# Large Language Models for Cryptocurrency Transaction Analysis: A Bitcoin Case Study

Yuchen Lei and Yuexin Xiang, *IEEE Graduate Student Members*, Qin Wang, *IEEE Member*, Rafael Dowsley, Tsz Hon Yuen, Kim-Kwang Raymond Choo, *Senior Member, IEEE*, Jiangshan Yu, *IEEE Member*

*Abstract*—Cryptocurrencies are widely used, yet current methods for analyzing transactions often rely on opaque, black-box models. While these models may achieve high performance, their outputs are usually difficult to interpret and adapt, making it challenging to capture nuanced behavioral patterns. Large language models (LLMs) have the potential to address these gaps, but their capabilities in this area remain largely unexplored, particularly in cybercrime detection. In this paper, we test this hypothesis by applying LLMs to real-world cryptocurrency transaction graphs, with a focus on Bitcoin, one of the most studied and widely adopted blockchain networks. We introduce a three-tiered framework to assess LLM capabilities: foundational metrics, characteristic overview, and contextual interpretation. This includes a new, human-readable graph representation format, LLM4TG, and a connectivity-enhanced transaction graph sampling algorithm, CETraS. Together, they significantly reduce token requirements, transforming the analysis of multiple moderately large-scale transaction graphs with LLMs from nearly impossible to feasible under strict token limits. Experimental results demonstrate that LLMs have outstanding performance on foundational metrics and characteristic overview, where the accuracy of recognizing most basic information at the node level exceeds 98.50% and the proportion of obtaining meaningful characteristics reaches 95.00%. Regarding contextual interpretation, LLMs also demonstrate strong performance in classification tasks, even with very limited labeled data, where top-3 accuracy reaches 72.43% with explanations. While the explanations are not always fully accurate, they highlight the strong potential of LLMs in this domain. At the same time, several limitations persist, which we discuss along with directions for future research.

*Index Terms*—LLMs, Transaction Graph, Cybercrime Detection, Graph Representation, Cryptocurrency, Blockchain.

## I. INTRODUCTION

Large language models (LLMs) [1] have significantly boosted the productivity of daily life and have a huge impact on the research community, such as in natural language processing (NLP) [2], [3], computer vision (CV) [4], [5], and application research [6], [7], [8]. The applications of LLMs also extend beyond traditional domains, influencing areas with social and economic implications.

†Yuchen Lei and Yuexin Xiang contributed equally to this work.

‡Corresponding author (✉): Yuexin.Xiang@monash.edu.

Yuchen Lei is with the School of Cyber Science and Engineering, Wuhan University, Wuhan 430072, China.

Yuexin Xiang, Rafael Dowsley, and Tsz Hon Yuen are with the Faculty of Information Technology, Monash University, Clayton, VIC 3800, Australia.

Qin Wang is with CSIRO's Data61, Eveleigh, NSW 2015, Australia.

Kim-Kwang Raymond Choo is with the Department of Information Systems and Cybersecurity, University of Texas at San Antonio, TX 78249-0631, USA.

Jiangshan Yu is with the School of Computer Science, The University of Sydney, Camperdown, NSW 2006, Australia.

One such area is the cryptocurrency ecosystem. Its growing adoption in finance, retail, and entertainment has led to a surge in transaction volumes. However, the expansion also exposes the ecosystem to risks, such as scams and money laundering, enabled by its decentralized and pseudoanonymous nature. Current analysis methods rely on black-box models and struggle with interpretability of results and adaptability. In this context, applying LLMs to analyze cryptocurrency transactions offers a promising approach to bridging these gaps. By leveraging their capacity to interpret complex patterns and behaviors, LLMs can help identify illicit activities and enhance cybercrime detection efforts.

Although LLMs trained on massive datasets excel in NLP tasks, their application to graph analysis presents challenges due to structural differences between graph and text data. Recent studies investigated the possibilities of LLMs for handling graph data-related tasks, concluding affirmatively that they are capable of completing specific tasks with acceptable performance on graphs such as small graphs, citation graphs, or knowledge graphs (KGs) [9], [10], [11], [12]. Nevertheless, measuring LLMs' capability to understand and analyze cryptocurrency transaction graphs remains impractical. They contain different information compared with the other graph types such as KGs. Taking the Bitcoin transaction network as an example, the node represents the Bitcoin address or transaction, while the edge indicates the token flows among the address nodes and the transaction nodes [13].

In addition, due to the input token limit of LLMs, how to efficiently feed larger graph data into LLMs to gain more information for potentially improving the quality of generated answers to various questions relevant to Bitcoin transaction graphs (e.g., address type prediction) continues to be an open question. To bridge the gaps in applying LLMs to transaction graph analysis, we study Bitcoin networks and address three research questions (RQs):

- **RQ1:** What graph representation formats are effective in LLMs for Bitcoin transaction graphs?
- **RQ2:** How to measure LLMs' capacity to understand or analyze Bitcoin transaction graphs?
- **RQ3:** What are the key differences between using engineered graph features and raw graph data in analyses?

We adopt quantitative methods combined with qualitative analysis to answer those research questions. For RQ1, we investigate various graph representation formats and their feasibility for LLMs. To reduce the token consumption of raw graphs, we propose a novel representation format called

LLM4TG based on the characteristics of LLMs. As for RQ2 and RQ3, we propose three levels for measuring the understanding of transaction graph:

- **Level 1 - Foundational Metrics:** LLMs can determine the basic information of the graph such as the in-degree and output token amount of a node.
- **Level 2 - Characteristic Overview:** LLMs can figure out the highlighted characteristics of the graph, e.g., a node has a significantly large out-degree.
- **Level 3 - Contextual Interpretation:** LLMs can classify cryptocurrency address types for addresses without labels based on labeled address samples.

**Contributions.** To the best of our knowledge, this is the first systematic study of LLMs' capabilities in analyzing real-world cryptocurrency transaction graphs. We make the following contributions:

- We present a layered framework with three levels of understanding for measuring LLMs' ability to analyze transaction graphs in cryptocurrency networks.
- We propose a text-based graph representation format, denoted LLM4TG. It reduces redundant data and provides a human-readable syntax that naturally supports processing by LLMs.
- We design a Connectivity-Enhanced Transaction Graph Sampling algorithm, CETraS, for graph summarization, removing less important nodes in moderately large-scale transaction graphs while enhancing critical connections.
- We conduct both quantitative and qualitative evaluations of five representative models from three major LLM families (GPT, DeepSeek, and LLaMA) across multiple tasks, revealing their strengths and limitations in analyzing structural and behavioral patterns.

In the next section, we will review the extant literature.

## II. RELATED WORK

### A. Cryptocurrency Transaction Analysis

**Empirical Analysis.** Empirical analysis plays an important role in understanding the dynamics of the cryptocurrency ecosystem [14], [15], [16] and the behaviors of addresses/entities [17], [18]. For instance, Tovanich et al. [19] and Hou et al. [20] reveal that factors such as payout schemes and pool fees influence miners' behaviors in Bitcoin mining pools, and then impact the overall system performance.

