

Analysis on the ELECTRA-small Model for "Why" Questions and Adversarial Example Challenge Set

Matthew Fuerst

Zhixin Cao

Abstract

This paper describes two different analyses on question and answering datasets in the ELECTRA-small model. The first investigates what type of questions are failing the question and answering task the most. The second investigates how adversarial examples challenge the base model and look into a possible way to increase the model's robustness. The results will show that "why" questions seem to perform the worst regardless of what training set was used. A post-processing method was applied during evaluation to remove certain words from the beginning of each model prediction and gold label. This improved the exact_match score of "why" questions by 4.5% and in aggregate 2.7% for the dataset as a whole. For the second analysis, the team combined multiple adversarial dataset for training and was able to achieve very competitive extract_match and F1 score for both SQuAD and various adversarial dataset during evaluation. Lastly, a model was run combining the post-processing and the robust model which yielded an improvement of 5.07% compared to the base SQuAD model.

1 Introduction

Inspired by the "Universal Adversarial Trigger" and how they utilize it to attack Reading Comprehension model in (Wallace et al., 2019), the team started analyzing the ELECTRA-small model, which could be found here: <https://huggingface.co/google/electra-small-discriminator>, with the baseline of "SQuAD" dataset from (Rajpurkar et al., 2016), for dissecting the performance of various question types that are classified by the task's question stem.

1.1 Various Question Types

The team read all the suggested dataset documents listed in the final paper description and found the

SQuAD paper (Rajpurkar et al., 2016) the most interesting for the question and answering tasks. The team scanned through various papers on different analysis methods in the final project description document, including contrast sets, checklists, adversarial challenge sets, model ablations, and lastly "competency problems" frameworks. The team thought that the adversarial challenge examples looked promising and decided to research it more. From there the team read three out of the five suggested papers for adversarial challenge sets (Jia and Liang, 2017). (Wallace et al., 2019). (Bartolo et al., 2020).

The papers gave the team a better understanding on how adversarial examples work and the different forms adversarial examples can take. After the team finalized the dataset and the analytical method, the next step was to dig into the code to try and understand how to make the model run. Trying to run the model locally on the CPU was taking hours so the team utilized Goolge Colab to train the models more quickly.

1.2 Adversarial Dataset

While the team was investigating the performance on different question types, it also came to their notice that the overall performance of the ELECTRA-small model that trained on only the SQuAD dataset, was quite disappointing, especially when the model was evaluated based on the adversarial dataset.

Figure 1 shows the performance metrics of the various datasets that the team tested including: base SQuAD, SQuAD Adversarial, and the Adversarial QA datasets and their subsets.

Please note that *AddOneSent* and *AddSent* are adversarial dataset from (Jia and Liang, 2017) and *adversarialQA*, *dbidaf*, *dbert* and *droberta* are adversarial dataset from (Bartolo et al., 2020).

Additionally, according to Jia and Liang (2017),

Performances on squad_only_trained model		
	exact_match	f1
Squad	76	84.28
AddOneSent	59.87	67.05
AddSent	49.71	56.67
adversarialQA	16.3	25.99
dbidaf	25.9	36.93
dbert	10.1	18.99
droberta	12.9	22.06

Figure 1: ELECTRA-small model’s Performance on adversarial dataset

”While adversarially perturbed images punish model oversensitivity to imperceptible noise, our adversarial examples target model overstability—the inability of a model to distinguish a sentence that actually answers the question from one that merely has words in common with it.”

In other words, the authors suspected that the model that trained on SQuAD dataset did not have the ability to truly understand the context and question, instead, made the prediction based on the pattern on the context or the question.

Therefore, the team decided to further look into potential reasons that could explain why the ELECTRA-small model underperforms on adversarial dataset.

