# Robustness Analysis on Natural Language Processing Based AI Q&A Robots

Chengxiang Yuan[1], Mingfu Xue[1(✉)], Lingling Zhang[1], and Heyi Wu[2]

[1] College of Computer Science and Technology,
Nanjing University of Aeronautics and Astronautics,
Nanjing, China
{yuancx,mingfu.xue}@nuaa.edu.cn, bluezhll@126.com
[2] School of Cyber Science and Engineering,
Southeast University, Nanjing, China
why1988seu@126.com

**Abstract.** Recently, the natural language processing (NLP) based intelligent question and answering (Q&A) robots have been used in a wide range of applications, such as smart assistant, smart customer service, government business. However, the robustness and security issues of these NLP based artificial intelligence (AI) Q&A robots have not been studied yet. In this paper, we analyze the robustness problems in current Q&A robots, which include four aspects: (1) semantic slot settings are incomplete; (2) sensitive words are not filtered efficiently and completely; (3) Q&A robots return the search results directly; (4) unsatisfactory matching algorithms and inappropriate matching threshold settings. Then, we design and implement two types of evaluation tests, bad language and user's typos, to evaluate the robustness of several state-of-the-art Q&A robots. Experiment results show that these common inputs (bad language and user's typos) can successfully make these Q&A robots malfunction, denial of service, or speaking dirty words. Besides, we also propose possible countermeasures to enhance the robustness of these Q&A robots. To the best of the authors' knowledge, this is the first work on analyzing the robustness and security problems of intelligent Q&A robots. This work can hopefully help provide guidelines to design robust and secure Q&A robots.

**Keywords:** AI security · Question and answer robots · Robustness · Natural language processing

## 1    Introduction

In recent years, artificial intelligence (AI) techniques achieved major break-throughs and have been used ubiquitously. A representative application of AI is the natural language processing (NLP) techniques based intelligent question and answering (Q&A) robots, which are widely used in general application areas, professional business, and government applications.

Traditionally, people can search various information from a search engine, e.g., Google, Baidu. The search engine will return a ranked list of related web documents according to the user's input. The user cannot get the answers they want quickly and accurately from a large number of search results. Unlike the information retrieval system, the task of an intelligent Q&A robot is to give the user a precise and concise answer in several interactions with the user. Generally, the Q&A robots have the following two features: (1) users can query the Q&A robots in natural language; (2) the Q&A robot directly returns the answer that the user needs (rather than a ranked list of relevant documents). Recently, many companies have developed their own Q&A robots, which have made great progress in human-computer interaction, such as Google assistant [1], Cortana [2], Siri [3], Alexa [4], Watson [5], DuerOS [6], Ali Xiaomi [7], JD Instant Messaging Intelligence [8]. On the other hand, many government agencies have also provided Q&A robots for public business, e.g., 12306 (Q&A robots of Chinese railway system), tax bureau.

However, the robustness of these NLP based AI Q&A robots has not been studied yet in the literature. In this paper, first, we review the working principles of current Q&A robots. Then, we analyze the robustness problems of current intelligent Q&A robots, which include four aspects: (1) the semantic slots settings are incomplete; (2) it is difficult for Q&A robots to filter all the sensitive words in both the user's questions and the returned answers; (3) some Q&A robots directly return the results from a search engine; (4) unsatisfactory matching algorithms or inappropriate matching thresholds. Then, we design and implement two types of evaluation tests, bad language and user's typos, on several state-of-the-art Q&A robots. These common inputs (bad language and typos) can successfully make the Q&A robots malfunction, denial of service (DoS), or speaking dirty words. These consequences are disastrous to the Q&A robots, which will face a recall or withdraw from the market. Experiment results show that the semantic slots of the tested Q&A robots are incomplete, and the Q&A robots do not consider the contextual information in the multi-round interactions with the user. For sensitive words filtering, current Q&A robots still do not have an effective solution. The matching degree between the answer given by the Q&A robot and the user's question is unsatisfactory, and the accuracy of the answers needs to be improved. In addition, we propose several possible countermeasures for these robust problems of the intelligent Q&A robots.

## 2    Working Principles of AI Q&A Robots

Generally, there are three kinds of working principles used by current AI Q&A robots, using the knowledge base (KB), using information retrieval (IR), and using guiding questions. We will describe these three types of working principles in Sects. 2.1, 2.2, and 2.3, respectively. In Sect. 2.4, we will summarize and present the most complete working principle of state-of-the-art AI Q&A robots, which is a combination of the above three mechanisms.