Empirical analysis can also aid in identifying cryptocurrency scams. Li et al. [21] used *CryptoScamTracker* to analyze cryptocurrency giveaway scams. He et al. [22] developed *TxPhish* to detect Ethereum scams in which users lured by high profits to fake websites are tricked into signing transactions that allow scammers to steal their crypto assets. Gomez et al. [23] explored bidirectional transactions to map cybercrime networks. Wu et al. [24] focused on identifying Ethereum-based money laundering via asset flows. Chen et al. [25] combined on-chain and off-chain data to detect fake trading volumes of famous exchanges Huobi and Binance.

**Complex Network Analysis.** Complex networks use graph theory, centrality measures, and network topology analysis to examine patterns within networks. This method explores the structures of cryptocurrency transaction networks from a macro perspective, highlighting the interconnections between nodes and the network's overall structure. Nerurkar et al. [26] and Serena et al. [27] drew parallels between cryptocurrency systems and other complex systems and identified characteristics such as small-world property, indicating most nodes in the graph are not neighbors but most of them can be reached by every other within a few hops.

Moreover, Tao et al. [28] employed an innovative random walk with a flying-back sampling method on Bitcoin transaction graphs, uncovering phenomena such as the non-rich-club effect, i.e., that high-degree nodes are not more interconnected among themselves than with lower-degree nodes. Guo et al. [29] analyzed Ethereum transaction graphs, revealing heavy-tailed property in transaction networks, i.e., the majority of nodes have a relatively low degree while a small number of nodes have a very high degree. To mitigate the issue of money laundering on blockchain networks, *DenseFlow* framework, proposed by Lin et al. [30], uses dense subgraphs and the maximum flow algorithm to trace laundering activities. This approach improves precision compared to existing methods on Ethereum, demonstrating the effectiveness of network analysis in combating money laundering.

**Machine Learning Analysis.** Machine learning is used to achieve node- or graph-level classification and prediction tasks for concrete addresses. Chaudhari et al. [31] studied utilizing temporal features to detect Bitcoin address behavioral changes and identify money laundering activities. The proposed approaches based on the decision tree (DT) by Rathore et al. [32] show high accuracy rates in detecting illicit activities and phishing scams in cryptocurrencies. Wahrstatter et al. [33] also contributed by enhancing the detection of criminal activities in Bitcoin transactions using unsupervised learning.

Additionally, various machine learning methods, such as random forest (RF), multilayer perceptron (MLP), and graph neural network (GNN), are applied to benchmark cryptocurrency datasets, including the *Elliptic Data Set* by Weber et al. [34] and datasets by Xiang et al. [13], [35]. Besides, Gai et al. [36] proposed a transformer-based anomaly detection model *BlockGPT* for the Ethereum network that demonstrated acceptable utility.

### B. LLMs in Graph Analysis and Current Gaps

**LLMs in Graph Analysis.** Several studies introduced LLMs to graph analysis, using LLM as a classifier and GNN enhancement [37]. Wang et al. [11], [38] evaluated the basic capabilities of LLMs in natural language graph problem-solving. Both studies showed limitations of LLMs, particularly in solving complex graph structures and tasks. Complementing these insights, Tang et al. [39] and Hu et al. [10] assessed the performance of LLMs in graph data analysis and prediction, compared them with specialized GNNs. Likewise, Sui et al. [40] and Jiang et al. [41] explored the effectiveness of LLMs including GPT-3.5 and GPT-4 in processing structured data, such as tables and various structured data types, introducing innovative prompting methods for performance enhancement.

Besides, Das et al. [42], Chen et al. [43], [44], and Guo et al. [9] adopted a different approach by integrating LLMs with graph data, focusing on graph structure analysis, node classification, and a range of graph processing tasks. These studies investigated the potential and limitations of LLMs in more specialized and advanced graph analysis applications, offering new insights and directions for future research in LLMs and graph data analysis. Moreover, Sun et al. [45] uniquely studied the factual knowledge of LLMs, providing a broader perspective on their comprehension capabilities, especially for lesser-known entities and facts.

**Current Gaps.** We conclude three primary research limitations according to the abovementioned work:

- Existing research mainly focuses on knowledge graphs and randomly generated graphs [42], [10], [?], [11]. However, how to measure the LLMs' ability to understand and analyze real-world cryptocurrency transaction graphs is unresolved.
- Common graph representation formats, such as GEXF and GraphML, are not ideally suited for LLMs due to their inherent space constraints. This limitation explains why recent studies have focused exclusively on testing LLMs with smaller graphs [42], [11] (e.g., graphs containing ten nodes or marginally more [11]).
- In addition to applying raw graph data to LLMs, the effect of using engineered graph features for the cryptocurrency transaction graph analysis remains insufficiently studied.

## III. PRELIMINARIES

**Transaction Graphs for Cryptocurrencies.** Cryptocurrency transaction graphs [17], [34], [29], [38] represent the flow of digital currency between entities on blockchain networks such as Bitcoin. Each node in the graph typically represents a transaction or account/address, and edges represent the transfer of cryptocurrency between these nodes. These graphs are crucial for analyzing the behavior of users, identifying patterns such as fraud or money laundering, and understanding the overall dynamics of the cryptocurrency market.

**LLMs for Structured Data Analysis.** LLMs [1] like OpenAI's Generative Pre-trained Transformer (GPT) [46] series are advanced models trained on vast amounts of text data. They excel in generating coherent and contextually relevant text based on the input they receive and can perform a variety of tasks without specific task-oriented training. While primarily designed for natural language processing, LLMs can be adapted to non-textual tasks such as analyzing structured data, including graphs [47], [37]. By converting data into a format that mimics natural language or structured prompts, researchers can leverage LLMs' powerful generative and interpretative capabilities to perform complex analyses like those required for understanding cryptocurrency transaction graphs.

**Token Limitations in LLMs.** The term *token* indicates the smallest unit of text processed by LLMs. It can vary from individual characters, such as letters and punctuation, to more complicated units, such as words or subwords. LLMs analyze and generate text by processing token sequences. Each LLM has a maximum token limit per input, which poses a significant challenge when working with extensive data sets, such as large transaction graphs. The token limit in LLMs [46] constrains data analysis in a single query and poses significant challenges for analyzing large and complex transaction graphs, potentially leading to loss of context or omission of critical structures. Although strategies such as compression and sampling have been explored, they often compromise structural integrity or contextual completeness, especially in moderately large-scale graphs.

## IV. METHODOLOGY

### A. Framework Overview

We present our framework in Fig. 1. The first step is to construct the Bitcoin transaction graph using the historical on-chain data. The second step is to extract corresponding subgraphs (or the relevant subgraphs sampled by CETraS, and graph features to Bitcoin addresses with labels. Finally, raw graphs formatted as LLM4TG and graph features are input to LLMs with proper prompts according to the tasks on different understanding levels.

### B. LLM4TG

We introduce LLM4TG, a new graph representation format designed to optimize the analysis of transaction graphs using LLMs. This format is text-based and human-readable, minimizing syntactic noise/redundancy while reducing token usage and preserving data integrity.

In our approach, LLM4TG captures essential node information and integrates edge details directly within the nodes. It organizes nodes into layers based on their type, either address or transaction, thereby maintaining the structural integrity of the graph. This hierarchical layering provides a segmented and clear overview of the network's dynamics. Each layer categorizes nodes which are further defined by properties such as degrees and token amounts, simplifying the analysis and enhancing readability.

