

---

# Reproducibility report of CoSQA: 20,000+ Web Queries for Code Search and Question Answering (2) for ML Reproducibility Challenge 2021

---

Iman Barati\*, Yeganeh Morshedzadeh\*, Ghazaleh Mahmoodi\*  
School of Computer Engineering  
Iran University of Science and Technology, Tehran, Iran  
{iman\_barati, y\_morshedzadeh, gh\_mahmoodi}.iust.ac.ir

## Reproducibility Summary

### Abstract/Summary of the Original Paper (2)

The authors' goal in the paper is to find relevant code pieces that match a user's query. More specifically, they focused on two tasks, code question answering and code search. They gathered and used a new dataset named CoSQA since, for the described problem, a richer and bigger dataset than the previous ones is needed. Additionally, they introduced CoCLR, which incorporates code contrastive learning into the Siamese network with CodeBERT (1). In the code contrastive method, the goal is to learn the word representations in a way that similar objects are as close as possible, whereas the dissimilar objects are kept as far as possible. Therefore, the paper used two augmentation methods, In-Batch Augmentation (IBA) and Query-Rewritten Augmentation (QRA). In conclusion, using these models and datasets achieved better quality in query-code matching tasks.

### Scope of Reproducibility

- The paper's main claim that we tried to reproduce is the Siamese network with CodeBERT (1) and CoCLR models improve overall performance on both tasks (code search and code question answering) by using the CoSQA dataset, particularly on the CodeXGLUE WebQueryTest (3) <sup>2</sup>. Specifically, CoCLR increased the results by 15.6% the WebQueryTest.
- Therefore, in this work, we tried to show:
  1. The performance of the Siamese network with the CodeBERT model on the CodeSearchNet + CoSQA as the training and testing dataset. To do that, we used the same metrics as the paper that are accuracy for the code question answering task and MRR for the code search task.
  2. We performed previous experiments for the CoCLR (Siamese network + CodeBERT using code contrastive learning) model as well.
  3. We trained CoCLR with different variations of augmentations (IBA and QRA) for code search task:
    - (a) in-batch + query-rewritten (delete)
    - (b) in-batch + query-rewritten (copy)
    - (c) in-batch + query-rewritten (switch)
  4. We showed the performance of the CoCLR with various code components on the code search task:
    - (a) complete code
    - (b) header only
    - (c) doc only
    - (d) no header

---

\*All authors contributed equally to this work.

<sup>2</sup><https://github.com/microsoft/CodeXGLUE>

- 31 (e) no doc  
32 (f) no body

## 33 Methodology

- 34 • We used the authors' code provided on [Jun-jie-Huang's GitHub](#).  
35 • Our hardware resource was GPU.RTX2060 Super.xlarge and the specifications are as below:  
36 – RAM: 23.4GB  
37 – VCPUs: 6 VCPU  
38 – Disk: 200GB  
39 • We added ten bash files that include all the experiments done by authors. These can all be run by a single bash  
40 file which simplifies automating the running process significantly.  
• The total budget in terms of GPU hours per task is as below:

Task	Approximate time for each epoch	Number of epochs for training	Total approximate time
Code Question Answering	15 minutes	10	2.5 hours
Code Search	30 minutes	10	5 hour

Table 1: GPU computation time

41

## 42 Results

43 The results are summarized in the tables below:

- 44 • Overall results:
- 45 – Code Question Answering
- 46 \* Test hyperparameters:
- 47 · max\_seq\_length: 200
- 48 · per\_gpu\_eval\_batch\_size: 2
- 49 \* Train hyperparameters:
- 50 · max\_seq\_length: 200
- 51 · per\_gpu\_train\_batch\_size: 8
- 52 · per\_gpu\_eval\_batch\_size: 16
- 53 – Code Search
- 54 \* Test hyperparameters:
- 55 · max\_seq\_length: 200
- 56 · per\_gpu\_retrieval\_batch\_size: 67
- 57 \* Train hyperparameters:
- 58 · max\_seq\_length: 200
- 59 · per\_gpu\_train\_batch\_size: 8
- 60 · per\_gpu\_retrieval\_batch\_size: 67

Model	Code Question Answering (Accuracy)	Code Search (MRR)
CodeBERT (Paper's)	52.87	54.41
CodeBERT + CoCLR (Paper's)	63.38	64.66
CodeBERT (Our's)	46.75	54.60
CodeBERT + CoCLR (Our's)	66.92	64.28

Table 2: Evaluation on two tasks. The data is CodeSearchNet + CoSQA.

