

ClarQ: A large-scale and diverse dataset for Clarification Question Generation

Vaibhav Kumar

Language Technologies Institute
Carnegie Mellon University
vaibhav2@cs.cmu.edu

Alan W Black

Language Technologies Institute
Carnegie Mellon University
awb@cs.cmu.edu

Abstract

Question answering and conversational systems are often baffled and need help clarifying certain ambiguities. However, limitations of existing datasets hinder the development of large-scale models capable of generating and utilising clarification questions. In order to overcome these limitations, we devise a novel bootstrapping framework (based on self-supervision) that assists in the creation of a diverse, large-scale dataset of clarification questions based on post-comment tuples extracted from stackexchange. The framework utilises a neural network based architecture for classifying clarification questions. It is a two-step method where the first aims to increase the precision of the classifier and second aims to increase its recall. We quantitatively demonstrate the utility of the newly created dataset by applying it to the downstream task of question-answering. The final dataset, ClarQ, consists of ~ 2 M examples distributed across 173 domains of stackexchange. We release this dataset¹ in order to foster research into the field of clarification question generation with the larger goal of enhancing dialog and question answering systems.

1 Introduction

The ubiquitous nature of conversations has led to the identification of various interesting problems (Clark et al., 2019). One of these problems is the ability of a system to ask for clarifications (Rao and Daumé III, 2018; Aliannejadi et al., 2019) to a natural language question.

A user’s complex information need is often lost due to the brevity of the posed question. This leads to an under-specified question containing information gaps which lowers the probability of providing the correct answer. Thus, it would be an improvement if a conversational or a question answering system had a mechanism for refining

user questions with follow-ups (De Boni and Manandhar, 2003). In literature, such questions have been termed *Clarification Questions* (De Boni and Manandhar, 2003; Rao and Daumé III, 2018, 2019).

In the domain of question-answering, the major advantages of a clarification question are its ability to resolve ambiguities (Wang et al., 2018; Aliannejadi et al., 2019) and to improve the probability of finding the most relevant answer. For conversational systems, asking such questions help in driving the conversation deeper along with better engagement of the user (Li et al., 2016; Yu et al., 2016).

Recently, Rao and Daumé III (2018, 2019) have provided a dataset based on stackexchange and used it for clarification question retrieval as well as generation. They also modify a dataset based on Amazon Question-Answering and Product Reviews (McAuley et al., 2015; McAuley and Yang, 2016) to make it suitable for the same task. On the other hand, Aliannejadi et al. (2019) created a dataset (Qulac) built on top of TREC web collections.

However, there are several shortcomings to these datasets, which limit the development of generalizable and large-scale models aimed to tackle the problem of clarification question generation. The stackexchange dataset (Rao and Daumé III, 2018) is created by utilising simple heuristics. This adds a lot of noise, thereby reducing the number of actual clarification questions. It also limits the inclusion of diverse types of questions as it is collected from three similar domains (askubuntu, superuser and unix). The question generation model of Rao and Daumé III (2019) achieves a very low BLEU score when trained on this dataset. On the other hand, the dataset based on Amazon reviews is a poor proxy for clarification questions because product descriptions are not actual questions demanding an answer and

¹<https://github.com/vaibhav4595/ClarQ>

there is no information gap that needs to be addressed.

To overcome the shortcomings of existing datasets, we devise a novel bootstrapping framework based on self-supervision to obtain a dataset of clarification questions from various domains of stackexchange. The framework utilises a neural network based architecture to classify clarification questions. In a two step procedure, the framework first increases the precision of the classifier and then increases its recall. The first step is called down-sampling, where the classifier is iteratively trained on the most confident predictions (carried forward over from the previous iteration). The second step is the up-sampling procedure, where the classifier is iteratively trained by successively adding more positively classified examples. This step provides a boost in recall while restricting the drop in precision to a minimum. The classifier trained on the final iteration is then used for identification of clarification questions.

The overall process ensures that the final dataset is less noisy and, at the same time, consists of a large and diverse number of examples. We must emphasize that, given the large amount of data available on stackexchange, a classifier with moderate recall still serves our purpose. However, it is imperative that precision of the classifier be reasonably high.