We denote `T` as transaction and `A` as address. `<NodeID>` represents the node's ID, `<Number>` represents an integer. `<Float>` represents a real number. The syntax is displayed in Listing 1.

This format provides three key advantages for representing and analyzing transaction graphs: 1) It organizes nodes into
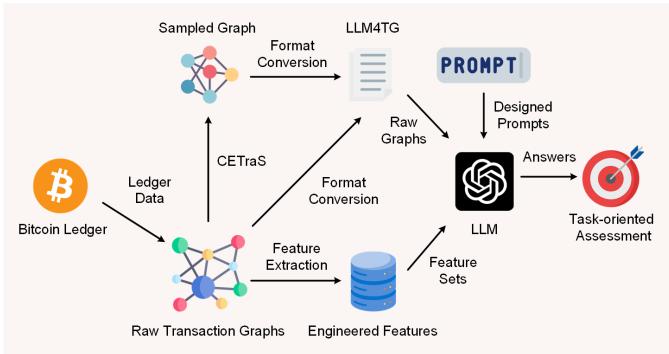


Fig. 1: LLM Evaluation Framework for Bitcoin Transaction

type-specific layers, mirroring the structure of original transaction graphs; 2) It efficiently utilizes the limited token budget of LLMs by allowing more data to be encoded; 3) It improves the interpretability of graph data for LLMs by organizing node attributes into closely associated key-value pairs.

To further demonstrate the effectiveness of LLM4TG, we compared its token consumption with other formats for the same graphs, as shown in Fig. 7. This comparison reveals that LLM4TG experiences a more gradual increase in token usage and consistently stays within the GPT-4/4o token limit across various graph sizes, making it a more efficient format, especially for larger graphs. Further details and discussion are presented in Section VI.

```
1  <LLM4TG> ::= <GraphLayer>+
2  <GraphLayer> ::= "Layer" <Number> ":" <NodeCount> <
       NodeType> "nodes" <NewLine> <Node>+
3  <NodeCount> ::= <Number>
4  <NodeType> ::= "address" | "transaction"
5  <Node> ::= <NodeA> | <NodeT>
6  <NodeA> ::= <NodeID> "address" ":" <PropertiesA> <
       NewLine>
7  <NodeT> ::= <NodeID> "transaction" ":" <PropertiesT>
       <NewLine>
8  <PropertiesA> ::= "{" <PropertyA> ("," <PropertyA>)*
       "}"
9  <PropertiesT> ::= "{" <PropertyT> ("," <PropertyT>)*
       "}"
10 <Property> ::= <InDegree> | <OutDegree> | <InValue>
       | <OutValue>
11 <PropertyA> ::= <Property> | <TimeRange>
12 <PropertyT> ::= <Property> | <InNodes> | <OutNodes>
13 <InDegree> ::= "in_degree:" <Number>
14 <OutDegree> ::= "out_degree:" <Number>
15 <InValue> ::= "in_value:" <Float>
16 <OutValue> ::= "out_value:" <Float>
17 <TimeRange> ::= "time_range:" <Number>
18 <InNodes> ::= "in_nodes:" "[" <NodeIDList> "]"
19 <OutNodes> ::= "out_nodes:" "[" <NodeIDList> "]"
20 <NodeIDList> ::= <NodeID> ("," <NodeID>)* | <Empty>
```

Listing 1: Graph Representation Definition

### C. CETraS

Despite LLM4TG's efficiency, some transaction graphs are too large for tasks like classification that involve few-shot learning, which processes multiple graphs at once. To tackle this, we introduce CETraS, a method that condenses mid-sized transaction graphs while maintaining essential structures.

We denote $I_{node}$ as the importance of the node. $a_{in/out}$ is the input/output token amount. $d_{in/out}$ is in/out-degree. $L_s$ is the shortest distance from the node to $n_0$. $\beta$ adjusts the relative significance of the node's degree. We set $\beta = 2$ as our scheme prioritizes graph connectivity. CETraS establishes a metric of importance for each node (with logic in Algorithm 1), calculated as:

$$I_{node} = \frac{log(a_{in} + a_{out} + 1) + \beta \cdot log(d_{in} + d_{out} + 1)}{L_s + 1}$$

In CETraS, nodes with lower importance are prioritized for elimination.to generate a subset of the nodes being preserved. The size of this retained subset is determined by a parameter that is flexible for specific demands. To maintain connectivity, paths connecting retained nodes are also preserved. Unlike other state-of-the-art graph summarization methods [48], [49], [50], [51] focusing on keeping accuracy for structure-relevant queries or computations on large-scale graphs (billion node-level), CETraS concentrates on accurately conveying transaction-relevant information to LLMs for moderately large-scale graphs that typically contain thousands of nodes.

---

**Algorithm 1:** CETraS

**Input:** Original transaction graph $G$; Target number of nodes to retain $N_{target}$

**Output:** Sampled transaction graph $G_{sampled}$

1 **Function** SampleGraph($G$, $N_{target}$):
2    $I_{node} \leftarrow \{v : I_{node}[v]$ for each $v$ in $V(G)\}$;
3    $P_{node} \leftarrow \frac{1}{I_{node}}$;
4    $P_{n_0} \leftarrow 0$;
5    $P_{sum} \leftarrow \sum P_{node}$;
6    **foreach** *node* $n \in G$ **do**
7      $P_n \leftarrow \frac{P_n}{P_{sum}}$;
8    **end**
9    $G_{subset} \leftarrow$ Sample from $V(G)$ with $P_{node}$ until $|G_{subset}| \geq \min\{N_{target}, |V(G)|\}$ ;
10    $G_{sampled} \leftarrow$ Initialize an empty graph;
11    **foreach** *node* $n \in G_{subset}$ **do**
12      $p \leftarrow$ Compute shortest path from $n_0$ to $n$ in $G$;
13      **foreach** *node* $m \in p$ **do**
14        Add node $m$ to $G_{sampled}$;
15      **end**
16      **foreach** *edge* $e \in p$ **do**
17        Add edge $e$ to $G_{sampled}$;
18      **end**
19    **end**
20    **return** $G_{sampled}$;
21 **return**

---

## V. EVALUATION AND ANALYSIS

### A. Settings

**Dataset.** We use two datasets for experiments, both constructed from the whole Bitcoin transaction graph, which spans a 22-month period (12 July 2019 – 26 May 2021). Specifically, BASD [35] includes eight types of subgraphs, each starting from a labeled address (denoted as $n_0$) and extending up to five hops with at most 3,000 nodes. It is one of the few datasets that provides structured graph-level representations of Bitcoin addresses. Meanwhile, BABD [13] contains labeled Bitcoin addresses, each associated with 148 engineered features. It is one of the largest publicly available datasets for Bitcoin address behavior analysis.

**LLM Selection.** We select five LLMs for evaluation: GPT-3.5 (gpt-3.5-turbo), GPT-4 (gpt-4-turbo), GPT-4o (gpt-4o), DeepSeek (deepseekv3), and LLaMA (llama3.3-70B). These models are chosen for their strong overall performance and representativeness in specific tasks. Due to differences in their input length capacities (e.g., up to 8,192 tokens for llama3.3-70B), not all models are evaluated on every task. All

models are accessed via application programming interfaces (APIs), which support prompt-based querying and response generation.

