
Beyond Text: A Deep Dive into Large Language Models' Ability on Understanding Graph Data

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Large language models (LLMs) have achieved impressive performance on many
2 natural language processing tasks. However, their capabilities on graph-structured
3 data remain relatively unexplored. In this paper, we conduct a series of experiments
4 benchmarking leading LLMs on diverse graph prediction tasks spanning node,
5 edge, and graph levels. We aim to assess whether LLMs can effectively process
6 graph data and leverage topological structures to enhance performance, compared
7 to specialized graph neural networks. Through varied prompt formatting and
8 task/dataset selection, we analyze how well LLMs can interpret and utilize graph
9 structures. By comparing LLMs' performance with specialized graph models,
10 we offer insights into the strengths and limitations of employing LLMs for graph
11 analytics. Our findings provide insights into LLMs' capabilities and suggest
12 avenues for further exploration in applying them to graph analytics.

13 1 Introduction

14 In recent years, there have been unprecedented advancements in large language models (LLMs) [28]
15 such as Transformers [33], BERT [7], GPT [4], and their variants. LLMs can be treated as foundation
16 models that can be readily applied to diverse downstream tasks with little adaptation [4, 16, 19]. These
17 models have achieved state-of-the-art results on many natural language processing tasks including
18 text classification, machine translation, sentiment analysis, and text summarization [42]. Significantly,
19 advancements in architectures and training methodologies have given rise to emergent capabilities,
20 setting state-of-the-art models like GPT-3.5 [4], GPT-4 [26], Claude-2 [2], BARD [10], LLaMA [31],
21 and LLaMA-2 [32] apart from their predecessors. For instance, in-context learning [22] and zero-shot
22 capabilities [16, 35] enable these models to generalize across tasks for which they were not explicitly
23 trained. This is confirmed by their excellent performance in complex activities such as mathematical
24 reasoning and Question Answering (QA) systems.

25 However, most of the tasks that Large Language Models (LLMs) surpassed previous benchmarks
26 are Natural Language Processing (NLP) tasks involving sequential data. Graph-structured data
27 presents additional complexity beyond sequences as it contains rich topological connections between
28 entities that must be modeled along with node, edge, and graph attributes. Graph-structured data is
29 ubiquitous across many domains, including social networks [24], knowledge graphs [27], molecular
30 structures [37], and transportation networks [3]. While LLMs have shown powerful reasoning
31 and generalization capabilities in sequential data, it remains unclear if they can handle structural
32 information beyond context when applied to graph-structured data. This raises a compelling research
33 question: Can the strengths of LLMs be extended to graph-structured data, enabling them to exhibit
34 significant predictive ability? Further, can they compete with state-of-the-art models specialized for
35 graph data, such as Graph Neural Networks (GNNs)?

36 To comprehensively study the capabilities of LLMs on graph-structured data, we conduct a series of
37 empirical experiments with leading LLMs on diverse graph-based tasks that span node-, edge-, and
38 graph-level predictions. By comparing their performance to specialized graph models like GNNs,
39 we aim to assess the potential strengths and limitations of LLMs in this domain. Critically, by
40 altering the input prompt formats, we aim to evaluate how effectively LLMs can extract and leverage
41 the underlying structural information from the graph to enhance their performance in subsequent
42 tasks. Additionally, we explore the importance of the structural data across different task dimensions
43 spanning node, edge, and graph levels as well as diverse dataset domains such as citation networks,
44 social networks, and chemical networks.

