

---

# Supplementary Material

---

Anonymous Authors

1 **1 Download URL**

2 You can download the datasets from Dropbox Link.

3 **2 Motivation**

4 **1. For what purpose was the dataset created?**

5 These datasets are intended to study the ability of large language models to solve graph-based  
6 tasks using Python API calls. Previous work has not explored using APIs for graph-based  
7 tasks, instead relying on direct responses from large language models.

8 **2. Any other comments?**

9 None.

10 **3 Composition**

11 **1. What do the instances that comprise the dataset represent?**

12 The benchmark dataset’s instance is a test case for LLM to evaluate the ability to use API to  
13 solve graph foundation tasks.

14 The document dataset’s instance is API info.

15 The Code(QA) dataset’s instance is a QA pair for finetuning open-source models.

16 The Doc+Code(QA) dataset’s instance is a QA pair for finetuning open-source models.

17 **2. How many instances are there in total (of each type, if appropriate)?**

18 The benchmark dataset contains 512 instances.

19 The document dataset contains 1413 instances.

20 The Code(QA) dataset contains 29,260 instances.

21 The Doc+Code(QA) dataset contains 29,260 instances.

Table 1: Statistics of datasets.

	<b>Benchmark</b>	<b>Document</b>	<b>Code (QA)</b>	<b>Doc+Code (QA)</b>
Category 1	312	1,115	23,324	23,324
Category 2	154	253	5,136	5,136
Category 3	46	45	800	800
Total	512	1,413	29,260	29,260

22 **3. Does the dataset contain all possible instances or is it a sample (not necessarily random)**  
23 **of instances from a larger set?**

24 It contains all possible instances.

25 **4. What data does each instance consist of?**

26 For the benchmark dataset, each instance includes a human-expert-annotated question, a

27 new role-play rephrased question from GPT-4-turbo-0409, the corresponding answers, the  
 28 reference code for those answers, the problem type (True/False, calculations, draw, multi),  
 29 the key APIs needed to solve the question, the number of lines in the reference code (exclud-  
 30 ing empty lines), the category of the question (basic graph theory, graph statistical learning,  
 31 graph embedding), and the difficulty level of the question.  
 32 For the document dataset, each instance contains the API info for an API.  
 33 For the Code(QA) dataset, each instance contains a question, a reference code corresponding  
 34 to the question, and an API name.  
 35 For the Doc+Code(QA) dataset, each instance contains a question, a reference code corre-  
 36 sponding to the question, and an API name.

37 **5. Is there a label or target associated with each instance? If so, please provide a descrip-**  
 38 **tion.**  
 39 For the benchmark dataset, the targets are the third column in the data.  
 40 For the document dataset, there is no target or label with each instance.  
 41 For the Code(QA) dataset, the targets are the second column in the data.  
 42 For the Doc+Code(QA) dataset, the targets are the second column in the data.

43 **6. Is any information missing from individual instances?**  
 44 None.

45 **7. Are relationships between individual instances made explicit (e.g., users' movie ratings,**  
 46 **social network links)?**  
 47 None.

48 **8. Are there recommended data splits (e.g., training, development/validation, testing)?**  
 49 We expect to use the benchmark dataset for testing, and other 3 datasets for enhancing or  
 50 finetuning.

51 **9. Are there any errors, sources of noise, or redundancies in the dataset?**  
 52 There are some instances that the code can't run properly in Code(QA) dataset and  
 53 Doc+Code(QA) dataset, but these instances can be used for finetuning open-source models to  
 54 enhance the ability to solve graph foundation tasks.

55 **10. Is the dataset self-contained, or does it link to or otherwise rely on external resources**  
 56 **(e.g., websites, tweets, other datasets)?**  
 57 The dataset is self-contained.

58 **11. Does the dataset contain data that might be considered confidential (e.g., data that is**  
 59 **protected by legal privilege or by doctor-patient confidentiality, data that includes the**  
 60 **content of individuals' non-public communications)?**  
 61 No.

62 **12. Does the dataset contain data that, if viewed directly, might be offensive, insulting,**  
 63 **threatening, or might otherwise cause anxiety?**  
 64 No.

65 **13. Does the dataset relate to people?**  
 66 Yes. The role-play rephrased questions contain profession information.

67 **14. Does the dataset identify any subpopulations (e.g., by age, gender)?**  
 68 No.

69 **15. Is it possible to identify individuals (i.e., one or more natural persons), either directly**  
 70 **or indirectly (i.e., in combination with other data) from the dataset?**  
 71 No.

72 **16. Does the dataset contain data that might be considered sensitive in any way (e.g., data**  
 73 **that reveals racial or ethnic origins, sexual orientations, religious beliefs, political**  
 74 **opinions or union memberships, or locations; financial or health data; biometric or**  
 75 **genetic data; forms of government identification, such as social security numbers;**  
 76 **criminal history)?**  
 77 No.

78 17. **Any other comments?**  
79 No.

## 80 **4 Collection Process**

81 You can see the details in the paper.

## 82 **5 Uses**

- 83 1. **Has the dataset been used for any tasks already?**  
84 The dataset has been used for evaluating LLMs' ability to solve graph foundation tasks  
85 using APIs.
- 86 2. **What (other) tasks could the dataset be used for?**  
87 The dataset could possibly be used for other LLM's evaluations and enhancements.
- 88 3. **Are there tasks for which the dataset should not be used?**  
89 None.
- 90 4. **Any other comments?**  
91 None.

## 92 **6 Distribution**

- 93 1. **Will the dataset be distributed to third parties outside of the entity (e.g., company,  
94 institution, organization) on behalf of which the dataset was created?**  
95 The dataset will be open-sourced later, but commercial use is prohibited.
- 96 2. **How will the dataset will be distributed (e.g., tarball on website, API, GitHub)?**  
97 The dataset is free for download at dropbox now, and will be possible at  
98 github.com/huggingface.com later.
- 99 3. **Will the dataset be distributed under a copyright or other intellectual property (IP)  
100 license, and/or under applicable terms of use (ToU)?**  
101 The dataset is licensed under a CC BY-NC 4.0 license.
- 102 4. **Have any third parties imposed IP-based or other restrictions on the data associated  
103 with the instances?**  
104 Not to our knowledge.
- 105 5. **Do any export controls or other regulatory restrictions apply to the dataset or to  
106 individual instances?**  
107 Not to our knowledge.
- 108 6. **Any other comments?**  
109 None.

## 110 **7 Maintenance**

- 111 1. **Who is supporting/hosting/maintaining the dataset?**  
112 All the authors.
- 113 2. **Is there an erratum?**  
114 None.
- 115 3. **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete  
116 instances)?**  
117 Yes. We will update our datasets when we open-source all content on Github.

118  
119  
120  
121  
122  
123

4. **If the dataset relates to people, are there applicable limits on the retention of the data associated with the instances (e.g., were individuals in question told that their data would be retained for a fixed period of time and then deleted)?**  
No.
5. **Any other comments?**  
None.