**Graph Formatting and Model Allocation.** We format raw transaction graphs as LLM4TG in the experiments. We use CETraS in level 2 and level 3 due to the token limit in few-shot prompts. We apply GPT-3.5 only in level 3 feature-based classification, given its particularly tight token limits (see Fig.7). Experiments are conducted on BASD dataset subsets and its corresponding BABD addresses.

**Prompt Engineering and Design.** We employ few-shot prompting [52] to interact with LLMs, providing only a small number of examples to guide task execution. To further align prompts with the model's capabilities and preferences, we use an LLM-feedback generation strategy, wherein the model is queried with task-specific requirements and subsequently designs the corresponding system prompts. All code, prompts, and results are publicly available at our GitHub.[1]

**Experimental Setup.** All experiments are conducted on a workstation equipped with an AMD Ryzen Threadripper PRO 5995WX and 256 GB of RAM, running Ubuntu 22.04.

### B. Level 1 - Foundational Metrics

We evaluate LLMs on sampled transaction graphs with selected nodes across layers. The evaluation adopts twelve metrics covering three perspectives: response metrics, global metrics, and node metrics:

- **Response Metrics.** To evaluate if LLM responses are correctly structured, *struct_correctness* is applied.
- **Global Metrics.** To assess the ability of LLMs to basic metrics understanding for the entire transaction graphs, LLMs need to find the node with the largest in/out-degree (*global_in/out_degree*), the node with the largest in/out-value (*global_in/out_value*), and the node with the largest difference between input and output values/degrees (*global_diff_degree/value*).
- **Node Metrics.** To investigate the capability of LLMs to understand foundational information of concrete nodes in transaction graphs, LLMs need to obtain the node's in/out-degree (*node_in/out_degree*), the node's in/out-value (*node_in/out_value*), and the node's special information (*node_special_info_a/t*). The special information for the address node is time interval; while for the transaction node is if a specific node exists in the input/output node sets.

**Results Analysis.** Table I demonstrates that LLMs are excellent at node metrics. Accuracy for most metrics is between 98.50% and 100.00%. *node_special_info_t* is the exception. This may be due to the limited capability of LLMs to match many structurally similar data in transaction graphs. Compared with node metrics, however, for the global metrics, the accuracy significantly drops, ranging from 24.00% to 58.00%. Especially, we find that compared with the other metrics (35.00% to 58.00%) in global metrics, difference-related metrics are relatively low (24.00% to 34.00%). The

[1] https://github.com/yuchen-lei/llm4tg

TABLE I: LLM Capability on foundational metrics

| Metrics | GPT-4 | GPT-4o |
|---|---|---|
| *struct_correctness* | 80.00% | 100.00% |
| *global_in_degree* | 50.00% | 44.00% |
| *global_out_degree* | 50.00% | 58.00% |
| *global_in_value* | 37.50% | 56.00% |
| *global_out_value* | 35.00% | 48.00% |
| *global_diff_degree* | 27.50% | 34.00% |
| *global_diff_value* | 27.50% | 24.00% |
| *node_in_degree* | 99.25% | 99.40% |
| *node_out_degree* | 100.00% | 99.80% |
| *node_in_value* | 98.50% | 99.00% |
| *node_out_value* | 98.50% | 99.00% |
| *node_special_info_a* | 99.08% | 100.00% |
| *node_special_info_t* | 70.88% | 67.63% |

reason for this may be that the capability of LLMs to calculate or compare is limited.

For different LLMs, i.e., GPT-4 and GPT4o, the most various points are *struct_correctness* in response metrics and *global_in/out_value* in global metrics. The enhancement of *struct_correctness* represents that the update and optimization of the LLMs improve the quality of the response format, which completely follows the requirements in prompts. Likewise, slightly improved *global_in/out_value* also illustrate the effectiveness of model upgrade. However, most of the metrics remain at similar levels, which shows the inherent flaws of LLMs for basic information understanding in transaction graphs, especially global metrics.

> **Level 1 Findings.** LLMs show consistently strong performance in node metrics, indicating a solid ability to extract explicit local information from transaction graphs. In contrast, performance on global metrics is noticeably weaker, especially for those requiring difference calculations. This disparity suggests that while models excel at retrieving localized details, they currently struggle with holistic graph reasoning and cross-node numerical comparisons.

### C. Level 2 - Characteristic Overview

To evaluate LLMs' ability to identify distinctive subgraph characteristics, we provide representative reference subgraphs and test samples without labels to avoid bias. We divide the quality of LLMs' responses into three quality levels:

- *High-quality* responses exclude irrelevant, flawed, or incorrect information.
- *Average-quality* responses include irrelevant or flawed but no incorrect information.
- *Low-quality* responses contain incorrect information.

In this context, irrelevant responses are correct but do not answer the question or provide useful information; Flawed responses have small mistakes or misleading parts but are mostly right; Incorrect responses contain false or completely wrong information.

**Overall Results.** Table II shows that GPT-4o delivers substantially higher response quality than GPT-4. High-quality

TABLE II: LLM capability on characteristic overview

| Metric | | GPT-4 | GPT-4o |
|---|---|---|---|
| High-quality | | **62.50%** | **82.50%** |
| Avg-quality | Total | 26.25% | 13.75% |
| | Flawed | _**7.50%**_ | _**12.50%**_ |
| | Irrelevant | _18.75%_ | _1.25%_ |
| Low-quality | | 11.25% | 3.75% |
| Meaningful | | **70.00%** | **95.00%** |



(a) bc1qah features high out-degrees for $n_1$, $n_2$, and both high in-/out-degree for $n_{14}$.

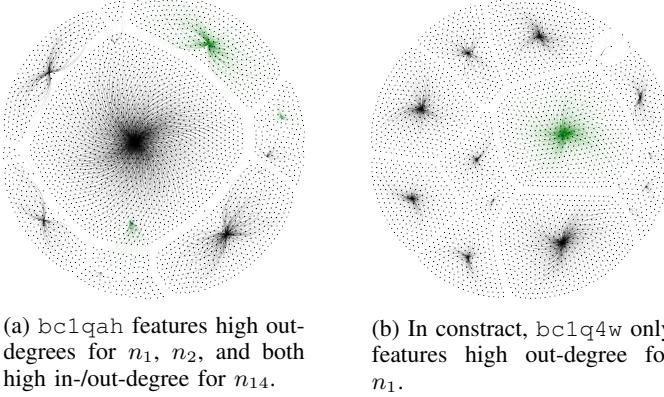(b) In constract, bc1q4w only features high out-degree for $n_1$.

Fig. 2: Examples with features

responses make up 82.50% of GPT-4o's responses, compared with 62.50% for GPT-4. For average-quality responses, GPT-4o scores 13.75%, with 12.50% flawed and 1.25% irrelevant; While GPT-4 records 26.25%, with 7.50% flawed and 18.75% irrelevant. Low-quality responses are 3.75% for GPT-4o and 11.25% for GPT-4. Overall, GPT-4o produces meaningful outputs in 95.00% of cases versus 70.00% for GPT-4.

**Illustrative Examples.** We select two examples bc1qah[2] and bc1q4w[3], to ensure each category is included to further illustrate our results, combined with visualization by *Gephi*. We apply the address identifier of node $n_0$ to represent the corresponding subgraph.