45 Broadly, this paper focuses on studying the central question of investigating the capabilities of LLMs
46 on graph-structured data from three perspectives:

- 47 • **Can LLMs effectively process graph analytics tasks even without explicit graph structure?**
48 Given that LLMs have already shown the capability to leverage contextual information for human-
49 like reasoning in many NLP tasks, it becomes intriguing to assess whether they can attain substantial
50 predictive performance on graph data tasks, even in the absence of structural information.
- 51 • **How well can LLMs interpret graph structures to enhance downstream task performance?** It
52 is essential to investigate to what extent LLMs can perceive and interpret important graph structures.
53 Furthermore, it is imperative to understand whether such recognition can influence and enhance
54 performance in subsequent tasks.
- 55 • **How do task dimensions and dataset domains affect LLMs’ ability to handle structured data?**
56 LLMs’ ability in identifying pivotal structural information for predictions can be influenced by
57 specific tasks and data domains. For example, node-level tasks may heavily rely on entity attribute
58 interpretation, while graph-level tasks may demand comprehensive understanding of intricate
59 inter-node interactions. Also, the distinct topologies properties to various dataset domains, whether
60 derived from intricate social networks or sophisticated molecular structures, further influence the
61 proficiency with which LLMs decipher and manage structured data.

62 The subsequent sections of this paper are structured as follows: We initiate with an extensive literature
63 review, highlighting the recent advancements of LLMs within graph domains. Subsequent to this, we
64 present our comprehensive findings on benchmarking LLMs on graph data, aiming to address the
65 aforementioned research questions. This is accompanied by a detailed discussion, delving into the
66 depth of our discoveries across varied experimental setups. We conclude by summarizing the key
67 points and proposing ideas for future explorations.

68 2 Related Works

69 **Large language models for graph-structured data.** In recent literature, a few preliminary stud-
70 ies [40, 5, 36, 11] have made attempts to uncover the potential of LLMs in handling graph-structured
71 data. Unfortunately, a comprehensive examination of LLMs’ capacity to extract and harness crucial
72 topological structures across diverse prompt settings, task levels, and datasets remains underexplored.
73 Both Chen et al.[5] and Guo et al.[11] proposed to apply LLMs directly on graph data. Their research
74 primarily focus on the node classification task, constrained to a selected few datasets within the
75 citation network domain, and thereby fails to offer a thorough exploration of LLMs’ ability over
76 diverse task levels and datasets. In addition, Ye et al.[40] fine-tuned LLMs on a designated dataset to
77 outperform GNN, underscoring a distinct research objective compared to our study which emphasizes
78 the intrinsic proficiency of LLMs in understanding and exploiting graph structures. Meanwhile, Wei
79 et al.[36] treated LLMs as autonomous agents within graph data, which is less relevant to the core
80 focus of our paper.

81 **Graph neural networks.** In recent years, graph neural networks (GNNs) [14, 6, 25, 9, 12, 38, 23,
82 41, 18] have emerged as a powerful deep learning approach for graph analysis and learning. GNNs
83 operate by propagating information along edges of the graph and aggregating neighborhood repre-
84 sentations for each node. The expressive power of GNNs to learn from graph structure makes them
85 well-suited for analyzing complex relational data [38, 43, 20]. Unlike standard deep neural networks
86 which operate on regular grids, GNNs can leverage the topological structure of graphs and have
87 achieved state-of-the-art performance on tasks such as node classification [14], link prediction [17],
88 and graph classification [8].

89 3 Experiments

90 **Datasets.** We conducted the experiments on 5 commonly used graph benchmark datasets for
 91 node-level, edge-level and graph-level tasks: CORA [30], PUBMED [30], OGBN-ARXIV [13],
 92 WORDNET18 [21] and REDDIT [12]. Brief descriptions of the datasets are shown in Table 1.

93 We selected these five datasets for our preliminary experiments due to their rich contextual information
 94 present in the attributes of nodes, edges, and graphs. Specifically, CORA, PUBMED and OGBN-
 95 ARXIV are citation network, where each node represents a research paper while an edge between two
 96 nodes indicates that there is a citation relationship between them. Edge in WORDNET18 links two
 97 synsets that are regarded as nodes. REDDIT came from Reddit posts, in which each node represents a
 98 post and two nodes are connected if the same user comments on two posts. The specifics regarding
 99 their textual features are as follows:

- 100 • CORA: Each node represents a paper in the domain of Artificial Intelligence, containing
 101 the information about its title and abstract. Each paper belongs to one of the following
 102 7 categories: ['Case_Based', 'Theory', 'Genetic_Algorithms', 'Probabilistic_Methods',
 103 'Neural_Networks', 'Rule_Learning', 'Reinforcement_Learning']. An edge from one node
 104 to another indicates the first paper cited the second one.
- 105 • PUBMED: Each node represents a scientific publication from PubMed database pertaining
 106 to diabetes. The node textual information contains keywords from its abstract and text body.
 107 Each paper belongs to one of the following 3 categories: ['Diabetes Mellitus, Experimental',
 108 'Diabetes Mellitus Type 1', 'Diabetes Mellitus Type 2']. An edge from one node to another
 109 indicates the first paper cited the second one.
- 110 • OGBN-ARXIV: Each node represents a research paper, containing the information about
 111 its title and abstract. Each paper belongs to one of 40 categories on arxiv.cs such as 'AI'
 112 (Artificial Intelligence). An edge leading from one node to another signifies that the first
 113 paper cites the second one.
- 114 • WORDNET18: Each node represents a synset, containing a description. An edge between
 115 two nodes indicate their relation such as 'furniture', 'includes', or 'bed'. Each edge belongs
 116 to one of 18 relationships.
- 117 • REDDIT: Each node corresponds to a post made by a user, which contains descriptions or
 118 discussions about a particular topic. Each graph symbolizes a subreddit (or community),
 119 with affiliations to one out of 29,651 distinct communities, for instance, 'math'.

120 **Choices of LLMs.** We opted to utilize OpenAI’s state-of-the-art models, GPT-3.5 (GPT) and
 121 GPT-4, via their API system, based on a balance between performance and cost considerations. We
 122 adopted GPT with the latest versions (*gpt-3.5-turbo-16k* and *gpt-4*) in experiments.

123 **Implementation Details.** For node classification task, we follow the same train-test split of CORA,
 124 PUBMED and OGBN-ARXIV as established in semi-supervised GNN methods [14, 38]. For link
 125 prediction on CORA, PUBMED and WORDNET18, a random 15% of the links from the graph and
 126 the same number of negative-edge node pairs are packed into the test sets. For graph classification,

Dataset	#Node	#Edge	#Task	Metric
CORA	2,708	5,278	7-class node classifi. & Link Prediction	Accuracy
PUBMED	19,717	44,324	3-class node classifi. & Link Prediction	Accuracy
OGBN-ARXIV	169,343	1,166,243	40-class node classifi.	Accuracy
Dataset	#Entity	#Relation	#Task	Metric
WORDNET18	40,943	18	18-class link classifi.	Accuracy
Dataset	#Node	#Subgraph	#Task	Metric
REDDIT	3,848,330	29,651	70-class subgraph classifi.	Accuracy

Table 1: Statistics of the datasets. For REDDIT, it actually contains 29,651 subreddits (classes). Here we only randomly sampled 70 communities for graph classification task in each run.

127 in each run, we randomly selected 70 communities. Experiments conducted on WORDNET18 and
 128 the retrieval test for CORA employed few-shot prompts. Conversely, all other experiments leveraged
 129 zero-shot prompts. We executed each experiment thrice, subsequently averaging the results."

130 **Comparison GNN Methods.** On node-level tasks, we choose the semi-supervised results from
 131 Graph Neural Network (GNN) [29], Graph Convolutional Network (GCN) [15] and Graph Attention
 132 Network (GAT) [34] to compare with performance from LLMs. On edge-level tasks, we consider
 133 Graph Auto-Encoder (GAE) [1], Graph InfoClust (GIC) [39]. It is worth noting this is not an abso-
 134 lutely fair comparison. Since LLMs operate under zero-shot or few-shot settings, where GNNs require
 135 a training set for parameter optimization. Additionally, potential data leakage during the LLMs'
 136 training process remains a concern. However, these studies aim to offer a foundational understanding
 137 of LLMs' proficiency in understanding graph data structures and forecasting downstream tasks.