2 Analysis

The analysis part could be divided into:

- Analysis on ”why” question type from the SQuAD dataset
- Analysis on possible reason(s) for the ELECTRA-small model to underperform on adversarial dataset

2.1 Analysis on ”why” questions from the SQuAD dataset

Another approach the team considered was checking what question types in the SQuAD dataset performed the worst. The team found that the questions could be separated into 8 main question groups: ”who”, ”what”, ”when”, ”where”, ”why”, ”which”, ”how”, and ”*other*”. ”*Other*” questions specify questions that did not start with any of the other words listed. The team found the various questions by printing out the question statements

Starting word in question:	Word Count:	% of Total:
who	1096	10.37%
what	4740	44.84%
when	691	6.54%
where	429	4.06%
why	150	1.42%
which	454	4.30%
how	1087	10.28%
other	1923	18.19%
	10570	100.00%

Figure 2: Count of Question Types in SQuAD

Starting word in question:	Number of successful exact_match predictions	Number of failed exact_match predictions	Total number of examples per question word	exact_match _success Rate
who	906	190	1096	82.7%
what	3394	1346	4740	71.6%
when	599	92	691	86.7%
where	296	133	429	69.0%
why	78	72	150	52.0%
which	340	114	454	74.9%
how	817	270	1087	75.2%
other	1444	479	1923	75.1%
	7874	2696	10570	74.5%

Figure 3: Exact Match Success Rate of SQuAD Question Words

from the training and validation examples in the code.

From the SQuAD dataset the count of the questions can be found in Figure 2.

As Figure 2 shows ”what” questions comprised the biggest share of the overall dataset at 44.84%, followed by ”*other*” questions at 18.19%, ”who” questions at 10.37%, and ”how” questions at 10.28%. These four make up 83.68% of all the examples. Before the team decided to pick on a specific question type to improve, the team wanted to analyze the exact_match success rate of each question type so that the team could focus on questions that had the lowest exact_match success rate. The exact_match success rate is computed by checking if the model prediction is exactly the same as one of the gold labels. The team modified the code to track the starting question word for each correct model prediction and each incorrect model prediction in terms of the exact_match score. Results can be found in Figure 3.

”Why” questions performed the worst at a 52.0% exact_match success rate. The next lowest was ”where” questions with a 69.0% success rate. The highest performing one was ”when” questions.

Once the team decided to focus on the ”why”

questions the next step was to try and figure out why the “why” questions were failing. To do so the team printed out the context, question, gold label, and model prediction for each “why” question that did not get an exact match with one of the gold labels. It became quickly evident that the model actually works well. The team categorized the incorrect responses into five main categories. Many of the differences between the gold labels and model prediction were minor and the major chunk of the correct answer was given in most cases.

Below are examples of the categorized incorrect model predictions (please note that these examples were taken from the SQuAD dataset and the squad_adversarial dataset, including *AddSent* and *AddOneSent* subsets):

1. The model was almost correct with the exception of an article or word at the beginning of the gold label or model prediction.

Valid gold label answers:

- southern China withheld and fought to the last
- because southern China withheld and fought to the last before caving in
- withheld and fought to the last

Model Prediction:

- southern China withheld and fought to the last before caving in

2. The model predicted too much following a comma in a sentence.

Valid gold label answers:

- solid economic growth
- solid economic growth
- economic growth

Model Prediction:

- solid economic growth, an increase in foreign investment as well as funding from the European Union

3. The model predicted some of the answer correct but not enough.

Valid gold label answers:

- a result of increasing crime and poverty
- increasing crime and poverty in the Hyde Park neighborhood

Starting word in question:	Number of successful exact_match predictions	Number of failed exact_match predictions	Total number of examples per question word	exact_match _success Rate
who	164	127	291	56.4%
what	891	820	1711	52.1%
when	113	84	197	57.4%
where	52	54	106	49.1%
why	16	47	63	25.4%
which	87	79	166	52.4%
how	209	161	370	56.5%
other	337	319	656	51.4%
	1869	1691	3560	52.5%

Figure 4: Exact Match Success Rate of SQuAD Adversarial *AddSent* Question Words

- increasing crime and poverty in the Hyde Park neighborhood

Model Prediction:

- increasing crime and poverty

4. The model combined multiple gold labels into an answer that was not considered valid.

Valid gold label answers:

- transgenes in these plastics cannot be disseminated by pollen
- chloroplasts are not inherited from the male parent
- transgenes in these plastics cannot be disseminated by pollen

Model Prediction:

- chloroplasts are not inherited from the male parent, transgenes in these plastics cannot be disseminated by pollen

5. The model predicted something completely incorrect.

Valid gold label answers:

- to guard against armed groups
- guard against armed groups
- guard against armed groups

Model Prediction:

- it infringed on democratic freedoms

The team looked at all the missed “why” examples in the three different datasets. The results of the exact match success rate for each of the following datasets can be found in Figures 3, 4, and 5. In the following section in Implementation Details the team will discuss what was done to improve this values.

Starting word in question:	Number of successful exact_match predictions	Number of failed exact_match predictions	Total number of examples per question word	exact_match _success Rate
who	120	64	184	65.2%
what	511	312	823	62.1%
when	75	38	113	66.4%
where	34	23	57	59.6%
why	10	22	32	31.3%
which	50	32	82	61.0%
how	112	72	184	60.9%
other	186	126	312	59.6%
	1098	689	1787	61.4%

Figure 5: Exact Match Success Rate of SQuAD Adversarial *AddOneSent* Question Words

2.2 Analysis on possible reason(s) that the ELECTRA-small model underperforms on the adversarial dataset - Disadvantage on the construction of SQuAD dataset?

It is evident that the ELECTRA-small model does not perform well on adversarial datasets, comparing with its performance on SQuAD dataset. While it is puzzling to detect what might have caused it, the research from [Jia and Liang \(2017\)](#) could provide some initial clues:

"A model that relies on superficial cues without understanding language can do well according to average F1 score, if these cues happen to be predictive most of the time. Weissenborn et al. (2017) argue that many SQuAD questions can be answered with heuristics based on type and keyword-matching.... Consider the example in Figure 1: the BiDAF Ensemble model originally gives the right answer, but gets confused when an adversarial distracting sentence is added to the paragraph."

As illustrated, Figure 6 is the Figure 1 referred by [Jia and Liang \(2017\)](#). With the doubt in mind, the team had continue exploring and analyzing on the SQuAD with *AddSent* from [Jia and Liang \(2017\)](#).

Looking into this example in Figure 7 from *AddSent* dataset, it is clear that the same context were given four times, with the exception of the last sentence, which appears to be an adversarial sentence. In one of the cases, the last sentence was "Hamster should abide by the general public and be avoided", which is apparently distracting and does not align with any other part of the context; yet, the model still provides the prediction result as "Hamster", given the last sentence has a very similar structure as the question stem.

Article: Super Bowl 50

Paragraph: "Peyton Manning became the first quarterback 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: An example from the SQuAD dataset. The BiDAF Ensemble model originally gets the answer correct, but is fooled by the addition of an adversarial distracting sentence ([in blue](#)).

Figure 6: Example of Figure 1 reproduced from [Jia and Liang \(2017\)](#)

The more data examples the team looked into, the more convincing it is that the ELECTRA-small model is indeed trained primarily based on the structure, types, or keyword-matching, without actually understanding the language itself. In order to further confirm this error case, the team decided to adopt the training methods suggested in [Jia and Liang \(2017\)](#), where they trained the initial model on adversarial dataset and discovered training on adversarial dataset could improve the performance of model to certain degree. For which, it will be further discussed in the following session.