2 Methodology

Stackexchange is a network of online question answering websites. On these websites, users may comment on the original post with content such as third party URLs, clarifying questions, etc. We only want to select comments which act as clarifying questions and remove the rest as noise. To this end, we devise a bootstrapping framework for training a classifier capable of identifying clarifying questions.

The bootstrapping method utilises a neural network based classifier \mathcal{L} which is posed with the task of clarification question detection. Formally, given a tuple (p, q) , where $p \in P$ is a post and $q \in q_p$ is a comment made on p , the task is to predict whether q is an actual clarification question for p . This makes it a binary classification problem, where a label 1 indicates q being an actual clarification question and 0 indicates otherwise.

2.1 Data Collection

We first utilise the stackexchange data dump available at <https://archive.org/details/stackexchange>. We extract the posts and the comments made by users on those posts from 173 different domains. We remove all posts which did not have a provided answer. The comments made on the posts act as a potential candidate for clarifying question. This leads to 6,186,934 tuples of (p, q) .

2.2 Bootstrapping

First, we initialise a seed dataset that is used to train \mathcal{L} using the process of iterative refinement as described later. Iterative-refinement itself is subdivided into two parts: (1) Down-Sampling (2) Up-Sampling.

2.2.1 Classifier \mathcal{L}

We utilise a neural network based architecture for the classifier \mathcal{L} . Inspired by [Lowe et al. \(2015\)](#), \mathcal{L} utilises a dual encoder mechanism i.e it uses two separate LSTMs ([Hochreiter and Schmidhuber, 1997](#)) for encoding a post p and a question q . The dual encoder generates hidden representations h_p and h_q for p and q respectively. The resulting element-wise product of h_p and h_q is further passed on to fully connected layers before making predictions via softmax. More formally, the entire process can be summarised as follows:

$$h_p = LSTM_P(p) \quad (1)$$

$$h_q = LSTM_Q(q) \quad (2)$$

$$h_{pq} = \phi(h_p \odot h_q) \quad (3)$$

$$\hat{y} = Softmax(h_{pq}) \quad (4)$$

where, \odot represents the element-wise product, ϕ represents the non-linearity introduced by the fully connected layers and ψ represents the final classification layer.

2.2.2 Seed Selection

In order to select seeds for the bootstrapping procedure, we consider all the collected posts but only use the last comment made on these posts as clarifying questions. We make the assumption that the comments act as a proxy for a clarification question. Later, we remove all (p, q) tuples where q does not have a question mark. Intuitively, the last comment can be a better signal for identifying

clarifying questions as it has more chances of capitalizing the requirements of the original post. It can also be more opinionated than others. We then randomly sample a question from the same domain as that of the post and treat it as an instance of a negative clarification question. Thus each question gets paired with a positive and a negative clarification question. We denote this seed dataset as D_0 .

Algorithm 1 Iterative Refinement

```

1:  $N \leftarrow 5$ 
2:  $D_0 \leftarrow$  Seed Data
3:  $\mathcal{T} \leftarrow$  Annotated Ground Truth
4: for  $i = 1, 2, \dots, N$  do  $\triangleright$  down-sampling
5:    $\mathcal{L} \leftarrow$  Classifier
6:   train  $\mathcal{L}$  on  $D_{i-1}$ 
7:    $D_{temp} \leftarrow []$ 
8:   for  $(p, q) \in D_{i-1}$  do
9:      $y \leftarrow \mathcal{L}(p, q)$ 
10:    if  $y$  is true positive then
11:      add  $(p, q)$  to  $D_{temp}$ 
12:    end if
13:  end for
14:  Sort  $D_{temp}$  using prediction confidence
15:   $D_i \leftarrow$  top 40% of  $D_{temp}$ 
16:  Randomly sample Negatives for  $D_i$ 
17: end for
18:  $S_N \leftarrow D_N$ 
19: for  $i = N, N - 1, \dots, 0$  do  $\triangleright$  up-sampling
20:    $\mathcal{L} \leftarrow$  Classifier
21:   train  $\mathcal{L}$  on  $S_N$ 
22:    $S_{temp} \leftarrow []$ 
23:   for  $(p, q) \in D_{i-1}$  do
24:      $y \leftarrow \mathcal{L}(p, q)$ 
25:     if  $y$  is true positive then
26:       add  $(p, q)$  to  $S_{temp}$ 
27:     end if
28:   end for
29:    $S_{i-1} \leftarrow S_{temp}$ 
30:   Randomly Sample Negatives for  $S_{i-1}$ 
31: end for
32:  $\mathcal{L}_{best} \leftarrow$  Classifier
33:  $\mathcal{L}_{best}$  on  $S_0$ 
34: Use  $\mathcal{L}_{best}$  to classify remaining data

```

2.2.3 Iterative Refinement

The procedure is described in Algorithm 1. This entire process can be segmented into two parts.