Model	Node-level		
	CORA	PUBMED	ARXIV
GCN-64*	0.814	0.790	0.731
GAT	0.832	0.790	0.742
GNN	0.829	0.738	0.721
GPT-3.5	0.627	0.673	0.516
GPT-4	0.432	0.821	0.642
GPT-3.5*	0.647	0.712	0.509
GPT-4*	0.531	0.833	0.673
GPT-3.5[⊕]	0.656	0.704	0.445
GPT-4[⊕]	-	0.814	-
GPT-3.5[•]	0.054	-	-
GPT-4[•]	0.047	-	-

138

Table 2: Average accuracy on node classification tasks. **[No structure information]** **GPT-3.5** and **GPT-3.5*** mean zero-shot and few-shot prompt strategy. **[With structure information]** **GPT-3.5[⊕]** and **GPT-3.5[•]** mean prompts contain neighbors' information and retrieval requires, respectively.

Model	Edge-level		
	CORA	PUBMED	WORDNET
GAE	0.793	0.923	-
GIC	0.812	0.775	-
GNN	0.739	0.528	-
GPT-3.5	0.512	0.116	0.134
GPT-4	0.578	0.132	0.169
GPT-3.5[◦]	0.646	0.114	-
GPT-4[◦]	0.967	0.143	-

Table 3: Average accuracy on link prediction tasks. **GPT-3.5[◦]** means adding structural information into prompt like Table 6. There are only triples for entries in WORDNET18, which makes there is no connected graph structure for it.

Model	Number of labels									
	1	5	10	15	20	30	40	50	60	70
GPT-3.5	0.773	0.662	0.618	0.594	0.628	0.604	0.638	0.536	0.618	0.507
GPT-4	0.957	0.886	0.895	0.843	0.795	-	-	-	-	-

Table 4: Average performance of GPT-3.5 and GPT-4 on REDDIT when possible labels increase from 1 to 70. Results on GPT-4 with more labels are not available due to input limit of prompt length.

139 3.1 Node-level task

140 Driven by the goal of investigating LLMs' capabilities in discerning patterns within textual graphs and
 141 leveraging this for downstream tasks, we crafted three distinct prompts for our node-level prediction
 142 task experiments: (1) absence of graph topology descriptions; (2) straightforward presentation of
 143 all neighborhood data to the LLM; and (3) a retrieval-based prompt guiding the LLM to extract
 144 task-centric structural details. Examples of these prompts can be found in Table 6.

145 **LLMs' zero-shot or few-shot ability on node classification tasks is still usually weaker than the**
 146 **semi-supervised GNN performance.** This may arguably suggests that general LLMs still can not
 147 surpass the specialized graph models on node classification task. As indicated in Table 2, GPT-4
 148 outperforms the GNN models in zero-shot and few-shot settings on PUBMED, but this superiority isn't
 149 observed on CORA and OGBN-ARXIV. We hypothesize three potential reasons for this discrepancy:
 150 1. Fewer categories; 2. Less semantic overlap between categories; 3. Questionable groundtruth
 151 categories. GPT's 'wrong' predictions about citation networks are mostly reasonable. Particularly,

Node-level Task	Structure?	Prompt to GPT
Zero-shot (CORA & PUBMED & OGBN-ARXIV)	No	The title of one paper is <Title> and its abstract is <Abstract>. Here are possible categories: <Categories>. Which category does this paper belong to? You can only output one category.
	Yes	The title of one paper is <Title> and its abstract is <Abstract>. This paper cited following papers: <TitleList> and abstracts of these papers are <AbstractList>. Here are possible categories: <Categories>. Which category does this paper belong to? You can only output one category.
	Yes	The title of one paper is <Title> and its abstract is <Abstract>. This paper is cited by following papers: <TitleList1>. Each of these papers belongs to one categories in: <Categories>. You need to 1.Analyse the papers’ topic based on the give title and abstract; 2.Analyse the pattern of citation information based on their titles, and retrieve the citation information you think is important to help you determine the category of the first given paper. Now you need to combine the information from 1 and 2 to predict the category of the first given paper. You should only output one category.
Few-shot	Yes	Here is a list of labeled papers: The title and abstract of the first paper are <Title1> and <Abstract1> respectively, and this paper belongs to <Category1>... Here is a new paper whose title is <Title> and its abstract is <Abstract>. Here are possible categories: <Categories>. Which category does this paper belong to? You can only output one category.