### 2.1    Knowledge Base Based Q&A Robots

The process of knowledge base (KB) based Q&A robots is shown in Fig. 1. After semantic parsing of the input questions, a standard structured query languages will be generated which will be searched directly in the KB. Then, it returns the answer to the user. There are several large-scale knowledge bases, such as DBpedia [9], Freebase [10], and YAGO [11], which store a large amount of valuable information in the form of Resource Description Framework (RDF) triples [12]. The key of KB based Q&A robots is to transform users' natural language questions into standard structured query formats.
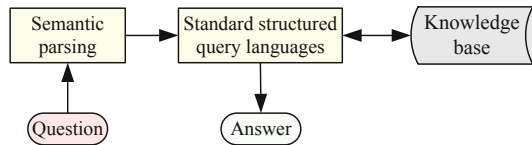


**Fig. 1.** The process of KB-based Q&A robots.

The early Q&A robot is a natural language interface to database (NLIDB) that allows the user to query and access information from a database in natural language questions [13]. However, this method can only be used in domain-specific applications, while in open-domain question answering, the database will be extremely complex. In order to solve this problem, semantic parser is used to analyze the user's natural language questions. Semantic parsing includes four steps [14]: paraphrase extraction, mapping to formal meaning representation, semantic combination, data retrieval from the knowledge base. There are many semantic parsing methods [15–18]. In [15], Zettlemoyer *et al.* use a learning algorithm to map natural language sentences to a lambda calculus encoding of their semantics. Zelle *et al.* [16] use shift-reduce derivations to map sentences into database queries. Wang *et al.* [17] exploit statistical machine translation techniques to generate logical forms. Lu *et al.* [18] construct a generative model based on a hybrid tree whose nodes contain natural language words and meaning representation tokens.

Traditional semantic parsers require annotated logical forms as supervision and have only a few logical predicates [14]. Therefore, semantic understanding

methods based on deep learning have been proposed recently. Yih *et al.* [19] first decompose each question into an entity mention and a relation pattern. Then, they use convolutional neural network models to measure the similarity between entity mentions and entities, and the similarity between relation patterns and relations in the KB. Yih *et al.* [20] use an entity linking system and deep convolutional neural network for question answering.

### 2.2   Information Retrieval Based Q&A Robots

Information retrieval (IR) based Q&A robots search for unstructured text documents and reorganize them into answers. These unstructured text documents, which are obtained from the web page using search engines, may contain the answer of the questions. Figure 2 presents the framework of a IR based Q&A robot, which includes three steps: question processing, passage retrieval and answer processing [21]. The question processing module consists of query formulation and answer type detection. The query formulation extracts one or more keywords from natural language questions. The answer type detection mainly extracts key information from the question and further explains the question in order to predict the answer type [14]. The passage retrieval module uses these keywords to retrieve the information from the web and obtain relevant documents that potentially contain the answer [14]. The relevant documents are ranked according to the matching degree of the question. Then, passages that contain potential answers are extracted and ranked. The answer processing module determines the type of the answer, and selects the most appropriate answer from the candidate answers by using a answer extraction algorithm. Pattern matching and N-gram tiling are two typical algorithms for extracting answers [14]. For example, Ravichandran *et al.* [22] develop a method to automatically learn patterns which can be used to find answers. [23] and [24] exploit the N-gram tiling method to extract answer.
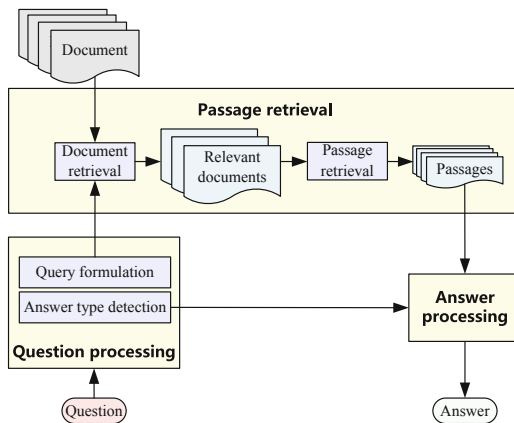


**Fig. 2.** The framework of IR-based Q&A robots [21].

## 2.3 Using Guiding Questions

When a Q&A robot can only understand a part of the question or consider the question to be ambiguous, it can further ask the user some guiding questions. These guiding questions can help the Q&A robot understand the user's intentions and return a more accurate answer. For example, suppose we want to query the 12306 Q&A robot for the train information from Shanghai to Beijing. If the Q&A robot does not use the guiding questions, it will return all the train information from Shanghai to Beijing. On the contrary, if the Q&A robot applies the guiding question technique, it can ask the user what time period of the trains the user wants to query. In this way, the returned result will be more accurate.

