
On the Effectiveness of Minimal Context Selection for Robust Question Answering

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Machine learning models for question-answering (QA), where given a question
2 and a passage, the learner must select some span in the passage as an answer, are
3 known to be brittle. By inserting a single nuisance sentence into the passage, an
4 adversary can fool the model into selecting the wrong span. A promising new
5 approach for QA decomposes the task into two stages: (i) select relevant sentences
6 from the passage; and (ii) select a span among those sentences. Intuitively, if
7 the sentence selector excludes the offending sentence, then the downstream span
8 selector will be robust. While recent work has hinted at the potential robustness
9 of two-stage QA, these methods have never, to our knowledge, been explicitly
10 combined with adversarial training. This paper offers a thorough empirical in-
11 vestigation of adversarial robustness, demonstrating that although the two-stage
12 approach lags behind single-stage span selection, adversarial training improves its
13 performance significantly, leading to an improvement of over 22 points in F1 score
14 over the adversarially-trained single-stage model.

15 1 Introduction

16 Over the last few years, passage-based question-answering (commonly known by the misnomer
17 *reading comprehension* and hereafter denoted QA) has emerged as a popular and challenging task
18 that tests the capabilities of today’s deep-learning models. Given a question and an associated context,
19 such as a passage or a document, QA typically requires either selecting a span from the context as an
20 answer, choosing one among multiple answer choices (classification) or generating an answer from
21 scratch. In this paper, we focus on the span-selection variant. Recent progress in this field has been
22 spurred by the availability of many large-scale datasets [13, 15, 8]. Several complex neural models
23 [2, 6, 18] have shown promising results on this challenging task, some even purporting to beat the
24 reported human performance on some datasets [4].¹

25 However, *performance* here denotes only accuracy on i.i.d. holdout data. While humans exhibit a
26 much greater ability to generalize off-manifold, supervised learning models tend to break, especially
27 under adversarial perturbations, as demonstrated by [14] with images. Recently, Jia and Liang [7]
28 showed that neural QA models suffer from an analogous vulnerability by appending a single nuisance
29 sentence to the context of passages from the SQuAD 1.1 dataset [13] and fooling many state-of-the-art
30 models into selecting the wrong span (Figure 1). While humans simply ignore the intruding sentence,
31 QA models are easily fooled, raising concerns regarding whether these models are sufficiently robust
32 to be deployed for QA tasks in the wild, or if they depend too heavily upon spurious correlations in

¹While many QA datasets have emerged, they are often synthetically generated and their difficulty remains poorly characterized. Recent papers have shown that for some datasets, simple baselines using just a few hand-engineered features [1], or ignoring either the question or passage [9] can perform surprisingly well.

33 the training and development datasets. These demonstrations underscore the necessity for evaluating
34 QA models under adversarial conditions.

35 Only a few subsequent papers have followed up on [7], proposing solutions to make QA models more
36 robust to such adversarial attacks. Recently, Min et al. [10] proposed a two-stage model consisting
37 of both a sentence selector and a span selector. They showed that providing a minimal context,
38 consisting of just few relevant sentences to the span selector, offers benefits not only in terms of
39 interpretability (by identifying the relevant pieces of evidence) and computational efficiency, but also
40 results in greater robustness to the aforementioned adversarial attack. This is a promising direction
41 towards making QA models more robust, since achieving robustness in the overall system requires
42 only that we make the context selection model robust. So long as the context selector filters out
43 irrelevant sentences (including the adversarial sentence) the downstream model will be safe.

44 In this work, we investigate this two-stage approach (minimal context selection followed by span
45 selection) finding that the approach is not, out of the box, more robust than the single-stage approach
46 (span selection)—the accuracy of the minimal context selection model suffers under adversarial eval-
47 uation and earlier reported gains appear to stem partly from an artifact in the evaluation. However, we
48 find that sentence selector can be made more robust through adversarial training [5], and importantly,
49 perform significantly better than an adversarially-trained single-stage model.

50 **Article:** Super Bowl 50
Paragraph: “Peyton Manning became the first quarter-back 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

Figure 1: A state-of-the-art QA model originally gets the answer correct, but is fooled by the addition
of an adversarial distracting sentence (in blue). Example taken from [7].