- Augmentation results:
  - Test hyperparameters:
    - \* max\_seq\_length: 200
    - \* per\_gpu\_retrieval\_batch\_size: 67
  - Train hyperparameters:
    - \* max\_seq\_length: 200
    - \* per\_gpu\_train\_batch\_size: 8
    - \* per\_gpu\_retrieval\_batch\_size: 67

Augmentations	Train (MRR)	Test (MRR)
Delete (Paper's)	-	63.41
Copy (Paper's)	-	63.97
Switch (Paper's)	-	64.66
Delete (Our's)	62.15	63.87
Copy (Our's)	63.32	64.39
Switch (Our's)	64.34	64.28

Table 3: Performance of CodeBERT with different augmentations (in-batch + query-rewritten) in COCLR on code search.

- Code components results:
  - Test hyperparameters:
    - \* max\_seq\_length: 200
    - \* per\_gpu\_retrieval\_batch\_size: 67
  - Train hyperparameters:
    - \* max\_seq\_length: 160
    - \* per\_gpu\_train\_batch\_size: 8
    - \* per\_gpu\_retrieval\_batch\_size: 67

Augmentations	Our's (MRR)	Paper's (MRR)
Complete code	64.28	64.66
w/o header	60.06	62.01
w/o body	58.31	59.11
w/o documentation	57.57	58.54
w/o header & body	51.67	52.89
w/o header & documentation	42.16	43.35
w/o body & documentation	43.47	42.71

Table 4: Performance of CoCLR-incorporated CodeBERT trained and tested with different code components on code search.

- Reducing the max sequence length resulted in the model's poor performance.
- Changing batch size had little impact on the model's performance, requiring less RAM but increasing the running time. In addition, in the case of the question-answering task, it resulted in worsening the accuracy.

#### What was easy

- Good README file and step-by-step instructions on how to use the code they provided in their [GitHub repository](#).
- Running commands as a bash file which helps automate the testing and training process.
- They used command-line arguments to change the training and testing conditions/inputs, which reduce unintentional mistakes and confusion.
- Providing a log file for each command gives good information about the running process.

#### What was difficult

- Our resource was not as powerful as their resource since we had to change hyper-parameters such as decreasing the batch size in order to train the models. Specifically, we got poorer results for the question-answering task.
- There is not enough instruction to write the command with the correct argument to get the first four rows of results in table 6, which uses the IBA or QRA methods.
- The confusion between the names of the code and their corresponding part in the paper. For example, the term vanilla model was used as the base model, which is equal to the term CodeBERT trained with CodeSearchNet(CSN) + CoSQA.
- Difficulty to download the datasets since they were on google drive. Wget/curl command was not available to download these materials.

#### Conclusion

In conclusion, we achieved quite the same results as the original paper. Our reproducibility study concurs with their main claim in most cases. More precisely, except in one case, that is for the CodeBERT model on code question answering, we obtained poorer results due to modifications to some hyperparameters to make the code executable and model trainable on our hardware. For all the other cases, what we achieved was similar to what was reported in the paper.

## References

- [1] FENG, Z., GUO, D., TANG, D., DUAN, N., FENG, X., GONG, M., SHOU, L., QIN, B., LIU, T., JIANG, D., AND ZHOU, M. Codebert: A pre-trained model for programming and natural languages, 2020.
- [2] HUANG, J., TANG, D., SHOU, L., GONG, M., XU, K., JIANG, D., ZHOU, M., AND DUAN, N. CoSQA: 20,000+ web queries for code search and question answering. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)* (Online, Aug. 2021), Association for Computational Linguistics, pp. 5690–5700.
- [3] LU, S., GUO, D., REN, S., HUANG, J., SVYATKOVSKIY, A., BLANCO, A., CLEMENT, C. B., DRAIN, D., JIANG, D., TANG, D., LI, G., ZHOU, L., SHOU, L., ZHOU, L., TUFANO, M., GONG, M., ZHOU, M., DUAN, N., SUNDARESAN, N., DENG, S. K., FU, S., AND LIU, S. Codexglue: A machine learning benchmark dataset for code understanding and generation. *CoRR abs/2102.04664* (2021).