**Level 2 Example - `bc1qah`.**
**GPT-4:**
- High in-degree and out-degree of transactions; Transactions like $n_1$ and $n_2$ have high in-degrees and out-degrees. [Average - Flawed]
- Significant value differences; The in- and out-value of certain transactions show significant volume and discrepancies, such as in $n_1$ and $n_2$. [Low]

**GPT-4o:**
- High in-degree and out-degree of addresses in Layer 2; The addresses in Layer 2 have exceptionally high in-degree and out-degree values, indicating a large number of transactions. For example, address $n_{14}$ has an in-degree of 198 and out-degree of 188. [High]
- High in-value and out-value in transactions; The transactions in this graph involve significant amounts of value. For instance, transaction $n_1$ has an in-value of 77.29740945 and out-value of 77.29452845. [High]

[2]bc1qahe54yxl33clnwdtleuh4cw0fw4df62t0tnuk2
[3]bc1q4w090gzj7m9hy39l8e48gzz4uvdfn2x9y9g70h

**Level 2 Example - `bc1q4w`.**
**GPT-4:**
- High out-degree in transaction nodes; A noticeable feature in this graph is the significant out-degree in some transaction nodes. For instance, $n_1$ has an out-degree of 600. [High]
- Value discrepancy between in-value and out-value in transactions; The transaction node $n_1$ shows a discrepancy between the total incoming value (in-value: 27.69691553) and the outgoing value (out-value: 27.69184153). [Average - Irrelevant]

**GPT-4o:**
- High out-degree transactions; The transaction node $n_1$ has an extremely high out-degree of 600, meaning it distributes funds to 600 different address nodes. [High]
- Large volume of transactions; The transaction node $n_1$ also handles a substantial volume of in-value and out-value (27.69691553). [High]

*(i) Example - `bc1qah`.* The responses provided by GPT-4 exhibit average or low quality. The first response correctly identifies that nodes $n_1$ and $n_2$ have high out-degrees (Fig. 2a), but it flawedly describes their in-degrees. Although both nodes have higher in-degrees than many in the reference subgraph (typically in-degree is 1, while 4 and 3 for $n_1$ and $n_2$), their in-degrees are not particularly high within this graph. For example, node $n_{14}$ has an in-degree of 198, significantly exceeding that of $n_1$ and $n_2$. In the second response, most transactions have similar in-value and out-value; while the differences between in-value and out-value of $n_1$ and $n_2$ are both about 0.003, which is trivial.

In contrast, the responses by GPT-4o are both high-quality. In the first response, $n_{14}$ does have high in-degree and out-degree (Fig.2a), while other address nodes in Layer 2, such as $n_{13}$ and $n_{19}$, also have high in-degree and out-degree. Though the second response of GPT-4o focuses on similar characteristics as GPT-4, the description of GPT-4o is accurate and reveals the high transaction values.

*(ii) Example - `bc1q4w`.* The responses by GPT-4 are high-quality or average-quality. The first response is accurate and illustrates high out-degree of $n_1$ in Layer 1 (Fig.2b). Also, transaction nodes, including $n_{886}$ and $n_{1228}$ in Layer 3, have high out-degree that exceed 400. The second response is correct and demonstrates the difference between in-value and out-value. However, it is irrelevant since many transactions follow this pattern, which is ineffectual for transaction analysis.

In comparison, the responses by GPT-4o both have high quality. The first response is not only accurate and nearly identical to that of GPT-4, but also introduces the meaning of high out-degree for transaction nodes. The second response focuses on the same feature in the transaction node. It is accurate while meaningful for transaction analysis since the values are significantly high.

**Level 2 Findings.** Despite occasional irrelevant or incorrect responses, LLMs can extract informative structural and value-based features from Bitcoin transaction subgraphs without labels. GPT-4o shows clearer judgment in distinguishing genuinely significant features from trivial ones, and provides interpretations that reflect their analytical relevance; Whereas GPT-4 more often misjudges feature significance or emphasizes patterns of limited utility.

## D. Level 3 - Contextual Interpretation

We evaluate the LLMs' contextual interpretation capability under two settings: one using *graph features* and the other using *raw graphs*. In both settings, we adopt a few-shot prompting strategy, where the model is given several labeled subgraphs (or their features) as references and one unlabeled subgraph for explainable classification in each iteration of LLM query.
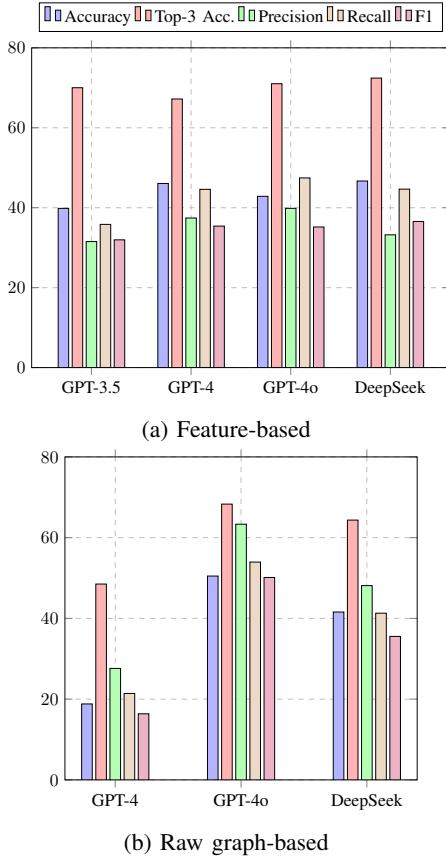


(a) Feature-based



(b) Raw graph-based

Fig. 3: Classification via different LLMs (x for LLM models, y for percentage (%); bars 1st-5th represents accuracy, top-3 accuracy, precision, recall and F1 score, respectively)

**Graph Feature-based Classification.** In this setting, we provide LLMs with subgraph features and their textual descriptions, including both labeled and unlabeled samples. Specifically, we use the ten most important features identified in BABD [13] (Table III). We then evaluate the models' performance across classes using accuracy, macro precision, recall, and F1 scores.

The overall accuracy of LLMs (Fig.3a) ranges from 39.83% to 46.68%, with similarly limited performance in macro precision, recall, and F1 scores. Although these results appear modest, they are partly attributable to the scarcity of reference samples, which constrains precise classification. Notably, the top-3 accuracy remains high, i.e., 67.20% to 72.43%, indicating that while exact class prediction is challenging under limited data, LLMs can still reliably narrow down the correct category within the top three candidates.

TABLE III: Summary of important metrics in analysis

| Label | Metric Description |
|---|---|
| S2-2 | Maximum out-degree in subgraphs |
| S1-6 | Standard deviation of in- and out-degree in subgraphs |
| S1-2 | Standard deviation of in-degree in subgraphs |
| S3 | Degree correlation of subgraphs |
| PAIa21-1 | Ratio of the minimum input token amount of an address node to the total input token amount of the address node |
| PTIa41-2 | Minimum transaction time interval of an address node |
| S6 | Longest distance between any two nodes in the subgraph |
| S5 | Closeness centrality of the subgraph |
| CI3a32-2 | Maximum change ratio in in-degree to each transaction time interval for the address node in chronological order |
| S7 | Density of the subgraph |

Beyond the overall metrics, Fig. 4 reveals uneven, category-dependent strengths. GPT-4 and GPT-4o generally set the pace, pairing standout precision in mining pools at 80.00% and 95.00% with near-saturated recall on the darknet market at 98.46% and 96.92%. DeepSeek remains competitive where a recall-weighted balance matters, yielding the strongest F1 on pools and a slight edge on exchange. GPT-3.5, though weaker overall, leads on money laundering because all models share the same recall there and its precision is highest, producing the top F1. This suggests that LLM performance is complementary rather than absolute, with relative strengths distributed unevenly across categories in this context.