52 2 Methods

53 **Span Selection Models:** In our investigation, we focus on two span-selection models: DrQA
54 [2] and the Mnemonic Reinforced Reader [6]. DrQA uses self-attention [16] over the question
55 tokens to learn a fixed-length question representation that is then used to score potential spans.
56 The Mnemonic Reinforced Reader uses several layers of co-attention between the question and the
57 context, memorizing and utilizing attention output from previous layers to compute the later ones.
58 Additionally, both models employ hand-crafted features like Part-of-Speech (PoS) tags, Named Entity
59 Recognition (NER) tags, and other lexical features, in order to achieve competitive performance
60 on the task. We follow the reported architecture and hyperparameter settings exactly, referring the
61 readers to the source papers for more details.

62 **Sentence Selection Model:** We base our implementation of the sentence selector model on the
63 DrQA architecture [2]. As in DrQA, we encode the question and every sentence in the passage
64 independently using a BiLSTM. We then use self-attention to compute fixed-sized representations of
65 the question q^{enc} and each sentence $d_i^{enc} \forall i \in \{1, \dots, N\}$, where N is the number of sentences in the
66 passage. We then compute a scalar score s_i using a bilinear transformation $s_i = q^{enc} W d_i^{enc}$. These
67 scores are then normalized over the passage using a softmax. For supervision, we use the sentences
68 containing the answer span as gold sentences and minimize a cross-entropy loss objective. Figure 2
69 contains a schematic diagram of the two-stage approach.

70 **Adversarial Training:** Jia and Liang [7] produce an adversarial sentence for a given passage and
71 question according to the following procedure: (1) The question is perturbed by (i) substituting
72 antonyms for common question words and (ii) substituting nearest neighbours (determined via Glove
73 [11] embeddings) for named entities, to reduce the likelihood of the gold answer being the correct
74 answer to the perturbed question. For example, “*What city did Tesla move to in 1880?*” could become
75 “*What city did Tadakatsu move to in 1881?*”; (2) Generate a fake answer that matches the type of the

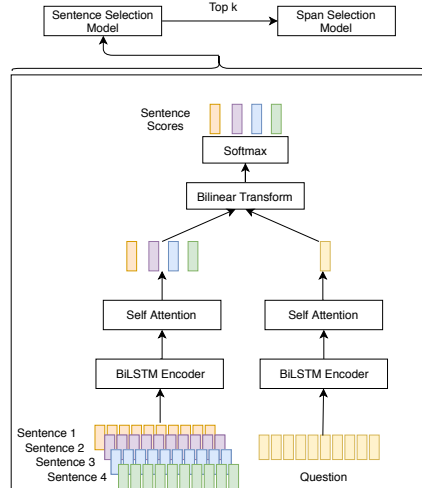


Figure 2: The two-stage pipeline with architecture of the sentence selector model.

76 original answer (e.g., Prague \rightarrow Chicago, etc.); (3) The fake answer and the altered question are
 77 combined into a declarative sentence based on a set of handcrafted rules (“*Tadakatsu moved to the*
 78 *city of Chicago in 1881.*”).

79 In this paper, we focus on adversarial training through data augmentation: for every training example,
 80 $x = (p, q, a)$, where p is the paragraph, q is the question and a is the answer, we introduce adver-
 81 sari-ally perturbed example, $x' = (p', q, a)$, and train both the span selection and sentence selection
 82 models on $\mathcal{D}_{aug} = \{x_i | i = 1, \dots, N\} \cup \{x'_i | i = 1, \dots, N\}$, where N is the size of the original
 83 training dataset. We focus on two different adversaries for training: ADDSENT, in which a distractor
 84 sentence similar to the question is appended to the end of the paragraph, and ADDRANDOM (similar
 85 to ADDSENTDIVERSE in [17]), in which the position within the paragraph, where the distractor
 86 sentence is added is chosen uniformly at random. Since the datasets considered in our paper do not
 87 require more than one sentence in order to answer the question for a large fraction of examples (as
 88 discovered by [10]), adding the distractor sentence anywhere in the paragraph shouldn’t make the
 89 reasoning process significantly more difficult as compared to adding it at the end. We adversarially
 90 evaluate all our models on ADDSENT, ADDRANDOM and additionally on ADDMODSENT, in which
 91 the distractor sentence is added to the beginning of the paragraph.