Table 5: Examples of different prompts used in node classification experiments. For few-shot tasks, we randomly sampled two papers for each category due to the limit of input length.

152 papers in OGBN-ARXIV with lable of Computation and Language always are classified into other
153 categories like Artificial Intelligence and Machine Learning. These prediction error papers always
154 mentioned many machine learning concepts in their abstracts. We also argued weather datasets
155 are ‘out of fashion’ compared with the information in the GPT’s corpus. We prompted GPT to use
156 pre-2000 categorization criteria on CORA, but this does not lead to improvements. Intriguingly,
157 GPT-4 did not consistently surpass GPT-3.5 in terms of predictive accuracy, hinting that the extent of
158 pre-trained knowledge could significantly influence predictions in zero-shot scenarios.

159 **Incorporating structural information can slightly enhance the performance of GPT on node-**
160 **level tasks to a certain degree.** As evidenced in Table 2, the inclusion of neighborhood information
161 enhances node classification performance in certain instances. However, this improvement lacks
162 consistency across different LLM selections and datasets. Such observations could suggest that
163 structural information might not be a pivotal factor in node-classification tasks. Additionally, we
164 assessed the capability of GPT to retrieve information by incorporating retrieval requirements into
165 the prompts for CORA. Regrettably, this led to both GPT-3.5 and GPT-4 struggling significantly,
166 rendering them largely unable to provide accurate predictions.

167 3.2 Edge-level task

168 **Contrary to node-level tasks, the structural information of a graph seems to be more crucial for**
169 **edge-level tasks.** When only node data, excluding neighborhood information, is presented to GPT, the
170 link prediction accuracy on CORA stands at 51.2% for GPT-3.5 and 57.8% for GPT-4. Remarkably,
171 these figures significantly increased to 64.3% and 96.7% respectively when we incorporate the nodes’
172 neighbors information. Notably, the zero-shot result of GPT-4 even surpass the performances of all
173 trained GNN models. It is worth noting that some wrongly predicted links can be attributed to the
174 absence of neighbor information for these nodes in the dataset. Table 6 illustrates the difference
175 between prompts used on link prediction tasks. It is also interesting that when we further increase the
176 information of neighboring nodes (e.g., the abstract of an article), the prediction accuracy becomes
177 worse than only with information of titles. For the link prediction task on WORDNET18, we randomly
178 selected 5 entries for each relationship and presented this information to GPT. Unfortunately, both
179 GPT-3.5 and GPT-4 struggled to achieve a high predictive accuracy based on the provided data. A

Edge-level Task	Structure?	Prompt to GPT
Zero-shot (CORA & PUBMED)	No	There are two papers. Title of the first paper is <code><Title1></code> and its abstract is <code><Abstract1></code> . Title of the second paper is <code><Title2></code> , and its abstract is <code><Abstract2></code> . You need to predict whether the second paper or not. Answer 'YES' or 'NO'.
	Yes	There are two papers. Title of the first paper is <code><Title1></code> and its abstract is <code><Abstract1></code> . Title of the second paper is <code><Title2></code> , and its abstract is <code><Abstract2></code> . The first paper cited following papers: <code><Titles></code> . You need to predict whether the second paper or not. Answer 'YES' or 'NO'.
Few-shot (WORDNET18)	No	We define two descriptions should have a relationship, such as furniture <code><includes></code> bed. There are some samples: <code><Entries></code> . Here are possible relations: <code><Relationships></code> . The first entry is <code><Entry1></code> and the second entry is <code><Entry2></code> . You must output only one relationship between these two entries.

Table 6: Examples of different prompts used in link prediction experiments. Structural information plays an important role in link prediction task.