## 2.4 The Complete Working Principle of the State-of-the-art AI Q&A Robots

In this section, we will summarize and present the most complete working principle of the state-of-the-art AI Q&A robots, which is a combination of the above three mechanisms, as shown in Fig. 3. The Q&A robot first searches for similar questions from a question and answer database. The question and answer database contains a large number of questions and corresponding answers. The Q&A robot needs to perform text segmentation on the user's question. Keywords are extracted from the word segments, which are represented as a vector. The vector is used to match the answer in the Q&A database. The Q&A robot ranks the answers according to the matching degree, and returns the top $k$ answers to the user. For example, if there are 4000 questions and 4000 corresponding answers in a Q&A database, the vector is used to match the keywords of 4000 questions, and the top 3 or top 5 answers are returned to the user. If the match
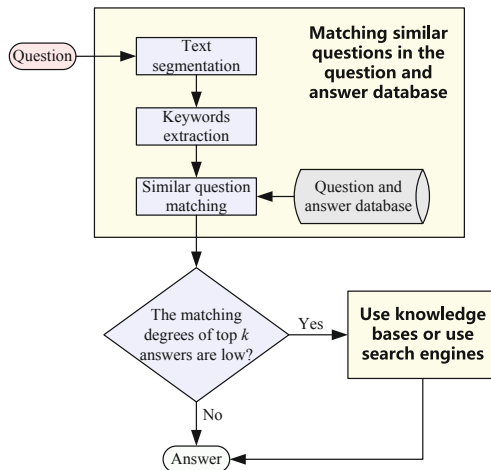


**Fig. 3.** The complete working principle of the state-of-the-art Q&A robots.

degrees between the user's question and the top $k$ answers are very low, the Q&A robot will use other methods, e.g., using knowledge bases, using search engines, or using guiding questions, to further determine the answer.

## 3   Robustness Problems of Current Q&A Robots

Although there have been many studies on techniques used by the Q&A robots in the literature, there are still many robustness problems in the practical Q&A robots. For example, the news reported that the intelligent Q&A robot XXX1 insulted a user [25], as shown in Fig. 4. Robustness/security problems in Q&A robots will lead to serious consequences. We analyze the robustness problems of current Q&A robots in the following subsections.



**Fig. 4.** News: the XXX1 Q&A robot insulted a user [25].

### 3.1   Semantic Slot Settings Are Incomplete

The slot filling process extracts key information from the user's questions, which has a very important impact on the quality of the whole question and answer system. In order to understand the user's questions correctly, the Q&A robot needs to use a semantic frame that contains different slots, and fills the key information in the semantic slots by interacting with the user.

Many semantic template-based methods are used in early slot filling techniques. In these method, recognition templates are constructed manually and used to extract the key information of sentences [26]. However, the construction of templates requires a lot of manual efforts, and the scope of application of the template is limited. A major method to solve the slot filling problem is to employ the sequence labeling model [27]. This method uses the sequence labeling model to mark the label for each word of a sentence, and generates the final slot filling

result based on these labels. Recently, researches achieve better performance in slot filling by using recurrent neural network, such as [28,29]. However, there are still many problems in slot filling of the practical Q&A robots. Their semantic slots settings are incomplete, and the performance of many slot filling methods is unsatisfactory. Besides, in the multi-round interactions with the user, it can't effectively translate the user's initial intentions into explicit instructions.

### 3.2   Sensitive Words Are Not Filtered Effectively and Completely

In many cases, there are sensitive words in the questions that users enter into the Q&A robot, such as foul words, violent words, sensitive political words. The Q&A robot will directly use these sensitive words to search for answers in the knowledge base or in the Internet, and reply related results that are also likely to have these sensitive words to the user. This could be a disaster. However, many Q&A robots do not filter or cannot effectively filter sensitive words entered by the user. The reasons include the following two aspects. On the one hand, due to the variety of sensitive words, the number of sensitive words is extremely large, while the Q&A robot lacks a complete database of sensitive words. On the other hand, when detecting sensitive words, normal words may be incorrectly identified as sensitive words.