When comparing models, GPT-4o attains the highest recall of 47.45% and macro precision of 39.84%, suggesting greater consistency across categories. DeepSeek records the strongest overall accuracy at 46.68%, top-3 accuracy at 72.43%, and F1 score of 36.56%, though accompanied by notably lower precision of 33.22%. GPT-4 yields comparatively balanced precision and recall, whereas GPT-3.5 trails in most metrics. Taken together, these results show that performance advantages are distributed unevenly, with different models favoring different trade-offs.

Focusing specifically on the GPT series, GPT-4o demonstrates better stability, with fewer categories having precision, recall, or F1 scores significantly low. These results suggest that despite updates leading to marginal improvements in graph feature-based classification, GPT-4o still offers mild but limited advantages over GPT-3.5 and GPT-4.

**Raw Graph-based Classification.** For raw graph-based classification, we rely on fewer subgraphs, as the higher computational cost limits the scale compared to feature-based setting.

The overall accuracy of GPT-4 (Fig.3b) is significantly lower than desired. In contrast, GPT-4o surpasses all feature-based classification tasks, achieving an accuracy of 50.49%. The macro precision, recall, and F1 scores of GPT-4 on raw graph-based classification tasks fall below those of GPT-3.5, GPT-4, and GPT-4o on graph feature-based classification tasks. Conversely, the metrics for GPT-4o on graph-based classification tasks are notably higher than those on graph feature-based classification tasks, all exceeding 50.14%. As for top-3 accuracy, GPT-4 is still lower than average, while GPT-4o achieves similar performance on raw graph-based classification compared with those on feature-based. While DeepSeek achieves competitive accuracy of 41.58% and F1
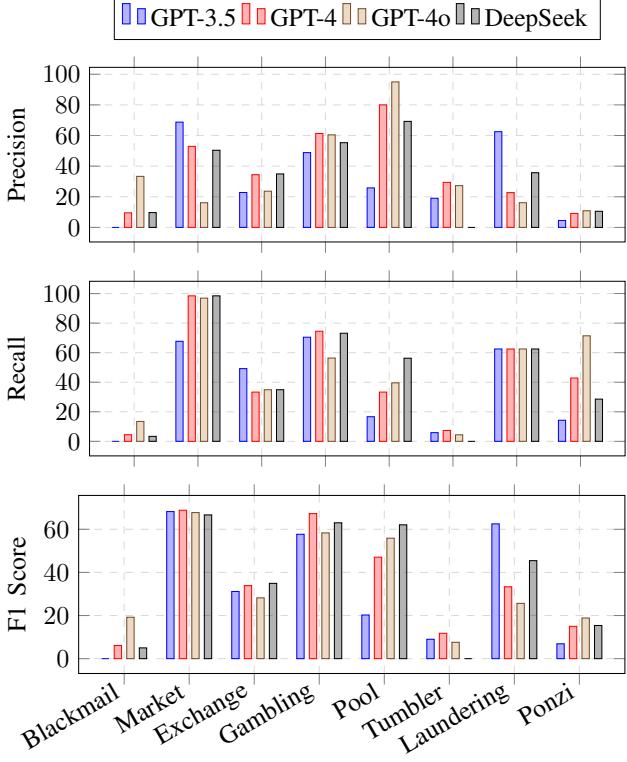
Fig. 4: LLMs' performance in contextual interpretation using graph features (x axis for category, y for rate (%); GPT-3.5 in blue bar, GPT-4 in red, GPT-4o in brown, DeepSeek in gray).



Fig. 5: LLMs' performance in contextual interpretation using raw graphs (x axis for category, y for rate (%); GPT-4 in red, GPT-4o in brown, DeepSeek in gray).

score of 35.54% using raw graphs, clearly outperforming GPT-4 yet still trailing GPT-4o.

For the specific class as shown in Figure 5, GPT-4o performs better than GPT-4 in almost all metrics, in some classes such as darknet market, the differences of F1 scores are huge. One remarkable exception is that both GPT-4 and GPT-4o recorded zero scores in all metrics on blackmail; however, for GPT-4, classes including gambling and Ponzi scheme recorded zero scores across all metrics, reflecting a huge failure in identifying relevant instances. In contrast, DeepSeek achieves competitive performance in several categories, particularly darknet market, exchange, and gambling, where its F1 scores are markedly higher than GPT-4, although it still struggles in classes such as pool and Ponzi scheme, where GPT-4o maintains clear advantages.

To sum up, GPT-4o consistently achieves substantially higher performance than GPT-4 across all metrics in raw graph-based tasks and outperforms it in most specific classes. DeepSeek also delivers competitive results, showing notable strengths in categories such as darknet market and gambling, though it remains less stable than GPT-4o overall. These findings suggest that model optimization may enhance raw graph-based classification, while different LLM architectures may offer complementary advantages across categories.

**Comparison with Traditional Models.** As illustrated in Fig.6, in the context of a very small number of reference samples, LLMs based on graph features generally perform significantly better than s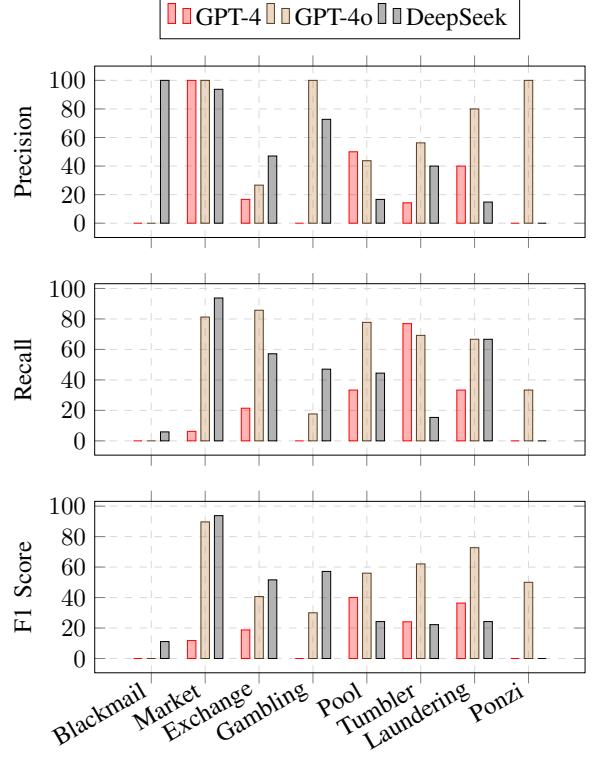upport vector machine (SVM) and MLP, yet remain inferior compared to DT, RF, CatBoost, and GNN. Among feature-based LLMs, DeepSeek delivers the strongest results, with accuracy 46.68% and F1 score 36.56% surpassing GPT-3.5 and GPT-4, though still behind tree models and GNN.