180 plausible explanation for this could be GPT’s heavy reliance on its pre-trained knowledge, especially
181 when not fine-tuned for specific tasks.

182 3.3 Graph-level task

183 For graph-level tasks, we only conducted experiments on REDDIT due to its text richness and semantic
184 ambiguity. Only GPT-3.5 was tested since the information of one community is large even we have
185 summarized the information from each user. We selected *top-k* post summaries of the most replied
186 users as representative information of a community. As shown in Table 4, when GPT needs to make
187 predictions from full 70 communities, the accuracy was 50.7%. The accuracy decreased from 77.3%
188 to 50.7% when possible communities in the `<SubReddits>` list increased from 1 to 70.

Graph-level Task	Structure?	Prompt to GPT
Zero-shot (REDDIT)	Yes	There are texts from representative users of one Reddit community: <code><Posts></code> . There are following communities: <code><SubReddits></code> . Which community does these texts belong to? You should only output one community from given communities.

Table 7: Example prompt used in graph classification experiments. Structural information is given by a list of *top-k* important nodes a graph.

189 4 Conclusion and Future Works

190 This research presented a systematic empirical evaluation of leading LLMs on diverse graph learning
191 tasks spanning node, edge, and graph levels. Through careful variation of prompt design and dataset
192 selection, we assessed the capabilities of models such as GPT-3.5 and GPT-4 in interpreting and
193 leveraging graph structural information to enhance predictive performance. Our results demonstrate
194 that while LLMs exhibit reasonable node classification capabilities even without explicit graph data,
195 likely by relying on contextual clues, their zero-shot performance continues to lag behind state-of-
196 the-art GNNs specialized for this domain. However, incorporating graph topology information can
197 significantly boost performance on edge-level link prediction tasks, with GPT-4 even surpassing
198 certain GNNs in select cases. On more complex graph classification tasks, limitations emerge
199 in handling increased output complexity. In summary, this research provides valuable evidence
200 that LLMs have promising capabilities on graph analytics, while also revealing clear areas for
201 improvement compared to specialized graph models.

202 Our future work should explore more rigorous benchmarking LLMs on graph learning tasks with
203 graph specialized models, novel prompt designs to focus on topological structures, evaluating on
204 additional graph tasks, and even fine-tuning open-sourced LLMs on graphs. By exploring these
205 avenues, the full potential of large language models for advancing graph representation learning and
206 analytics can be more promising.