Down-Sampling: The aim of this step is to increase the precision of the classifier. In the first

iteration of this step, the classifier \mathcal{L} is trained on the seed dataset D_0 . After training is complete, \mathcal{L} classifies D_0 and the most confident 40% of the positives are selected to train \mathcal{L} in the next iteration. This process is continued for N iterations. Each iteration leads to a new dataset D_i (which is smaller in size than D_{i-1}). Intuitively, the precision of \mathcal{L} on the task of selecting actual clarification question should increase at the end of each iteration as it is successively trained only on the examples which it was more confident about in the previous round.

Up-Sampling: This step is intended to improve the recall of \mathcal{L} while restricting the loss of precision to a minimum. In the first iteration, \mathcal{L} is trained on $S_N = D_N$ i.e the data obtained at the last iteration of the down-sampling procedure. After training is complete, \mathcal{L} is used for classifying D_{N-1} (which is obtained during the second-list iteration of the down-sampling process). The tuples which get classified as positive are used for training \mathcal{L} in the next round. This process continued for N iterations. Note that this procedure has two major differences to the iterative procedure of the down-sampling process. First, instead of using \mathcal{L} for classifying the same dataset which it was trained on, it is used for classifying an up-sampled version of the current dataset. Second, it relaxes the condition of selecting 40% of the most confident examples. Intuitively, this relaxation should help in increasing the recall of the classifier and at the same time should not drastically hamper the precision (as it operates only on the examples which it classifies as positives).

Note that, in order to provide the classifier with examples of non-clarifying questions, we randomly sample negative examples at the end of each iteration (during both up and down-sampling). This is similar to the way in which the D_0 is created.

2.2.4 Classifying Remaining Data

At the end of the iterative refinement procedure, we obtain a dataset on which \mathcal{L} can achieve a good precision and moderate recall on the task of classifying clarification questions. Thus, \mathcal{L} is finally trained on S_0 and used for classifying the 6,186,934 tuples of (p, q) extracted from stackexchange. We again emphasize that it is more important to obtain better precision, as it reduces the amount of noise added to the dataset. Given that there are a large number of (p, q) tuples, a moder-

Iteration	Precision	Recall	F1
1	0.736	0.601	0.662
2	0.758	0.561	0.645
3	0.771	0.390	0.518
4	0.827	0.286	0.426
5	0.829	0.270	0.407

Table 1: Performance of the classifier on the annotated test set at the end of each iteration of the down-sampling procedure.

Iteration	Precision	Recall	F1
1	0.829	0.270	0.407
2	0.835	0.262	0.434
3	0.800	0.270	0.404
4	0.82	0.344	0.488
5	0.82	0.414	0.550

Table 2: Performance of the classifier on the annotated test set at the end of each iteration of the up-sampling procedure.

ate recall can still ensure the incorporation of large and diverse types of (p, q) tuples.

3 Experimental Results

This section describes the results of the iterative refinement strategy.

Test Set Creation: We first create a manually annotated test set to evaluate the effectiveness of the classifier at each step of the iterative refinement process. For this, we randomly sample 100 (p, q) tuples each from 7 different domains (Apple, cooking, gaming, money, photography, scifi, travel). These questions are either the last, second last or the third last comments of their corresponding posts. The annotated test set has a 7:3 ratio of positives to negatives.

Seed Dataset: It is created based on the method described in Section 2.2.2. It consists of 1,800,000 (p, q) tuples, amongst which 50% are randomly sampled negative instances. The classifier is then iteratively trained based on Algorithm 1.

3.1 Results of Iterative Refinement

The results of the down-sampling and the up-sampling procedure are discussed below:

3.1.1 Down-Sampling

Table 1 describes the performance of the classifier on the annotated test set during the down-sampling process. It can be clearly observed that the precision of the classifier increases with each iteration.