For raw-graph inputs, performance varies widely. GPT-4 only slightly outperforms support vector machine and multilayer perceptron, whereas GPT-4o achieves substantial improvements, with accuracy 50.49% and precision 63.33% comparable to tree models. This high precision is crucial for illegal address detection, as misclassifying legitimate addresses as illegal could lead to severe consequences such as account suspension or fund freezing. DeepSeek, however, lags behind GPT-4o in both accuracy and F1, though it still demonstrates relatively strong precision 48.12%, close to some tree models.

Overall, GPT-4o analyzing sampled raw graphs attains performance close to carefully engineered tree-based methods, while DeepSeek shows that feature-driven methods can also be highly competitive.

**Analysis of Classification Explanation.** Unlike traditional results of classification tasks, LLM-based outcomes can include detailed explanations that may be valuable for further analysis. However, our initial investigation indicates that these explanations are not always accurate. We choose two near-accurate explanations for `1Bsjsa`[4] and `124mpe`[5] from GPT-4o to show the potential of LLMs and how they utilize different types of data.

---

[4] `1BsjsaHST2Qohs8ZHxNHeZ1UfWhtxoKHEN`
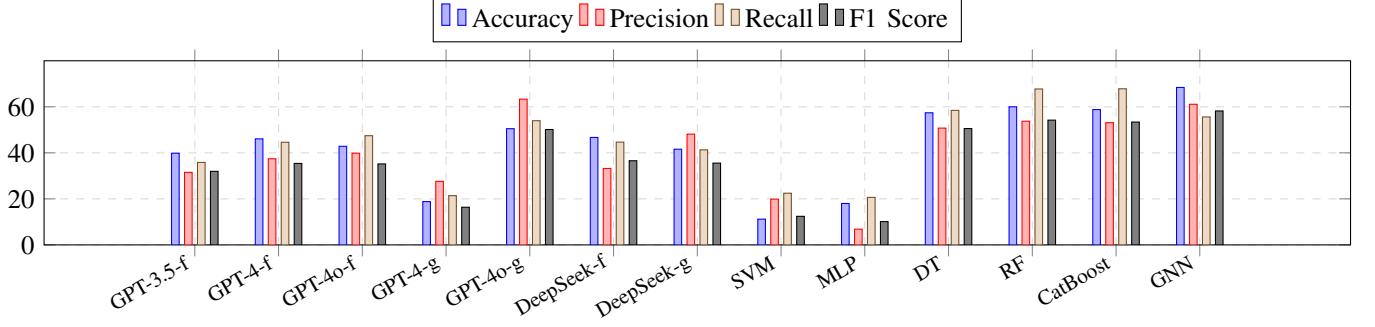[5] `124mpePGM2vEqcHUK96wQXur9vS7Vn7Kdj`

Fig. 6: Evaluation on different models (x axis for category, y for corresponding rates (%))

**Level 3 Example - `1Bsjsa`.**
**Ground Truth:** Ponzi
**Feature-based Result:** Ponzi, blackmail, tumbler
**Reason:**
- *CI3a32-2* value is extremely high, which is a distinctive characteristic of Ponzi graphs.
- *S7* is also higher than average, another indicative feature for Ponzi. Additionally, *S5* and *S6* align well with Ponzi characteristics.
- Blackmail and tumbler categories are also considered due to similar *S2-2* and *S6* values.

**Graph-based Result:** Ponzi, tumbler, pool
**Reason:**
- The graph has a structure where a single address node ($n_0$) has a high degree of both incoming and outgoing transactions (300 in-degree and 197 out-degree). This pattern is indicative of schemes like Ponzi, tumbler, and pool, where a central node is involved in numerous transactions with other nodes.
- The high number of small-value transactions and the presence of multiple transaction nodes with varying degrees further support these categories.

**Level 3 Example - `124mpe`.**
**Ground Truth:** gambling
**Feature-based Result:** gambling, darknet market, exchange
**Reason:**
- The high value of *PTIa41-2* is a strong indicator of gambling as similar high values are observed in reference graphs.
- The *S2-2*, *S1-6*, *S1-2*, and *S3* values are within ranges observed in darknet market and exchange categories, making them the next most probable categories.

**Graph-based Result:** gambling, pool, exchange
**Reason:**
- The structure and value distributions of the graph suggest that it is most likely a gambling, pool, or exchange category. The presence of large in-degree transactions and multiple out-degree addresses is characteristic of gambling and pool activities, as seen in the reference graphs.
- The significant transaction values and multiple address connections also align with the patterns found in exchange-related graphs

As demonstrated above, GPT-4o can almost accurately categorize `1Bsjsa` from two perspectives. From a graph feature perspective, *CI3a32-2* and *S7* are the primary factors influencing GPT-4o to classify it as a Ponzi address. *S5* and *S6* also play significant roles, where the values of these features compared with the reference graph features are crucial in this determination. From a raw graph perspective, the high in-degree and out-degree with massive small-value transactions is the key reason for the conclusion.

Similarly, GPT-4o is capable of accurately identifying `124mpe` from two perspectives. From a graph feature perspective, *PTIa41-2* is the primary factor influencing GPT-4o to classify it as a gambling address. *S2-2*, *S1-6*, *S1-2*, and

*S3* also play significant roles. The values of these features compared with reference graph features are crucial in determining its classification. From a raw graph perspective, the high in-degree transactions and multiple out-degree addresses with significant transaction values are key reasons for this conclusion, strongly indicating gambling, pool, or exchange activities.

These examples illustrate that for feature-based classification, adequate quantity and quality samples with labels are required to improve classification consequences since the selected features are already processed. In contrast, for graph-based classification, the quality of samples themselves seems to be also essential. The unsuitable compression of the graphs leading to much information loss may negatively affect the results. Also, the chain of thought (CoT) prompting [52] could be used for further explanations, with detailed reasoning and better classification interpretability.

**Level 3 Findings.** LLMs exhibit strong capabilities in contextual interpretation, achieving high top-3 accuracy rates even with limited data sets. This suggests that they can capture broad relational patterns and generate plausible explanations. However, their overall accuracy in exact classification remains moderate, trailing behind feature-engineered tree models and GNNs. This gap indicates that while LLMs are promising for exploratory analysis and hypothesis generation, further optimization is required before they can rival specialized models in precise detection tasks.

## VI. DISCUSSION

### A. Efficiency of Graph Representation and Processing

**Runtime Comparison of LLM4TG and CETraS.** Table IV shows that both CETraS and LLM4TG scale nearly linearly with graph size. CETraS is consistently faster, whereas LLM4TG incurs extra overhead from its enriched representation schema. Nevertheless, both methods remain efficient, with execution times within tens of seconds even for graphs exceeding 2,000 nodes.