207 **References**

- 208 [1] Seong Jin Ahn and MyoungHo Kim. Variational graph normalized autoencoders. In *Proceedings*
209 *of the 30th ACM international conference on information & knowledge management*, pages
210 2827–2831, 2021.
- 211 [2] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn
212 Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, et al. Training a helpful and harmless
213 assistant with reinforcement learning from human feedback. *arXiv preprint arXiv:2204.05862*,
214 2022.
- 215 [3] Jayanth R Banavar, Amos Maritan, and Andrea Rinaldo. Size and form in efficient transportation
216 networks. *Nature*, 399(6732):130–132, 1999.
- 217 [4] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhari-
218 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal,
219 Ariel Herbert-Voss, Gretchen Krueger, T. J. Henighan, Rewon Child, Aditya Ramesh, Daniel M.
220 Ziegler, Jeff Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin,
221 Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Rad-
222 ford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. *ArXiv*,
223 abs/2005.14165, 2020.
- 224 [5] Zhikai Chen, Haitao Mao, Hang Li, Wei Jin, Hongzhi Wen, Xiaochi Wei, Shuaiqiang Wang,
225 Dawei Yin, Wenqi Fan, Hui Liu, et al. Exploring the potential of large language models (llms)
226 in learning on graphs. *arXiv preprint arXiv:2307.03393*, 2023.
- 227 [6] Michaël Defferrard, Xavier Bresson, and Pierre Vandergheynst. Convolutional neural networks
228 on graphs with fast localized spectral filtering. *arXiv preprint arXiv:1606.09375*, 2016.
- 229 [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of
230 deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Confer-*
231 *ence of the North American Chapter of the Association for Computational Linguistics: Human*
232 *Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis,
233 Minnesota, June 2019. Association for Computational Linguistics.
- 234 [8] Federico Errica, Marco Podda, Davide Bacciu, and Alessio Micheli. A fair comparison of graph
235 neural networks for graph classification. *arXiv preprint arXiv:1912.09893*, 2019.
- 236 [9] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural
237 message passing for quantum chemistry. In *International Conference on Machine Learning*,
238 pages 1263–1272. PMLR, 2017.
- 239 [10] Google AI. BARD: Deepmind’s conversational ai assistant. [https://blog.google/
240 products/search/introducing-bard-google-ai](https://blog.google/products/search/introducing-bard-google-ai), 2022. Accessed: October 1, 2023.
- 241 [11] Jiayan Guo, Lun Du, and Hengyu Liu. Gpt4graph: Can large language models understand graph
242 structured data? an empirical evaluation and benchmarking. *arXiv preprint arXiv:2305.15066*,
243 2023.
- 244 [12] William L Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large
245 graphs. *arXiv preprint arXiv:1706.02216*, 2017.
- 246 [13] Weihua Hu, Matthias Fey, Marinka Zitnik, Yuxiao Dong, Hongyu Ren, Bowen Liu, Michele
247 Catasta, and Jure Leskovec. Open graph benchmark: Datasets for machine learning on graphs.
248 *Advances in neural information processing systems*, 33:22118–22133, 2020.
- 249 [14] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional
250 networks. *arXiv preprint arXiv:1609.02907*, 2016.
- 251 [15] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional
252 networks. In *International Conference on Learning Representations*, 2017.
- 253 [16] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large
254 language models are zero-shot reasoners. *Advances in neural information processing systems*,
255 35:22199–22213, 2022.

- 256 [17] Ajay Kumar, Shashank Sheshar Singh, Kuldeep Singh, and Bhaskar Biswas. Link prediction
257 techniques, applications, and performance: A survey. *Physica A: Statistical Mechanics and its*
258 *Applications*, 553:124289, 2020.
- 259 [18] Chen Ling, Junji Jiang, Junxiang Wang, My T Thai, Renhao Xue, James Song, Meikang
260 Qiu, and Liang Zhao. Deep graph representation learning and optimization for influence
261 maximization. In *International Conference on Machine Learning*, pages 21350–21361. PMLR,
262 2023.
- 263 [19] Chen Ling, Xujiang Zhao, Jiaying Lu, Chengyuan Deng, Can Zheng, Junxiang Wang, Tanmoy
264 Chowdhury, Yun Li, Hejie Cui, et al. Domain specialization as the key to make large language
265 models disruptive: A comprehensive survey. *arXiv preprint arXiv:2305.18703*, 2023.
- 266 [20] Meng Liu, Hongyang Gao, and Shuiwang Ji. Towards deeper graph neural networks. In
267 *Proceedings of the 26th ACM SIGKDD international conference on knowledge discovery &*
268 *data mining*, pages 338–348, 2020.
- 269 [21] George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*,
270 38(11):39–41, 1995.
- 271 [22] Sewon Min, Mike Lewis, Luke Zettlemoyer, and Hannaneh Hajishirzi. Metaicl: Learning to
272 learn in context. *arXiv preprint arXiv:2110.15943*, 2021.
- 273 [23] Christopher Morris, Martin Ritzert, Matthias Fey, William L Hamilton, Jan Eric Lenssen,
274 Gaurav Rattan, and Martin Grohe. Weisfeiler and leman go neural: Higher-order graph neural
275 networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages
276 4602–4609, 2019.
- 277 [24] Mark EJ Newman, Duncan J Watts, and Steven H Strogatz. Random graph models of social
278 networks. *Proceedings of the national academy of sciences*, 99(suppl_1):2566–2572, 2002.
- 279 [25] Mathias Niepert, Mohamed Ahmed, and Konstantin Kutzkov. Learning convolutional neural
280 networks for graphs. In *International conference on machine learning*, pages 2014–2023.
281 PMLR, 2016.
- 282 [26] OpenAI. Gpt-4 technical report. *ArXiv*, abs/2303.08774, 2023.
- 283 [27] Heiko Paulheim. Knowledge graph refinement: A survey of approaches and evaluation methods.
284 *Semantic web*, 8(3):489–508, 2017.
- 285 [28] Xipeng Qiu, Tianxiang Sun, Yige Xu, Yunfan Shao, Ning Dai, and Xuanjing Huang. Pre-trained
286 models for natural language processing: A survey. *Science China Technological Sciences*,
287 63(10):1872–1897, 2020.
- 288 [29] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini.
289 The graph neural network model. *IEEE transactions on neural networks*, 20(1):61–80, 2008.
- 290 [30] Prithviraj Sen, Galileo Mark Namata, Mustafa Bilgic, Lise Getoor, Brian Gallagher, and Tina
291 Eliassi-Rad. Collective classification in network data. *AI Magazine*, 29(3):93–106, 2008.
- 292 [31] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timo-
293 thé Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez,
294 Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation
295 language models. *ArXiv*, abs/2302.13971, 2023.
- 296 [32] Hugo Touvron, Louis Martin, Kevin R. Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei,
297 Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Daniel M. Bikel, Lukas
298 Blecher, Cristian Cantón Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes,
299 Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony S.
300 Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian
301 Khabsa, Isabel M. Kloumann, A. V. Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut
302 Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov,
303 Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta,