In addition, there also may be sensitive words in the process of extracting and generating the answers. Since the answers are typically obtained from the Internet, if the system cannot correctly identify those sensitive words, the answers returned to the user may contain these sensitive words, which will lead to a bad user experience.

### 3.3   Return the Search Results Directly

In general, after the Q&A robot obtains relevant documents that contain potential answers from the Internet, the Q&A robot needs to rank these relevant documents according to the matching degree of the question [14]. Some Q&A robots use a single feature to match and rank these relevant documents, while others rank these relevant documents by combining different features. However, we found that many Q&A robots on the market do not have this process. These Q&A robots directly return the results obtained from the search engine to the user, which results in a low matching degree between the answer and the question. Therefore, these answers are always unable to meet the users' needs.

### 3.4   Unsatisfactory Matching Algorithms and Inappropriate Matching Threshold Settings

In the process of matching the answers to the question, a predefined threshold is used to remove documents with low relevance to the question, and those documents that are highly correlated with the question are reserved [21]. However, the matching degrees between all the relevant documents and the question

may be lower than the predefined threshold. In this case, the answer processing module cannot extract an answer to the question. Most current Q&A robots have two solutions to this issue. One solution is to tell the user directly that the question cannot be answered. Another solution is to extract the answer from the previously ranked documents and return it to the user. However, both these two solutions are unsatisfactory. The first solution will give users a bad user experience. The second solution will give users an irrelevant answer to the question.

Similarly, Q&A robots will also have this problem when ranking passages. Although numerous features are used to rank passages, these features may still be not good enough. The candidate answers may be inappropriate answers.

## 4  Bad Language Experiments

We design and implement five bad language tests in different ways to evaluate current Q&A robots. It is shown that most of the Q&A robots have more than one robustness problems. Note that, in order to protect the privacy and interests of these platforms, the names of these Q&A robots will be hidden in this paper, and replaced by numbers XXX2–XXX6 (XXX1 has been illustrated in Sect. 3). We also want to apologize for using some rude words to induce the system errors in the experiments.

### 4.1  Evaluation on the XXX2 Q&A Robot

As shown in Fig. 5, in the conversation with the XXX2 Q&A robot, we want to inquire the operation information of train No. G7097. Firstly, we only input the train number to the Q&A robot. The robot cannot query information about train No. G7097. At that time, the G7097 was a train that suddenly stop running due to irresistible factors, such as weather. Then, we ask the Q&A robot when does the train start running again. Because the Q&A robot lacks information of the train, it does not answer directly. To further understand the user's question, the robot asks the user to enter the train number. It also provides a list of other frequently asked questions to the user. We enter the train number again. The robot still cannot query information about the train. This experiment shows that, the XXX2 Q&A robot does not adopt slot filling technique and the guiding question technique. Its ability to understand the user's questions is unsatisfactory. Besides, the Q&A robot also does not consider the contextual information in the multi-round interactions with the user. Compared with the usual human customer service, the performance of this Q&A robot system is unsatisfactory.