Metric	Without CQ	With CQ
P@1	0.751	0.791
P@2	0.399	0.416
P@3	0.278	0.287
P@4	0.214	0.220
P@5	0.174	0.178
MRR	0.791	0.816

Table 3: Performance on the task of question-answer retrieval. CQ stands for clarification question. P@k represents the precision at the kth position of the ranked list. MRR represents the Mean Reciprocal Rank.

Even though there is a substantial decline in recall, the down-sampling procedure helps in increasing the overall precision.

3.1.2 Up-Sampling

Table 2 describes the performance of the classifier on the annotated test set during the up-sampling process. It can be clearly observed that recall of the classifier increases with each iteration, although the final recall (i.e at iteration 5) is lower than the recall obtained in the first iteration of the down-sampling process. Given that there are a large number of (p, q) tuples, a drop in recall will not hamper the quality nor the diversity of the dataset. At the end of the process, we also observe that there is only a marginal drop in precision. Thus, at the end of the last iteration we are able to obtain a classifier which has a high precision and a reasonable recall.

3.2 Downstream Utility

We evaluate the utility of the clarification question in ClarQ by using it for the task of reranking answers. We first randomly sample 1000 (p, q) tuples from 11 different domains (Apple, askubuntu, biology, cooking, english, gaming, money, puzzling, scifi, travel, unix). Corresponding to each tuple, we randomly sample a list of 99 answers (from the same domain as that of the post) and append the actual answer to this list. We first rerank the answers based on the post alone. Later, we rerank the answers by concatenating the post and the clarifying question. Based on the results from Table 3, we observe that concatenating the clarification question to the post does help in improving the performance. The success of this experiment depicts the usefulness of our created dataset.

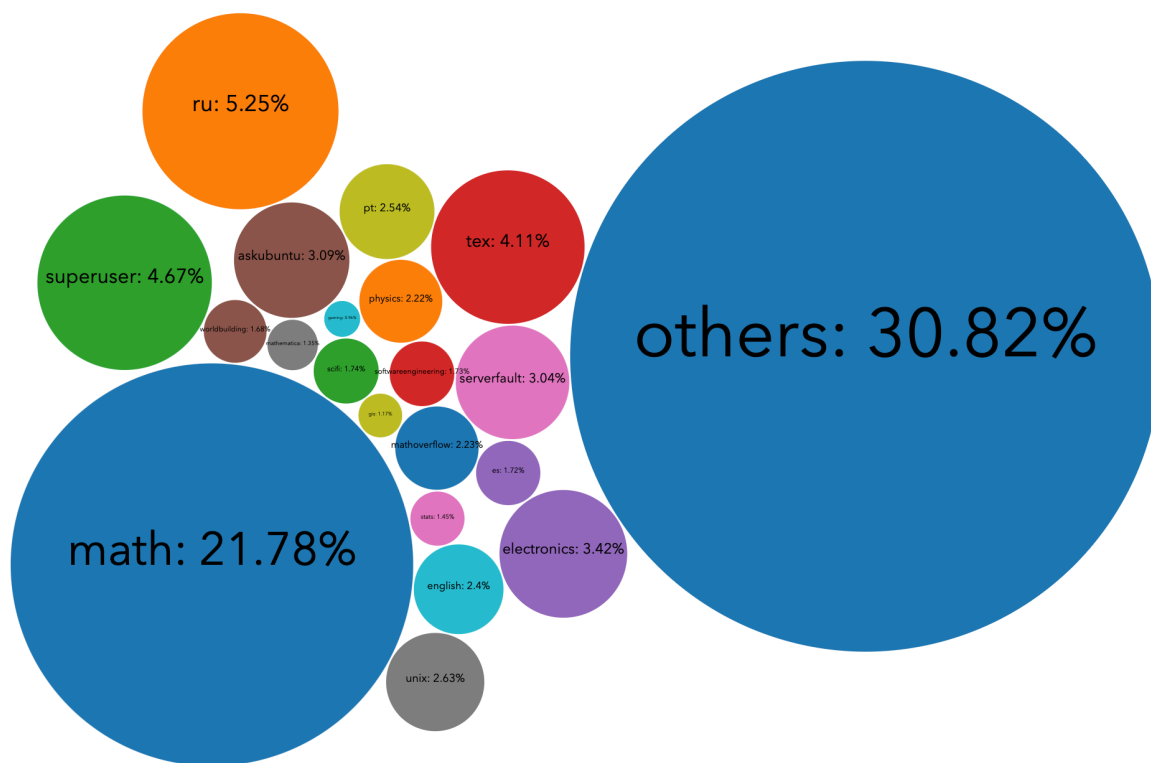


Figure 1: Distribution of the Clarifying Questions across different domains. The figure depicts the top 20 domains. Rest of the domains are clubbed at the end of the spectrum in "others".