3 Implementation Details

3.1 Post-processing

One thing that came to the team's attention is there were a lot of examples where the only difference between a gold label and a model prediction was just one word at the beginning of the sentence. Here is an example where having the word "the" appear in the prediction would have yielded an exact match:

Valid gold label answers:

- the defection of a number of Kenyan athletes to represent other countries
- defection of a number of Kenyan athletes
- a number of Kenyan athletes to represent other countries

context = One theory is that, while disobedience may be helpful, any great amount of it would undermine the law by encouraging general disobedience which is neither conscientious nor of social benefit. Therefore, conscientious lawbreakers must be punished. Michael Bayles argues that if a person violates a law in order to create a test case as to the constitutionality of a law, and then wins his case, then that act did not constitute civil disobedience. It has also been argued that the intent behind the breaking of the law, rather than the breaking of the law itself, is what constitutes civil disobedience. An act of breaking the law for self-gratification, in the case of a homosexual or cannabis user who does not direct his/her intent towards the reform of the law, is not civil disobedience. Likewise, an act of breaking the law in order to escape punishment by committing the crime covertly and avoiding attribution, or by denying having committed the crime, or by fleeing the jurisdiction, is generally viewed as not being a civil disobedient. Obedience by the general public should abide by the general principles of society.

question = Why should disobedience by the general public be avoided?

answer = [text: ['neither conscientious nor of social benefit', 'neither conscientious nor of social benefit', 'any great amount of it would undermine the law', 'conscientious lawbreakers'], answer_start: [147, 147, 55, 283]]

predictions[0][text] = not being a civil disobedient

.....

context = One theory is that, while disobedience may be helpful, any great amount of it would undermine the law by encouraging general disobedience which is neither conscientious nor of social benefit. Therefore, conscientious lawbreakers must be punished. Michael Bayles argues that if a person violates a law in order to create a test case as to the constitutionality of a law, and then wins his case, then that act did not constitute civil disobedience. It has also been argued that the intent behind the breaking of the law, rather than the breaking of the law itself, is what constitutes civil disobedience. An act of breaking the law for self-gratification, in the case of a homosexual or cannabis user who does not direct his/her intent towards the reform of the law, is not civil disobedience. Likewise, an act of breaking the law in order to escape punishment by committing the crime covertly and avoiding attribution, or by denying having committed the crime, or by fleeing the jurisdiction, is generally viewed as not being a civil disobedient. Obedience by the general public should abide by the general principles of society.

question = Why should disobedience by the general public be avoided?

answer = [text: ['neither conscientious nor of social benefit', 'neither conscientious nor of social benefit', 'any great amount of it would undermine the law', 'conscientious lawbreakers'], answer_start: [147, 147, 55, 283]]

predictions[0][text] = Hamster

.....

context = One theory is that, while disobedience may be helpful, any great amount of it would undermine the law by encouraging general disobedience which is neither conscientious nor of social benefit. Therefore, conscientious lawbreakers must be punished. Michael Bayles argues that if a person violates a law in order to create a test case as to the constitutionality of a law, and then wins his case, then that act did not constitute civil disobedience. It has also been argued that the intent behind the breaking of the law, rather than the breaking of the law itself, is what constitutes civil disobedience. An act of breaking the law for self-gratification, in the case of a homosexual or cannabis user who does not direct his/her intent towards the reform of the law, is not civil disobedience. Likewise, an act of breaking the law in order to escape punishment by committing the crime covertly and avoiding attribution, or by denying having committed the crime, or by fleeing the jurisdiction, is generally viewed as not being a civil disobedient. Obedience by the general public should abide by the general principles of society.

question = Why should disobedience by the general public be avoided?

answer = [text: ['neither conscientious nor of social benefit', 'neither conscientious nor of social benefit', 'any great amount of it would undermine the law', 'conscientious lawbreakers'], answer_start: [147, 147, 55, 283]]

predictions[0][text] = Obedience

.....