### 4.2  Evaluation on the XXX3 Q&A Robot

Figure 6 shows our tests on the XXX3 Q&A robot. We input the sensitive word "shit" to the Q&A robot. Unexpectedly, the Q&A robot returns a sensitive sentence "Where to go shit" to the user without filtering them at all. Then, we input the sensitive sentence "You are a dog eating shit". The Q&A robot
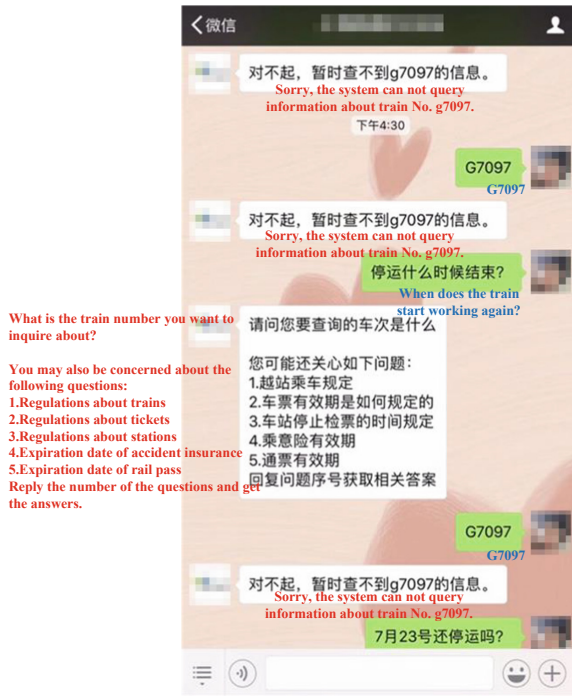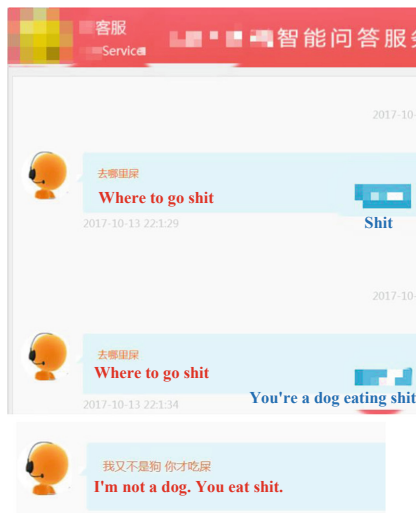
**Fig. 5.** Evaluation on the XXX2 Q&A robot.



**Fig. 6.** Evaluation on the XXX3 Q&A robot

answers "I am not a dog. You eat shit." The Q&A robot incorrectly fills the "shit dog" into different semantic slots in the slot filling process, resulting in very serious results. It is shown that if the input question contains sensitive words, the answers returned by the Q&A robot also contain sensitive words. To make things worse, inappropriate slot filling methods can even lead the Q&A robot to insult users, which will cause a very bad user experience.

### 4.3    Evaluation on the XXX4 Q&A Robot

Figure 7 shows our tests on the XXX4 Q&A robot. This Q&A robot is a very mature commercial product with powerful functions, which has been widely used. We input a sensitive word in Korean "You are an asshole". The robot replies "What's the problem with you in the process of 'asshole' (in Korean)?", in which the sensitive word is also replied in Korean. The Q&A robot does not filter other forms or other languages of sensitive words, and directly returns the sensitive word to the user in Korean. Therefore, when facing various sensitive words, the solutions of current popular business Q&A robots are unsatisfactory.



**Fig. 7.** Evaluation on the XXX4 Q&A robot.

### 4.4    Evaluation on the XXX5 Q&A Robot

Figure 8 shows our tests on the XXX5 Q&A robot, which includes three rounds. In the first round, as shown in Fig. 8(a), we enter an sensitive sentence "You are an asshole" into the Q&A robot in a different language. Like the XXX3 Q&A robot, the Q&A robot also does not filter the sensitive words and returns the translation of this sentence in Chinese to the user directly. Additionally, the Q&A robot replies a length of voice to the user. It is shown that this Q&A robot doesn't filter sensitive words and directly uses the translation results on the Internet. It also doesn't process the search results but returns it to the user directly.

After that, the developers of the XXX5 Q&A robot have fixed the problem of sensitive words. In our second round of experiment, as shown in Fig. 8(b), we enter the same sensitive sentence in different languages. The Q&A robot replies "Dear, I am glad to talk to you. If you have any questions, I am glad to help you.", "Little AI is thinking. Mom said that the child who likes to think will be more clever.", and a length of voice. Then, we input the same sensitive sentence again. The robot replies four periods, a length of voice, and "Hello, welcome to follow this account", and another piece of voice. Obviously, there is no correlation between the answers returned by the Q&A robot and the input questions. In this round, the Q&A robot filters the sensitive words but use results with inappropriate matching algorithms or low matching thresholds.

In the third round of experiment, as shown in Fig. 8(c), we still enter the same sensitive sentence in different languages. In this time, the system stops the service to the user, which shows "System error, please try again later".

This experiment illustrates that some Q&A robots return the search results directly, don't filter the sensitive words, and use unsatisfactory matching algorithms or inappropriate matching thresholds.



**Fig. 8.** Three rounds of tests on the XXX5 Q&A robot: (a) The first round; (b) The second round; (c) The third round.

## 4.5    Evaluation on the XXX6 Q&A Robot

As shown in Fig. 9, we test the XXX6 Q&A robot. This Q&A robot is also a mature commercial version, which has been widely used. We input "You are an asshole" in different languages one after another. The Q&A robot replies "What are you not satisfied with?", "How to bind your cell phone to your account? How to logout your account?", "Sorry, I am still learning. I can't understand your question right now. I can provide you with the following services.", respectively. There are no sensitive words in the answers returned by this Q&A robot. This Q&A robot also asks the users guiding questions to help it understand the user's intentions. However, it gives different answers each time, and all these answers are irrelevant to the question.