4 Dataset Statistics

The classifier obtained at the end of iterative refinement procedure is used for classifying the initially collected (p, q) tuples of 6,186,934. The classifier predicts 2,079,300 tuples as actual clarification questions. As can be seen from Figure 1, these tuples are unequally distributed across 173 different domains. The top 20 domains account for 69.18% of the total (p, q) tuples in the dataset. The remaining 155 domains account for the remaining 30.82% of the total number of tuples.

It is noteworthy that our provided dataset also comprises of actual answers to each post. This would help researchers in evaluating the quality of the clarification questions in a standalone perspective and at the same time with respect to the downstream task of question-answering.

5 Conclusion and Future Work

In this paper, we present a diverse, large-scale dataset (**ClarQ**) for the task of clarification ques-

tion generation. It is created by a two-step iterative bootstrapping framework based on self-supervision. ClarQ consists of $\sim 2M$ post-question tuples spanning 173 different domains. We hope that this dataset will encourage research into clarification question generation and, in the long run, enhance dialog and question-answering systems.

Acknowledgments

We would like to extend our sincere gratitude to Abhimanshu Mishra, Mrinal Dhar and Yash Kumar Lal for helping us understand the structure of the comments and their distribution across domains.

References

Mohammad Aliannejadi, Hamed Zamani, Fabio Crestani, and W Bruce Croft. 2019. Asking clarifying questions in open-domain information-seeking conversations. In *Proceedings of the 42nd International ACM SIGIR Conference on Research and De-*

- velopment in *Information Retrieval*, pages 475–484. ACM.
- Leigh Clark, Nadia Pantidi, Orla Cooney, Philip Doyle, Diego Garaialde, Justin Edwards, Brendan Spillane, Emer Gilmartin, Christine Murad, Cosmin Munteanu, et al. 2019. What makes a good conversation?: Challenges in designing truly conversational agents. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, page 475. ACM.
- Marco De Boni and Suresh Manandhar. 2003. An analysis of clarification dialogue for question answering. In *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology-Volume 1*, pages 48–55. Association for Computational Linguistics.
- Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Jiwei Li, Alexander H Miller, Sumit Chopra, Marc’Aurelio Ranzato, and Jason Weston. 2016. Learning through dialogue interactions by asking questions. *arXiv preprint arXiv:1612.04936*.
- Ryan Lowe, Nissan Pow, Iulian Serban, and Joelle Pineau. 2015. The ubuntu dialogue corpus: A large dataset for research in unstructured multi-turn dialogue systems. *arXiv preprint arXiv:1506.08909*.
- Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton Van Den Hengel. 2015. Image-based recommendations on styles and substitutes. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 43–52. ACM.
- Julian McAuley and Alex Yang. 2016. Addressing complex and subjective product-related queries with customer reviews. In *Proceedings of the 25th International Conference on World Wide Web*, pages 625–635. International World Wide Web Conferences Steering Committee.
- Sudha Rao and Hal Daumé III. 2018. Learning to ask good questions: Ranking clarification questions using neural expected value of perfect information. *arXiv preprint arXiv:1805.04655*.
- Sudha Rao and Hal Daumé III. 2019. Answer-based adversarial training for generating clarification questions. *arXiv preprint arXiv:1904.02281*.
- Yansen Wang, Chenyi Liu, Minlie Huang, and Liqiang Nie. 2018. Learning to ask questions in open-domain conversational systems with typed decoders. *arXiv preprint arXiv:1805.04843*.
- Zhou Yu, Ziyu Xu, Alan W Black, and Alexander Rudnicky. 2016. Strategy and policy learning for non-task-oriented conversational systems. In *Proceedings of the 17th annual meeting of the special interest group on discourse and dialogue*, pages 404–412.