context = One theory is that, while disobedience may be helpful, any great amount of it would undermine the law by encouraging general disobedience which is neither conscientious nor of social benefit. Therefore, conscientious lawbreakers must be punished. Michael Bayles argues that if a person violates a law in order to create a test case as to the constitutionality of a law, and then wins his case, then that act did not constitute civil disobedience. It has also been argued that the intent behind the breaking of the law, rather than the breaking of the law itself, is what constitutes civil disobedience. An act of breaking the law for self-gratification, in the case of a homosexual or cannabis user who does not direct his/her intent towards the reform of the law, is not civil disobedience. Likewise, an act of breaking the law in order to escape punishment by committing the crime covertly and avoiding attribution, or by denying having committed the crime, or by fleeing the jurisdiction, is generally viewed as not being a civil disobedient. Obedience by the general public should abide by the general principles of society.

question = Why should disobedience by the general public be avoided?

answer = [text: ['neither conscientious nor of social benefit', 'neither conscientious nor of social benefit', 'any great amount of it would undermine the law', 'conscientious lawbreakers'], answer_start: [147, 147, 55, 283]]

predictions[0][text] = Hamster

Figure 7: Example of how additional distracting sentence could interrupt the prediction from ELECTRA-small model that only trained on SQuAD

Model Prediction:

- defection of a number of Kenyan athletes to represent other countries

Here is an example where having the word “to” removed from the prediction would have yielded an exact match:

Valid gold label answers:

- rebuild St. Peter’s Basilica
- rebuild St. Peter’s Basilica
- rebuild St. Peter’s Basilica

Model Prediction:

- to rebuild St. Peter’s Basilica

The team believed that removing articles from the beginning of a sentence did not impact the true meaning of the answer. The team settled on 5 different words: “a”, “an”, “to”, “the”, and “because”. Next the code was augmented two ways to see if removing all of these “article words” at the beginning of each answer or prediction would have an impact on the exact match score. The first was the gold labels were checked to see if the beginning of any answer started with any of these five “articles words”. If so the word was removed and the answer was updated. The second way the code was augmented was by performing the same task

Starting word in question:	Number of successful exact_match predictions	Number of failed exact_match predictions	Total number of examples per question word	exact_match success Rate	Number of improved exact_matches
who	930	166	1096	84.85%	+24
what	3573	1167	4740	75.38%	+179
when	610	81	691	88.28%	+11
where	321	108	429	74.83%	+25
why	85	65	150	56.67%	+7
which	355	99	454	78.19%	+15
how	822	265	1087	75.62%	+5
other	1498	425	1923	77.90%	+54
all	8194	2376	10570	77.52%	+320

an increase in ... 3.03%

Figure 8: Exact Match Success Rate of SQuAD Question Words with Post-Processing Applied

Starting word in question:	Number of successful exact_match predictions	Number of failed exact_match predictions	Total number of examples per question word	exact_match success Rate	Number of improved exact_matches
who	167	124	291	57.4%	+3
what	933	778	1711	54.5%	+42
when	113	84	197	57.4%	0
where	57	49	106	53.8%	+45
why	19	44	63	30.2%	+3
which	87	79	166	52.4%	0
how	209	161	370	56.5%	0
other	353	303	656	53.8%	+16
all	1938	1622	3560	54.4%	+69

an increase in ... 1.94%

Figure 9: Exact Match Success Rate of SQuAD Adversarial AddSent Question Words with Post-Processing Applied

of removing “article words” except this time for the model predictions. The code then continued as previously run. The results can be found in Figures 8, 9, and 10 on the various datasets that the team tested.

The original base SQuAD dataset had an exact match improvement of 3.03% after the post processing rules were applied. The *AddSent* Adversarial examples had an improvement of 1.94%. Lastly, the *AddOneSent* adversarial examples had an improvement of 2.29%.