After several rounds of conversation with the user, the Q&A robot still does not understand the user's question, and it gives up answering this question. Therefore, this Q&A robot, although has sensitive words filtering and uses guiding questions, it still has the following problems: use unsatisfactory matching algorithms; replies irrelevant answers to the user; and gives up answering the question when failed several times.



**Fig. 9.** Evaluation on the XXX6 Q&A robot.

## 5    User's Typo Experiments

When a user interacts with a Q&A robot, the input questions are likely to contain some typos. A tiny typo may cause the Q&A robot to misunderstand the question and thus unable to answer the question correctly. In this section, we evaluate the robustness of Q&A robots when facing tiny typos.

### 5.1 Experimental Setup

**Datasets:** Three Q&A datasets, WebQuestionsSP [30], CuratedTREC [31] and WikiMovies-10k [32] are used to evaluate the accuracy and robustness of Q&A robots. The WebQuestionsSP dataset contains semantic parses for the questions from the WebQuestions dataset. The CuratedTREC dataset is collected from TREC1999, TREC2000, TREC2001 and TREC2002 data. The WikiMovies-10k dataset contains questions related to movies. The number of questions contained in the three datasets are 4737, 2180 and 10000, respectively.

**Target Q&A Robots:** In the experiment, two state-of-the-art Q&A robots, Siri [3] and Zo [33], are used as the target Q&A robots. Siri is a virtual assistant designed by Apple, which provides a user interface to answer questions. Users can ask Siri questions in natural language to obtain answers. Similarly, Zo is an intelligent chatting robot developed by Microsoft. Users can chat with Zo or ask Zo questions. In the experiment, we use the question and answer functions of Siri and Zo to evaluate their robustness.

### 5.2 Experimental Results of User's Typos

First, we evaluate the accuracy of Siri and Zo answering the original questions in the three Q&A datasets. The accuracy of Siri and Zo answering questions is shown in Table 1. For the questions in the three datasets, both Siri and Zo can only correctly answer a small number of the questions. This shows that the accuracy of current Q&A robot in answering questions is still low.

**Table 1.** The accuracy of Siri and Zo answering questions on the three Q&A datasets.

| Target Q&A robot | WebQuestionSP | CuratedTREC | WikiMovies-10k |
|---|---|---|---|
| Siri | 20.61% | 22.68% | 10.57% |
| Zo | 26.24% | 35.23% | 13.72% |

Since it is meaningless to generate typos with those original questions that the Q&A robot cannot answer, only the questions in the three datasets that the Q&A robot can answer correctly are selected to form new datasets. After that, we slightly modify these questions that the Q&A robot can answer correctly to generate questions with tiny typos. We use two methods to modify the original questions: (1) randomly replace one or more letters in a word; (2) replace the words in the original question with words that are spelled similarly. These generated questions with typos are used to evaluate the robustness of the Q&A robots. Table 2 presents three examples of questions with typos and the corresponding answers returned by the target Q&A robots, where the underlined letters represent the difference between the modified question and the original question. Although the target Q&A robots can correctly answer the original questions,

they cannot give the correct answer when facing the questions with tiny typos. Moreover, the answer to the modified question is very different from the original answer. Table 3 shows the accuracy of Siri and Zo answering questions on the three modified Q&A datasets. It is shown that the typos in the questions will significantly reduce the accuracy of the Q&A robot answering questions. Compared with Zo, Siri is less robust to these typos in the questions. For many questions that contain typos, Siri cannot return correct answers. Obviously, the robustness of current Q&A robots needs to be further improved.

