

---

# Text-space Graph Foundation Models: Comprehensive Benchmarks and New Insights: Supplementary Information

---

Zhikai Chen<sup>1</sup>, Haitao Mao<sup>1</sup>, Jingzhe Liu<sup>1</sup>, Yu Song<sup>1</sup>, Bingheng Li<sup>1</sup>,  
Wei Jin<sup>2</sup>, Bahare Fatemi<sup>3</sup>, Anton Tsitsulin<sup>3</sup>, Bryan Perozzi<sup>3</sup>,  
Hui Liu<sup>1</sup>, Jiliang Tang<sup>1</sup>

<sup>1</sup>Michigan State University, <sup>2</sup>Emory University, <sup>3</sup>Google Research

## 1 Overview of Supplementary Information

In this section, we will first briefly summarize the content included in the supplementary information. The technical appendices in the main text have already covered most of the experimental details, such as the computation environment, hyperparameter settings, and the introduction of the datasets. Therefore, we will mainly provide a more comprehensive datasheet and instructions for code reproduction based on the original appendices.

## 2 Datasheet

### 2.1 Motivation

- For what purpose was the dataset created?** *Our benchmark dataset was created to serve as a foundation for evaluating the effectiveness of text-space graph foundation models (GFMs) across diverse domains of data and tasks.*
- Who created the dataset and on behalf of which entity?** *The dataset was developed by ML researchers listed in the author list.*
- Who funded the creation of the dataset?** *Funding sources of authors will be listed in the acknowledgment section of the main text.*

### 2.2 Distribution

- Will the dataset be distributed to third parties outside of the entity (e.g., company, institution, organization) on behalf of which the dataset was created?** *Yes, the dataset is open to the public.*
- How will the dataset will be distributed (e.g., tarball on website, API, GitHub)?** *The dataset will be distributed through Huggingface <https://huggingface.co/datasets/zkchen/tsgfm>, and the code will be used for developing baseline models through GitHub.*
- Have any third parties imposed IP-based or other restrictions on the data associated with the instances?** *No.*
- Do any export controls or other regulatory restrictions apply to the dataset or to individual instances?** *No.*

## 28 2.3 Maintenance

- 29 1. **Who will be supporting/hosting/maintaining the dataset?** *DSE Lab from Michigan State*  
30 *University will support, host, and maintain the dataset.*
- 31 2. **How can the owner/curator/manager of the dataset be contacted (e.g., email address)?**  
32 *The owner/curator/manager(s) of the dataset can be contacted through the following emails:*  
33 *Zhikai Chen(chenzh85@msu.edu)*
- 34 3. **Is there an erratum?** *No. If errors are found in the future, we will release errata on the*  
35 *project's GitHub page <https://github.com/CurryTang/TSGFM>.*
- 36 4. **Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete**  
37 **instances)?** *Yes, the datasets will be updated as needed to ensure accuracy. Announcements*  
38 *regarding updates will be posted on the project's main webpage [## 47 2.4 Composition](https://github.com/</a></i><br/>39 <i>CurryTang/TSGFM.</i></li><li>40 5. <b>If the dataset relates to people, are there applicable limits on the retention of the data</b><br/>41 <b>associated with the instances (e.g., were the individuals in question told that their data</b><br/>42 <b>would be retained for a fixed period of time and then deleted?)</b> <i>N/A</i></li><li>43 6. <b>Will older version of the dataset continue to be supported/hosted/maintained?</b> <i>Yes,</i><br/>44 <i>older versions of the dataset will continue to be maintained and hosted.</i></li><li>45 7. <b>If others want to extend/augment/build on/contribute to the dataset, is there a mecha-</b><br/>46 <b>nisms for them to do so?</b> <i>Yes. They can submit a pull request on the Github page.</i></li></ol></div><div data-bbox=)*

- 48 1. **What do the instance that comprise the dataset represent (e.g., documents, photos,**  
49 **people, countries?)** *Each instance includes a Pytorch-geometric [1] like data object. It*  
50 *contains vector objects to store the structural and feature information of the graph dataset.*
- 51 2. **How many instances are there in total (of each type, if appropriate)?** *We include 23*  
52 *datasets in this benchmark. The detailed statistics of each one can be found in the main text*  
53 *or the Github page.*
- 54 3. **Does the dataset contain all possible instances or is it a sample of instances from a**  
55 **larger set?** *Yes.*
- 56 4. **Is there a label or target associated with each instance?** *Yes.*
- 57 5. **Is any information missing from individual instances?** *No.*
- 58 6. **Are there recommended data splits (e.g., training, development/validation, testing)?**  
59 *We write a split function in the code of the project. The split is deterministic if the random*  
60 *seed is fixed.*
- 61 7. **Are there any errors, sources of noise, or redundancies in the dataset?** *There may exist*  
62 *potential errors due to the annotation bias in the labeling process of the original dataset.*
- 63 8. **Is the dataset self-contained, or does it link to or otherwise rely on external resources**  
64 **(e.g., websites, tweets, other datasets)?** *The dataset is self-contained.*
- 65 9. **Does the dataset contain data that might be considered confidential?** *No.*
- 66 10. **Does the dataset contain data that, if viewed directly, might be offensive, insulting,**  
67 **threatening, or might otherwise cause anxiety?** *No.*