It should be noted the trained model was not changed to make these improvements. The post processing was updated with the rule that no answer or prediction can start with the words “a”, “an”,

Starting word in question:	Number of successful exact_match predictions	Number of failed exact_match predictions	Total number of examples per question word	exact_match success Rate	Number of improved exact_matches
who	123	61	184	66.85%	+3
what	535	288	823	65.01%	+24
when	75	38	113	66.37%	0
where	36	21	57	63.16%	+2
why	11	21	32	34.38%	+1
which	50	32	82	60.98%	0
how	112	72	184	60.87%	0
other	197	115	312	63.14%	+11
all	1139	648	1787	63.74%	+41

an increase in ... 2.29%

Figure 10: Exact Match Success Rate of SQuAD Adversarial AddOneSent Question Words with Post-Processing Applied

Performances on full_trained model		
	exact_match	f1
Squad	79.91	87.16
AddOneSent	97.14	98.95
AddSent	97.33	98.98
adversarialQA	72.36	82.71
dbidaf	76.6	86.23
dbert	72.3	82.07
droberta	68.2	79.81

Figure 11: ELECTRA-small model’s Performance after training on various adversarial dataset

“to”, “the”, or “because”.

The team was surprised with how much an improvement contributed by the removal of these 5 words had on the evaluation criteria of the dataset. After further looking into the data, it seemed that the training examples present an inconsistency where sometimes they contained articles like ”a” before the answer statement and sometimes they did not. This inconsistency could have attributed to the poor performance of the validation set with these examples.

3.2 Train on adversarial dataset

To train the ELECTRA-small model, as a method introduced by (Jia and Liang, 2017), the team gathered few Huggingface open adversarial dataset and concatenated them together, including `squad_adversarial` (https://huggingface.co/datasets/squad_adversarial) and `adversarial_qa` (https://huggingface.co/datasets/adversarial_qa), where `squad_adversarial` dataset is from (Jia and Liang, 2017) and `adversarial_qa` is from (Bartolo et al., 2020).

In particular, there are two subsets of data in `squad_adversarial` dataset, one is `AddSent`(1,803 examples) and another one is `AddOneSent` (3,591 examples). For `adversarial_qa` dataset, there are four categories of data subsets, which are `adversarialQA` (3,000 examples), `dbidaf` (1,000 examples), `dbert` (1,000 examples), and `droberta` (1,000 examples). With a combination of all training and validation dataset, including the original SQuAD, there were a total of 159,208 examples to train on. Although it took a while to train, the result was astonishing.

As indicated in Figure 11, comparing with the model’s performance when it was only trained on

Starting word in question:	Number of successful exact_match predictions	Number of failed exact_match predictions	Total number of examples per question word	exact_match success Rate	Number of improved exact_matches
who	954	142	1096	87.04%	+48
what	3648	1092	4740	76.96%	+254
when	618	73	691	89.44%	+19
where	315	114	429	73.43%	+19
why	92	58	150	61.33%	+14
which	375	79	454	82.60%	+35
how	850	237	1087	78.20%	+33
other	1558	365	1923	81.02%	+114
	8410	2160	10570	79.56%	+536

an increase in ... 5.07%

Figure 12: Exact Match Success Rate of Adversarial Trained SQuAD Question Words with Post-Processing Applied

Starting word in question:	Number of successful exact_match predictions	Number of failed exact_match predictions	Total number of examples per question word	exact_match success Rate	Number of improved exact_matches
who	289	2	291	99.3%	+125
what	1660	51	1711	97.0%	+769
when	194	3	197	98.5%	+81
where	96	10	106	90.6%	+44
why	54	9	63	85.7%	+38
which	164	2	166	98.8%	+77
how	343	27	370	92.7%	+134
other	633	23	656	96.5%	+296
	3433	127	3560	96.4%	+1564

an increase in ... 43.93%

Figure 13: Exact Match Success Rate of Adversarial Trained SQuAD Adversarial `AddSent` Question Words with Post-Processing Applied

SQuAD which could be found in Figure 1, the evaluation on the model based on the initial SQuAD, and various adversarial dataset have been dramatically improved. For example, the ELECTRA-small model achieves 79.91 extract-match and 87.16 F1 score on the initial SQuAD, while the model was only trained on SQuAD, the extract-match was 76 and F1 score was 84.28, with 3 epochs. Especially for the `squad_adversarial` dataset, the fully trained model was able to get both extract-match and F1 scores for `AddSent` and `AddOneSent` over 90. Not only does the fully trained model outperforms the initial model’s performance on `squad_adversarial` dataset, it also outperforms the evaluations on other dataset.

