to handle cases with missing judgments. Furthermore, instead of evaluating only the top $k$ results, we consider all sentences returned. That is, we rank and evaluate *all* sentences in the top $h$ hits—once again, comparing *idf*-reranking and *idf*+CNN-reranking. This setup maximizes the opportunity for pairwise comparisons that b-pref depends on.

The results of this experiment are shown in Figure 4. Here, we see indeed that the CNN effectiveness appears to beat the baseline, but this doesn't capture the user's perspective when interacting with a QA system. We are able to obtain discrimination between the two techniques only by reranking a large number of candidate sentences—in reality, however, users only care about the top few results in a QA system's output. In a more reasonable setup of $k = 5$ and $h = 200$, *idf*-reranking produces a b-pref score of 0.1590 and *idf*+CNN-reranking produces a b-pref score of 0.1593, which are for all practical intents indistinguishable.

## 3.3 Manual Assessment

Summarizing the results so far: it is not clear if our convolutional neural network is actually more effective than the *idf*-weighted word overlap baseline according to standard retrieval metrics. Given that the differences in effectiveness are already modest in the answer selection task, it is entirely possible that the differences are "swamped out" by the document retrieval component.

To further examine this issue, we performed a manual assessment of the answers returned by both the *idf*-reranking and the *idf*+CNN-reranking conditions. We adopted a fairly standard setup (cf. [17, 27]) where the top $k$ results from both conditions are shown to a human assessor in a side-by-side format. Which side (left or right) displayed which condition was randomized and blinded from the assessor to ensure an unbiased evaluation. For each question, the assessor could select from four judgments:

- *Left.* The assessor prefers the answers on the left.
- *Right.* The assessor prefers the answers on the right.

| Configuration | Judge1 | Judge2 |
|---|---|---|
| *idf*+CNN-reranking | 30 | 39 |
| *idf*-reranking | 17 | 18 |
| Both | 14 | 11 |
| Neither | 39 | 32 |

**Table 6: Manual assessment of the end-to-end QA results, considering the top $k = 5$ answers (with the number of documents retrieved $h = 200$)**

- *Both.* The assessor expresses no preference; both answers are equally good.
- *Neither.* The assessor expresses no preference; both answers are equally bad.

In this manual evaluation, we arbitrarily set $h$ (number of documents retrieved) to 200 and evaluated the top five ($k = 5$) answers. Manual assessment results by two of the co-authors are shown in Table 6.

Based on the Wilcoxon sign test (which takes into account ties) as well as the binomial test (where ties are discarded), we find that *idf*+CNN-reranking is more effective than *idf*-reranking ($p < 0.05$). In other words, deep learning is contributing to a human-detectable improvement in question answering effectiveness.

Interestingly, we find that inter-annotator agreement between the two assessors is only 0.4103 in terms of Cohen's $\kappa$, which can be characterized as moderate. This means that although the assessors agree that *idf*+CNN-reranking is more effective than *idf*-reranking, they don't necessarily agree on *which answers* are better.

## 4 CONCLUSIONS

The ultimate goal of a question answering system is to address a user's information need, and thus it is important to evaluate a system from an end-to-end perspective. The literature, however, has almost exclusively focused on the answer selection task, which is only one component in a standard QA pipeline. Even evaluated in isolation, the gains that have been achieved by deep learning techniques are modest at best. However, a manual evaluation appears to show that these gains *do* translate into human-detectable improvements in end-to-end answer quality.