## 68 2.5 Collection Process

- 69 1. **How was the data associated with each instance acquired?** *The data is collected from the*  
70 *raw version including [2, 3, 4, 5, 6, 7] by either crawling the raw texts or generating the*  
71 *corresponding text features. Detailed data processing functions can be seen on the GitHub*  
72 *page.*
- 73 2. **What mechanisms or procedures were used to collect the data (e.g., hardware apparatus**  
74 **or sensor, manual human curation, software program, software API)?** *Python code.*
- 75 3. **Who was involved in the data collection process (e.g., students, crowdworkers, con-**  
76 **tractors) and how were they compensated (e.g., how much were crowdworkers paid)?**  
77 *Listed authors.*

- 78 4. **Does the dataset relate to people?** *No.*  
79 5. **Did you collect the data from the individuals in questions directly, or obtain it via third**  
80 **parties or other sources (e.g., websites)?** *Third parties.*

## 81 2.6 Uses

- 82 1. **Has the dataset been used for any tasks already?** *No, this dataset has not been used for*  
83 *any tasks yet.*  
84 2. **What (other) tasks could be the dataset be used for?** *These datasets can be used for any*  
85 *tasks related to graph machine learning.*  
86 3. **Is there anything about the composition of the dataset or the way it was collected and**  
87 **preprocessed/cleaned/labeled that might impact future uses?** *No.*  
88 4. **Are there tasks for which the dataset should not be used?** *No.*

## 89 2.7 Accessibility

90 The datasets can be downloaded from <https://huggingface.co/datasets/zkchen/tsgfm> and  
91 the code can be accessed from <https://github.com/CurryTang/TSGFM>. The DOI of our dataset  
92 is 10.57967/hf/2455.

## 93 3 Instructions for reproducibility

94 Instructions for reproducing the benchmark results can be found in the readme of the Github Reposi-  
95 tory <https://github.com/CurryTang/TSGFM>.

## 96 References

- 97 [1] Matthias Fey and Jan Eric Lenssen. Fast graph representation learning with pytorch geometric.  
98 [arXiv preprint arXiv:1903.02428](https://arxiv.org/abs/1903.02428), 2019. Cited on page 2.
- 99 [2] Zhilin Yang, William Cohen, and Ruslan Salakhudinov. Revisiting semi-supervised learning  
100 with graph embeddings. In [International conference on machine learning](https://proceedings.mlr.press/v48/yang16.html), pages 40–48. PMLR,  
101 2016. Cited on page 2.
- 102 [3] Hao Liu, Jiarui Feng, Lecheng Kong, Ningyue Liang, Dacheng Tao, Yixin Chen, and Muhan  
103 Zhang. One for all: Towards training one graph model for all classification tasks. [arXiv preprint](https://arxiv.org/abs/2310.00149)  
104 [arXiv:2310.00149](https://arxiv.org/abs/2310.00149), 2023. Cited on page 2.
- 105 [4] Haiteng Zhao, Shengchao Liu, Chang Ma, Hannan Xu, Jie Fu, Zhi-Hong Deng, Lingpeng Kong,  
106 and Qi Liu. Gimlet: A unified graph-text model for instruction-based molecule zero-shot learning.  
107 [bioRxiv](https://arxiv.org/abs/2305.14920), pages 2023–05, 2023. Cited on page 2.
- 108 [5] Bharath Ramsundar, Peter Eastman, Patrick Walters, Vijay Pande, Karl Leswing, and Zhenqin  
109 Wu. [Deep Learning for the Life Sciences](https://www.amazon.com/Deep-Learning-Life-Sciences-Microscopy/dp/1492039837). O’Reilly Media, 2019. [https://www.amazon.com/](https://www.amazon.com/Deep-Learning-Life-Sciences-Microscopy/dp/1492039837)  
110 [Deep-Learning-Life-Sciences-Microscopy/dp/1492039837](https://www.amazon.com/Deep-Learning-Life-Sciences-Microscopy/dp/1492039837). Cited on page 2.
- 111 [6] Zhikai Chen, Haitao Mao, Hang Li, Wei Jin, Haifang Wen, Xiaochi Wei, Shuaiqiang Wang,  
112 Dawei Yin, Wenqi Fan, Hui Liu, and Jiliang Tang. Exploring the potential of large language  
113 models (llms) in learning on graphs. [ArXiv](https://arxiv.org/abs/2307.03393), abs/2307.03393, 2023. Cited on page 2.
- 114 [7] Hao Yan, Chaozhuo Li, Ruosong Long, Chao Yan, Jianan Zhao, Wenwen Zhuang, Jun Yin,  
115 Peiyan Zhang, Weihao Han, Hao Sun, et al. A comprehensive study on text-attributed graphs:  
116 Benchmarking and rethinking. [Advances in Neural Information Processing Systems](https://arxiv.org/abs/2307.03393), 36:17238–  
117 17264, 2023. Cited on page 2.