**Table 2.** Three examples of questions with typos and the corresponding answers returned by the target Q&A robots.

| 1 | Original question | What is the density of gold? |
|---|---|---|
| | Modified question | What is the destiny of gold? |
| | Original answer | 19.3 grams per cubic centimeter |
| | Siri answer | I didn't find anything for "What is the destiny of gold?" |
| | Zo answer | The gold with the power |
| 2 | Original question | Who does allen iverson play for now 2010? |
| | Modified question | Who does allen iverson pray for now 2010? |
| | Original answer | Philadelphia 76ers |
| | Siri answer | Allen Iverson played for the Denver Nuggets, Detroit Pistons, Memphis Grizzlies and Philadelphia 76ers in the NBA between 1996 and 2009 |
| | Zo answer | One is definitely good for Allen Iverson |
| 3 | Original question | What did ryan dunn died from? |
| | Modified question | What did ryan dunn diked from? |
| | Original answer | Traffic collision |
| | Siri answer | The result is searched from the web |
| | Zo answer | Woops, sorry having trouble reading that... |

**Table 3.** The accuracy of Siri and Zo answering questions on three modified Q&A datasets, where only the questions in the three datasets that the Q&A robot can answer correctly are retained and be used to generate questions with typos.

| Target Q&A robot | WebQuestionSP | CuratedTREC | WikiMovies-10k |
|---|---|---|---|
| Siri | 36.57% | 42.62% | 35.33% |
| Zo | 73.66% | 68.89% | 69.21% |

## 6   Proposed Countermeasures and Challenges

### 6.1   Use Explainable Machine Learning Models

When developers use machine learning to analyze the semantics of a question, they should try to choose or develop explainable models, e.g., linear models,

decision trees. In this way, when the answer of the Q&A robot is wrong, they can modify the model accordingly. If they use complex nonlinear model, such as Deep Neural Networks, they cannot tell the details of the model or explain the generation process of the results. Once the Q&A robot makes mistakes, it is difficult to find out the reasons of these errors, and these problems cannot be fixed.

## 6.2   Rule-Based Filtering

Developers can manually create rules based on the characteristics of sensitive words. According to these rules, the Q&A robot detects and filters sensitive words in the user's questions and output answers. However, the performance of filtering sensitive words of this method is poor. There are a large number of and various kinds of sensitive words, just like a 'word games'. It is extremely difficult to define complete rules. Once the sensitive words are slightly changed, the filtering rules will probably fail. On the other hand, sometimes some normal words are misidentified as sensitive words by these rules.

## 6.3   Fuzzy Matching

Developers can also generate a dictionary of sensitive words to fuzzy match the user's questions. However, the efficiency of this method is low. Like the rule-based filtering method, this method may also fail when facing various kinds of sensitive words and their variants. Therefore, current Q&A robots do not have an effective way to filter sensitive words.

## 6.4   Spell Checking

A spell checking module is needed for the Q&A robots. If the Q&A robot cannot obtain the correct answer after searching in the knowledge base and related documents, the Q&A robot should use the spell checking module to check the question for typos. If there are typos in the question entered by the user, the Q&A robot corrects the typos in the question and confirms it to the user. Then, the Q&A robot can search the answer based on the revised correct question.

# 7   Conclusion

In this paper, we present four robustness (security) issues of current NLP based AI Q&A robots, which will lead to serious consequences to the Q&A robots. We design and implement two types of tests, which simulated two types of common inputs, bad language and user's typos, to evaluate state-of-the-art Q&A robots. Experiment results show that these common inputs can successfully make these Q&A robots malfunction, DoS, or speaking dirty words. Besides, we propose possible countermeasures to these robustness problems. However, our analysis show that there are no satisfactory solutions currently, which requires further joint efforts by both researchers and industry engineers.