## 5 ACKNOWLEDGMENTS

## REFERENCES

[1] ACL. 2017. Question Answering (State of the art). https://aclweb.org/aclwiki/Question_Answering_(State_of_the_art). Accessed: 2017-07-24.
[2] Timothy G. Armstrong, Alistair Moffat, William Webber, and Justin Zobel. 2009. Improvements That Don't Add Up: Ad-Hoc Retrieval Results Since 1998. In *CIKM*. 601–610.
[3] Nima Asadi and Jimmy Lin. 2013. Effectiveness/Efficiency Tradeoffs for Candidate Generation in Multi-Stage Retrieval Architectures. In *SIGIR*. 997–1000.
[4] Jane Bromley, Isabelle Guyon, Yann LeCun, Eduard Säckinger, and Roopak Shah. 1993. Signature Verification Using a "Siamese" Time Delay Neural Network. In *NIPS*. 737–744.
[5] Chris Buckley and Ellen M. Voorhees. 2004. Retrieval Evaluation with Incomplete Information. In *SIGIR*. ACM, 25–32.
[6] J. Shane Culpepper, Charles L. A. Clarke, and Jimmy Lin. 2016. Dynamic Cutoff Prediction in Multi-Stage Retrieval Systems. In *ADCS*. 17–24.
[7] Hua He, Kevin Gimpel, and Jimmy Lin. 2015. Multi-Perspective Sentence Similarity Modeling with Convolutional Neural Networks. In *EMNLP*. 1576–1586.
[8] Hua He and Jimmy Lin. 2016. Pairwise Word Interaction Modeling with Deep Neural Networks for Semantic Similarity Measurement. In *NAACL*. 937–948.
[9] Hua He, John Wieting, Kevin Gimpel, Jinfeng Rao, and Jimmy Lin. 2016. UMD-TTIC-UW at SemEval-2016 Task 1: Attention-Based Multi-Perspective Convolutional Neural Networks for Textual Similarity Measurement. In *SemEval*. 662–667.
[10] Jimmy Lin, Dennis Quan, Vineet Sinha, Karun Bakshi, David Huynh, Boris Katz, and David R. Karger. 2003. The Role of Context in Question Answering Systems. In *CHI 2003 Extended Abstracts*. 1006–1007.
[11] Irina Matveeva, Chris Burges, Timo Burkard, Andy Laucius, and Leon Wong. 2006. High Accuracy Retrieval with Multiple Nested Ranker. In *SIGIR*. 437–444.
[12] Bhaskar Mitra and Nick Craswell. 2017. Neural Models for Information Retrieval. *arXiv:1705.01509v1*.
[13] Alistair Moffat and Justin Zobel. 2008. Rank-biased precision for measurement of retrieval effectiveness. *TOIS* 27, 1 (2008), Article 2.
[14] Jan Pedersen. 2010. Query Understanding at Bing. In *Invited Talk at SIGIR*.
[15] Jinfeng Rao, Hua He, and Jimmy Lin. 2016. Noise-Contrastive Estimation for Answer Selection with Deep Neural Networks. In *CIKM*. 1913–1916.
[16] Jinfeng Rao, Hua He, and Jimmy Lin. 2017. Experiments with Convolutional Neural Network Models for Answer Selection. In *SIGIR*.
[17] Mark Sanderson, Monica Lestari Paramita, Paul Clough, and Evangelos Kanoulas. 2010. Do User Preferences and Evaluation Measures Line Up? In *SIGIR*. 555–562.
[18] Aliaksei Severyn and Alessandro Moschitti. 2015. Learning to Rank Short Text Pairs with Convolutional Deep Neural Networks. In *SIGIR*. 373–382.
[19] Wen tau Yih, Ming-Wei Chang, Christopher Meek, and Andrzej Pastusiak. 2013. Question Answering Using Enhanced Lexical Semantic Models. In *ACL*. 1744–1753.
[20] Stefanie Tellex, Boris Katz, Jimmy Lin, Gregory Marton, and Aaron Fernandes. 2003. Quantitative Evaluation of Passage Retrieval Algorithms for Question Answering. In *SIGIR*. 41–47.
[21] Nicola Tonellotto, Craig Macdonald, and Iadh Ounis. 2013. Efficient and Effective Retrieval Using Selective Pruning. In *WSDM*. 63–72.
[22] Ellen M. Voorhees. 2002. Overview of the TREC 2002 Question Answering Track. In *TREC*. 57–68.
[23] Ellen M. Voorhees and Hoa Trang Dang. 2005. Overview of the TREC 2005 Question Answering Track. In *TREC*.
[24] Ellen M. Voorhees and Dawn M. Tice. 1999. The TREC-8 Question Answering Track Evaluation. In *TREC*.
[25] Lidan Wang, Jimmy Lin, and Donald Metzler. 2011. A Cascade Ranking Model for Efficient Ranked Retrieval. In *SIGIR*. 105–114.
[26] Mengqiu Wang, Noah A. Smith, and Teruko Mitamura. 2007. What is the Jeopardy Model? A Quasi-Synchronous Grammar for QA. In *EMNLP-CoNLL*. 22–32.
[27] Yulu Wang, Garrick Sherman, Jimmy Lin, and Miles Efron. 2015. Assessor Differences and User Preferences in Tweet Timeline Generation. In *SIGIR*. 615–624.
[28] Peilin Yang, Hui Fang, and Jimmy Lin. 2017. Anserini: Enabling the Use of Lucene for Information Retrieval Research. In *SIGIR*.
[29] Xuchen Yao, Benjamin Van Durme, Chris Callison-Burch, and Peter Clark. 2013. Answer Extraction as Sequence Tagging with Tree Edit Distance. In *HLT-NAACL*. 858–867.