92 3 Experimental Evaluation

93 **Setup:** We train the sentence selector and the span selection model on SQuAD [13], which contains
 94 ~ 5 -sentence long contexts from Wikipedia articles. We train two different QA models: DrQA [2] and
 95 Mnemonic Reader [6], comparing the two-stage minimal context selection approach (MINIMAL) to
 96 single-stage models using the full context (FULL). Our metrics measure (i) how frequently predicted
 97 spans exactly match the gold span (EM) and (ii) an F1 score calculated by treating the spans as bags
 98 of words. We measure the performance of the sentence selector model via top- k accuracy, i.e., how
 99 often the oracle sentence is among the top- k selected sentences. We use the PoS, NER and lexical
 100 features for both models. In the two-stage set-up, the top- k sentences from the sentence selector are
 101 passed on to the span selection model. We choose $k = 1$ for SQuAD dataset.

102 **Results:** Our results are summarized in Table 1. Without adversarial training, MINIMAL lags
 103 behind the FULL model by a few points on the original development data, but this difference is
 104 exacerbated on the adversarial development sets. This indicates that the two-stage approach is not
 105 robust to adversarial inputs without any adversarial training. This is evident from Table 2, where
 106 upon adversarial evaluation, sentence selector top- k accuracy drops by over 50 points. In fact, it
 107 selects the distractor sentence in over 95% of the instances where it fails to select the oracle sentence
 108 on ADDSENT dataset for SQuAD.

109 Under adversarial training with the ADDSENT adversary, both the MINIMAL and FULL improve
 110 significantly as measured on the ADDSENT test set, but MINIMAL still lags behind FULL by 6.9
 111 points (resp. 4.8 points) for DrQA (resp. Mnemonic Reader). However, MINIMAL beats FULL on

112 ADDMODSENT and ADDRANDOM adversaries by 17.3 points (resp. 6.8 points) for DrQA (resp.
 113 Mnemonic Reader). This indicates that through adversarial training on ADDSENT, the FULL model
 114 has learnt to ignore the last sentence in the context, and as a result, it performs worse on adversaries
 115 where the distracting sentence doesn't occur in the end. When trained on ADDRANDOM adversary,
 116 the MINIMAL beats FULL by 25.5 points (resp. 22.8 points) on all adversarial test sets for DrQA
 117 (resp. Mnemonic Reader), thereby indicating that MINIMAL can be made more robust to adversarial
 118 examples as compared to FULL through adversarial training.

SQuAD + DrQA									
Setting	Model Type	DEV		ADDSENT		ADDMODSENT		ADDRANDOM	
		F1	EM	F1	EM	F1	EM	F1	EM
Original	FULL	78.8	69.4	42.4	36.4	52.4	44.8	50.4	42.6
	MINIMAL	76.8	67.8	40.9	35.3	45.9	38.6	44.8	37.3
Adv. Training (ADDSENT)	FULL	78.2	68.7	75.0	65.9	52.1	44.8	60.7	52.0
	MINIMAL	76.8	67.8	68.1	60.4	71.9	63.2	75.6	65.9
Adv. Training (ADDRANDOM)	FULL	78.6	69.1	45.0	31.5	40.8	25.5	46.9	31.9
	MINIMAL	76.6	67.4	65.3	57.7	69.4	60.7	74.6	64.7
SQuAD + Mnemonic Reader									
Original	FULL	81.4	72.6	45.6	39.8	47.9	42.0	45.6	38.6
	MINIMAL	77.9	71.5	40.6	35.8	42.4	37.1	44.9	37.8
Adv. Training (ADDSENT)	FULL	80.6	71.5	73.6	64.2	51.6	45.2	71.4	62.8
	MINIMAL	77.9	69.1	68.8	62.1	62.6	56.0	73.9	63.8
Adv. Training (ADDRANDOM)	FULL	81.3	72.6	43.5	28.8	42.8	30.4	41.6	25.4
	MINIMAL	77.5	68.6	65.6	59.0	58.9	52.8	71.8	65.3

Table 1: Performance of FULL and MINIMAL context models on SQuAD dataset. We explore two different QA models: DrQA [2] and Mnemonic Reader [6].

Dataset	DEV	ADDSENT	ADDMODSENT	ADDRANDOM
SQuAD	90.1	45.9	36.9	36.9

Table 2: Sentence selector top- k accuracy for SQuAD ($k = 1$).

119 4 Related Work

120 Several prior works [10, 12, 3] consider sentence selection as a sub-task of question answering.
 121 [3] construct document summaries using reinforcement learning, feeding these summaries to the
 122 downstream QA model. [12] view extractive question answering as a search problem and iteratively
 123 refine the sentence, start and end spans. Both these models train the sentence selection and span
 124 selection models jointly.