- 304 Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, R. Subramanian, Xia Tan, Binh
305 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zhengxu Yan, Iliyan Zarov,
306 Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert
307 Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat
308 models. *ArXiv*, abs/2307.09288, 2023.
- 309 [33] Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N.
310 Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017.
- 311 [34] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua
312 Bengio. Graph attention networks. In *International Conference on Learning Representations*,
313 2018.
- 314 [35] Jason Wei, Maarten Bosma, Vincent Y Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan
315 Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. *arXiv*
316 *preprint arXiv:2109.01652*, 2021.
- 317 [36] Lanning Wei, Zhiqiang He, Huan Zhao, and Quanming Yao. Unleashing the power of graph
318 learning through llm-based autonomous agents. *arXiv preprint arXiv:2309.04565*, 2023.
- 319 [37] Zhenqin Wu, Bharath Ramsundar, Evan N Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S
320 Pappu, Karl Leswing, and Vijay Pande. Moleculenet: a benchmark for molecular machine
321 learning. *Chemical science*, 9(2):513–530, 2018.
- 322 [38] Keyulu Xu, Weihua Hu, Jure Leskovec, and Stefanie Jegelka. How powerful are graph neural
323 networks? *arXiv preprint arXiv:1810.00826*, 2018.
- 324 [39] Hong Yang, Shirui Pan, Peng Zhang, Ling Chen, Defu Lian, and Chengqi Zhang. Binarized
325 attributed network embedding. In *2018 IEEE International Conference on Data Mining (ICDM)*,
326 pages 1476–1481. IEEE, 2018.
- 327 [40] Ruosong Ye, Caiqi Zhang, Runhui Wang, Shuyuan Xu, and Yongfeng Zhang. Natural language
328 is all a graph needs. *arXiv preprint arXiv:2308.07134*, 2023.
- 329 [41] Zheng Zhang and Liang Zhao. Representation learning on spatial networks. *Advances in Neural*
330 *Information Processing Systems*, 34:2303–2318, 2021.
- 331 [42] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min,
332 Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv*
333 *preprint arXiv:2303.18223*, 2023.
- 334 [43] Jiong Zhu, Yujun Yan, Lingxiao Zhao, Mark Heimann, Leman Akoglu, and Danai Koutra. Be-
335 yond homophily in graph neural networks: Current limitations and effective designs. *Advances*
336 *in neural information processing systems*, 33:7793–7804, 2020.