3.3 Combining the Post-Processing with Adversarial Trained Data

One of the last tasks the team did was combine both of the approaches listed above and run an analysis on the result. This includes taking the robust model that was trained in section 3.2 and applying the post-processing as described in section 3.1. Results on the SQuAD and `squad_adversarial` datasets with the `AddSent` and `AddOneSent` subsets can be found in Figures 12, 13, and 14.

The model performed extremely well on the ad-

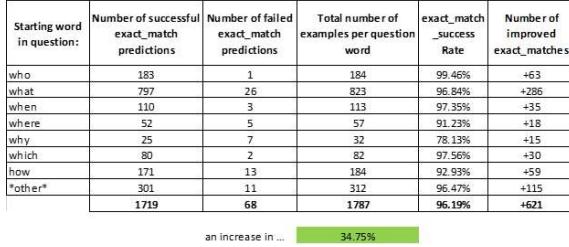


Figure 14: Exact Match Success Rate of Adversarial Trained SQuAD Adversarial *AddOneSent* Question Words with Post-Processing Applied

versarial challenge validation sets getting over a 96% exact_match score. The base SQuAD model when removing the article words and training on the adversarial examples had an increase in performance of 5.07%. The adversarial SQuAD: *AddSent* model had an improvement of 43.93% and the SQuAD: *AddOneSent* model had an improvement of 34.75%. Even after an extremely robust model from the adversarial examples the exact_match score was still further improved by the addition of the post-processing rule removing certain starting words.

4 Conclusion

Question and answering models are extremely important aspects of web search applications. This paper describes two different analyses performed on the SQuAD dataset. The first analysis looked into which types of questions were failing the most and found "why" questions had the lowest exact_match success rate. The team viewed all the incorrect "why" questions in depth and found that the model makes very similar mistakes on multiple examples. In many examples, the model prediction was almost the same as a gold answer label, except for an added article or word at the beginning of the model prediction. After removing various article words from both the gold labels and the model prediction and running our analysis again we had an aggregate 2.70% improvement on all types of questions, compared with the base model when trained on SQuAD only. The "why" questions specifically improved 4.5% after the post-processing rule of removing certain start words was applied.

The second analysis looked into how training on adversarial examples improved the dataset compared to a model that was not training on any adversarial examples. The team combined various adversarial datasets with SQuAD to create a larger

and more effective training set. After the model was trained on this combined datasets, the results showed that not only did this model improve it's F1 and exact_match score when validating adversarial challenge examples, it had also improved scores over the normal SQuAD validation set as well. This showed that training on adversarial examples could create a more robust model that performs better on variations in questions. Lastly, after the updated model was improved by the adversarial examples during training, and the post-processing was run the final model improved by 5.07% compared to the base SQuAD model.

Acknowledgments

Special acknowledgements to Dr. Durrett and all TAs who have been such a great help to guide us on conducting this paper. The TAs and fellow students were great resources to bounce ideas off of and where helpful in answering questions whenever we needed help.

References

- Max Bartolo, Alastair Roberts, Johannes Welbl, Sebastian Riedel, and Pontus Stenetorp. 2020. [Beat the AI: Investigating adversarial human annotation for reading comprehension](#). *Transactions of the Association for Computational Linguistics*, 8:662–678.
- Robin Jia and Percy Liang. 2017. [Adversarial examples for evaluating reading comprehension systems](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2021–2031, Copenhagen, Denmark. 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.
- Eric Wallace, Shi Feng, Nikhil Kandpal, Matt Gardner, and Sameer Singh. 2019. [Universal adversarial triggers for attacking and analyzing NLP](#). 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 2153–2162, Hong Kong, China. Association for Computational Linguistics.