# References

1. Google assistant. https://assistant.google.com
2. Microsoft Cortana personal assistant. https://www.microsoft.com/en-us/cortana
3. Apple Siri personal assistant. https://www.apple.com/ios/siri
4. Amazon Alexa. http://alexa.amazon.com
5. IBM Watson. http://www.ibm.com/watson
6. Baidu DuerOS. https://dueros.baidu.com
7. Alibaba Ali Xiaomi. http://www.alixiaomi.com
8. JD Instant Messaging Intelligence. http://open.jimi.jd.com
9. Auer, S., Bizer, C., Kobilarov, G., Lehmann, J., Cyganiak, R., Ives, Z.: DBpedia: a nucleus for a web of open data. In: Aberer, K., et al. (eds.) ASWC/ISWC -2007. LNCS, vol. 4825, pp. 722–735. Springer, Heidelberg (2007). https://doi.org/10.1007/978-3-540-76298-0_52
10. Bollacker, K., Evans, C., Paritosh, P., Sturge, T., Taylor, J.: Freebase: a collaboratively created graph database for structuring human knowledge. In: Proceedings of the ACM SIGMOD International Conference on Management of Data, pp. 1247–1250 (2008)
11. Suchanek, F.M., Kasneci, G., Weikum, G.: YAGO: a core of semantic knowledge. In: Proceedings of the 16th International Conference on World Wide Web, pp. 697–706 (2007)
12. West, R., Gabrilovich, E., Murphy, K., Sun, S., Gupta, R., Lin, D.: Knowledge base completion via search-based question answering. In: Proceedings of the 23rd International Conference on World Wide Web, pp. 515–526 (2014)
13. Androutsopoulos, I., Ritchie, G.D., Thanisch, P.: Natural language interfaces to databases-an introduction. Nat. Lang. Eng. **1**(1), 29–81 (1995)
14. Liu, X., Long, F.: A survey of multi-modal question answering systems for robotics. In: Proceedings of the 2nd International Conference on Advanced Robotics and Mechatronics (ICARM), pp. 189–194 (2017)
15. Zettlemoyer, L.S., Collins, M.: Learning to map sentences to logical form: structured classification with probabilistic categorial grammars. arXiv preprint arXiv:1207.1420 (2012)
16. Zelle, J.M., Mooney, R.J.: Learning to parse database queries using inductive logic programming. In: Proceedings of the National Conference on Artificial Intelligence, pp. 1050–1055 (1996)
17. Wong, Y.W., Mooney, R.: Learning synchronous grammars for semantic parsing with lambda calculus. In: Proceedings of the 45th Annual Meeting of the Association of Computational Linguistics, pp. 960–967 (2007)
18. Lu, W., Ng, H.T., Lee, W.S., Zettlemoyer, L.S.: A generative model for parsing natural language to meaning representations. In: Proceedings of the Conference on Empirical Methods in Natural Language Processing, pp. 783–792 (2008)
19. Yih, W.T., He, X., Meek, C.: Semantic parsing for single-relation question answering. In: Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics, pp. 643–648 (2014)
20. Yih, W.T., Chang, M.W., He, X., Gao, J.: Semantic parsing via staged query graph generation: question answering with knowledge base. In: Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing, pp. 1321–1331 (2015)
21. Jurafsky, D.: Speech & Language Processing. Pearson Education India (2000)

22. Ravichandran, D., Hovy, E.: Learning surface text patterns for a question answering system. In: Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, pp. 41–47 (2002)
23. Brill, E., Dumais, S., Banko, M.: An analysis of the AskMSR question-answering system. In: Proceedings of the ACL Conference on Empirical Methods in Natural Language Processing, pp. 257–264 (2002)
24. Lin, J.: An exploration of the principles underlying redundancy-based factoid question answering. ACM Trans. Inf. Syst. (TOIS) **25**(2), 6 (2007)
25. Zigong environmental protection bureau insulted a reporter. http://baijiahao.baidu.com/s?id=1603776993370451221
26. Wang, Y., Deng, L., Acero, A.: Semantic frame-based spoken language understanding. In: Spoken Language Understanding: Systems for Extracting Semantic Information from Speech, pp. 41–91 (2011)
27. Li, P., et al.: Dataset and neural recurrent sequence labeling model for open-domain factoid question answering. arXiv preprint arXiv:1607.06275 (2016)
28. Mesnil, G., et al.: Using recurrent neural networks for slot filling in spoken language understanding. IEEE/ACM Trans. Audio Speech Lang. Process. **23**(3), 530–539 (2014)
29. Liu, B., Lane, I.: Recurrent neural network structured output prediction for spoken language understanding. In: Proceedings of the NIPS Workshop on Machine Learning for Spoken Language Understanding and Interactions, pp. 1–7 (2015)
30. Yih, W.T., Richardson, M., Meek, C., Chang, M.W., Suh, J.: The value of semantic parse labeling for knowledge base question answering. In: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, vol. 2, pp. 201–206 (2016)
31. Baudiš, P., Šedivý, J.: Modeling of the question answering task in the YodaQA system. In: Mothe, J., et al. (eds.) CLEF 2015. LNCS, vol. 9283, pp. 222–228. Springer, Cham (2015). https://doi.org/10.1007/978-3-319-24027-5_20
32. Miller, A., Fisch, A., Dodge, J., Karimi, A.H., Bordes, A., Weston, J.: Key-value memory networks for directly reading documents. In: Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, pp. 1400–1409 (2016)
33. Zo. https://www.zo.ai