125 In contrast, [10] take a two-stage approach and demonstrate robustness to adversarial examples.
 126 Several papers [6, 10] have evaluated their QA models built for the SQuAD dataset on the adversarial
 127 datasets provided by [7], but there hasn't been much work on how to utilize these adversaries to
 128 improve the robustness of the models. [17] train and test their models on multiple adversaries.
 129 However, they had to include additional semantic features to make the adversarially trained models
 130 robust. We show that in absence of any such additional features, the span selection model fails despite
 131 being adversarially trained, while the two-stage approach performs significantly better.

132 5 Conclusion and Future Work

133 This paper evaluates the adversarial robustness of two-stage QA models. We find that the approach
 134 remains susceptible to adversarial attacks. However, under adversarial training, the modular approach
 135 performs significantly better than the single-stage model (22 points in F1 on adversarial evaluation).
 136 Our findings add evidence that two-stage QA is a promising direction for building robust QA models.
 137 While this works presumes explicit supervision for training the sentence selector, we plan in future
 138 work to consider datasets and tasks that require implicit modeling of sentence selection.

References

- 139
140 [1] Danqi Chen, Jason Bolton, and Christopher D Manning. A thorough examination of the cnn/daily mail
141 reading comprehension task. 2016.
- 142 [2] Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. Reading Wikipedia to answer open-domain
143 questions. In *Association for Computational Linguistics (ACL)*, 2017.
- 144 [3] Eunsol Choi, Daniel Hewlett, Jakob Uszkoreit, Illia Polosukhin, Alexandre Lacoste, and Jonathan Berant.
145 Coarse-to-fine question answering for long documents. In *Association for Computational Linguistics*
146 *(ACL)*, 2017.
- 147 [4] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirec-
148 tional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 149 [5] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples.
150 In *International Conference on Learning Representations (ICLR)*, 2014.
- 151 [6] Minghao Hu, Yuxing Peng, Zhen Huang, Xipeng Qiu, Furu Wei, and Ming Zhou. Reinforced mnemonic
152 reader for machine reading comprehension. In *International Joint Conference on Artificial Intelligence*
153 *(IJCAI)*, 2018.
- 154 [7] Robin Jia and Percy Liang. Adversarial examples for evaluating reading comprehension systems. In
155 *Empirical Methods in Natural Language Processing (EMNLP)*, 2017.
- 156 [8] Mandar Joshi, Eunsol Choi, Daniel S Weld, and Luke Zettlemoyer. TriviaQA: A large scale distantly
157 supervised challenge dataset for reading comprehension. In *Association for Computational Linguistics*
158 *(ACL)*, 2017.
- 159 [9] Divyansh Kaushik and Zachary C Lipton. How much reading does reading comprehension require? a
160 critical investigation of popular benchmarks. In *Empirical Methods in Natural Language Processing*
161 *(EMNLP)*, 2018.
- 162 [10] Sewon Min, Victor Zhong, Richard Socher, and Caiming Xiong. Efficient and robust question answering
163 from minimal context over documents. In *Empirical Methods in Natural Language Processing (EMNLP)*,
164 2018.
- 165 [11] Jeffrey Pennington, Richard Socher, and Christopher Manning. Glove: Global vectors for word repre-
166 sentation. In *Proceedings of the 2014 conference on empirical methods in natural language processing*
167 *(EMNLP)*, 2014.
- 168 [12] Jonathan Raiman and John Miller. Globally normalized reader. In *Empirical Methods in Natural Language*
169 *Processing (EMNLP)*, 2017.
- 170 [13] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. SQuAD: 100,000+ questions for
171 machine comprehension of text. In *Empirical Methods in Natural Language Processing (EMNLP)*, 2016.
- 172 [14] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and
173 Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- 174 [15] Adam Trischler, Tong Wang, Xingdi Yuan, Justin Harris, Alessandro Sordani, Philip Bachman, and Kaheer
175 Suleman. NewsQA: A machine comprehension dataset. In *Workshop on Representation Learning for NLP*,
176 2017.
- 177 [16] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
178 Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing*
179 *Systems*.
- 180 [17] Yicheng Wang and Mohit Bansal. Robust machine comprehension models via adversarial training. In
181 *North American Chapter of the Association for Computational Linguistics: Human Language Technologies*,
182 2018.
- 183 [18] Adams Wei Yu, David Dohan, Minh-Thang Luong, Rui Zhao, Kai Chen, Mohammad Norouzi, and
184 Quoc V Le. QANet: combining local convolution with global self-attention for reading comprehension.
185 *International Conference on Learning Representations (ICLR)*, 2018.