**Cross-Format Evaluation of Graph Representations.** We conduct the comparison on foundational metrics shown in Table V, as they provide a direct measure of how well each

TABLE IV: Performance of CETraS and LLM4TG

| CETraS | | LLM4TG | |
|---|---|---|---|
| Nodes | Time (s) | Nodes | Time (s) |
| 49 | 0.0173 | 10 | 0.0009 |
| 59 | 0.0172 | 171 | 0.0515 |
| 385 | 0.1440 | 1806 | 8.3211 |
| 598 | 2.6626 | 2307 | 12.8064 |
| 2311 | 13.8108 | 2402 | 48.9033 |



Fig. 7: Token consumption in different graph formats

format preserves graph structure and attributes. Unlike higher-level tasks, which conflate representation quality with model capacity, foundational metrics isolate the contribution of the input format itself, making them particularly suitable for cross-format evaluation.

Across three LLMs, LLM4TG achieves 100% accuracy on *struct_correctness* and preserves special attributes almost entirely, with *node_special_info_a* remaining near 100% across models. In contrast, these metrics drop in alternative formats, with *struct_correctness* declining to about 87–95% under LLaMA and DeepSeek, and *node_special_info_a* falling much more sharply to roughly 23% in GEXF, 15% in GML, and only 12% in GraphML under GPT-4o, and around 45% for LLaMA and DeepSeek.

At the node level, LLM4TG also delivers consistently stronger performance, with most metrics reaching 75–80% accuracy. For example, its *node_out_value* remains close to 80% across all three models, clearly surpassing the corresponding results of alternative formats. In addition, LLM4TG maintains balanced accuracy on both degree- and value-based node metrics, while alternative formats exhibit larger fluctuations and often about 60%-70%, underscoring the stability of its representation across models.

Global metrics show a more mixed pattern. LLM4TG delivers stable performance across models, avoiding the sharp fluctuations of other formats, though it is not always the top performer. In GPT-4o, *global_in_degree* reaches about 78% in GEXF/GML versus 42% in LLM4TG, and LLaMA performs better with GEXF/GraphML for *global_out_degree* with values above 70%. Still, LLM4TG remains competitive, with about 63% on *global_out_degree* for GPT-4o and 46% on *global_in_value* for DeepSeek, while keeping variance comparatively low. The *global_diff_value* metric is weak across all formats, never exceeding 25%, suggesting the task itself is the main challenge.

Overall, LLM4TG shows a balanced performance that achieves exact structural fidelity, preserves attributes with near-perfect accuracy, and maintains consistently strong node-level results. While global metrics remain mixed, LLM4TG avoids the sharp fluctuations of alternative formats, which exhibit uneven and task-specific strengths without offering stable benefits. These findings indicate that LLM4TG provides a concise, task-aligned representation that supports stable and robust transaction graph analysis with LLMs.

**Token Consumption in Various Graph Formats** We studied the differences in token consumption among various graph representation formats under the byte pair encoding-based tokenizer `cl100k_base`, the default GPT-3.5 and GPT-4
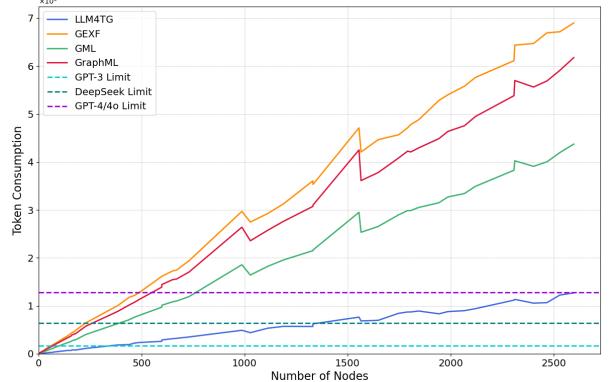
tokenizer. We chose three well-defined formats, i.e., GEXF, GML, and GraphML, for their efficiency and flexibility in representing transaction graphs. These widely used formats offer diverse encoding styles while supporting complex attributes of nodes and edges.

As in Fig.7, token usages of these formats exhibit significantly steeper curves, i.e., they consume a substantially larger number of tokens with an increasing number of nodes. These graph formats have extreme syntactic noise (i.e., great redundancy exists when representing graph data), which would result in numerous unnecessary tokens for LLMs to consume. The maximum token limits for GPT-3.5, DeepSeek, and GPT-4/4o (16,385, 64,000, and 128,000) are indicated by horizontal lines, where token usages of these formats surpass the GPT-4/4o limit when the number of nodes of a single graph exceeds approximately 500 to 750. Therefore, these graph representation formats are insufficient for LLMs to analyze transaction graphs.

We can also observe a counterintuitive phenomenon, i.e., at some points of the graphic, the number of necessary tokens decreases while the number of nodes increases. The explanation is that the total size of a graph is determined not just by the number of nodes, but also by the edges and their corresponding attributes (e.g., a graph may have more nodes but fewer edges than another).

### B. LLMs for Transaction Graphs: Benefits and Challenges

**Potential Advantages of LLMs.** We summarise four advantages observed in applying LLMs to transaction graph analysis.

- *Robust Performance under Limited Data.* LLMs can deliver meaningful results even with very limited labeled samples, achieving relatively high top-3 accuracy together with interpretable outputs. This makes them particularly useful for transaction types with scarce annotations, such as money laundering or Ponzi schemes, where traditional models often struggle due to data scarcity.
- *Context-aware Interpretation.* Beyond identifying statistical patterns, LLMs can capture aspects of human intent and behavioral context that are often overlooked by traditional data-driven approaches. Leveraging their extensive training on natural language and human behavior, LLMs provide more nuanced explanations of

TABLE V: LLM Capability on foundational metrics across graph formats and models

| Metrics | GPT-4o | | | | LLaMA | | | | DeepSeek | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | LLM4TG | GEXF | GML | GraphML | LLM4TG | GEXF | GML | GraphML | LLM4TG | GEXF | GML | GraphML |
| *struct_correctness* | 100.00% | 95.83% | 95.83% | 100.00% | 100.00% | 87.50% | 95.83% | 87.50% | 100.00% | 95.83% | 91.67% | 95.83% |
| *global_in_degree* | 41.67% | 78.26% | 78.26% | 75.00% | 54.17% | 71.43% | 69.57% | 76.19% | 50.00% | 69.57% | 50.00% | 60.87% |
| *global_out_degree* | 62.50% | 60.87% | 56.52% | 54.17% | 50.00% | 71.43% | 73.91% | 76.19% | 58.33% | 47.83% | 36.36% | 39.13% |
| *global_in_value* | 54.17% | 30.43% | 52.17% | 25.00% | 33.33% | 33.33% | 21.74% | 47.62% | 45.83% | 56.52% | 45.45% | 47.83% |
| *global_out_value* | 25.00% | 26.09% | 30.43% | 25.00% | 45.83% | 9.52% | 8.70% | 19.05% | 37.50% | 26.09% | 22.73% | 26.09% |
| *global_diff_degree* | 37.50% | 43.48% | 43.48% | 41.67% | 33.33% | 47.62% | 43.48% | 52.38% | 29.17% | 30.43% | 18.18% | 34.78% |
| *global_diff_value* | 8.33% | 8.70% | 8.70% | 8.33% | 25.00% | 0.00% | 4.35% | 4.76% | 16.67% | 8.70% | 4.55% | 4.35% |
| *node_in_degree* | 73.75% | 75.22% | 75.65% | 73.33% | 73.33% | 66.19% | 60.87% | 69.05% | 75.42% | 73.91% | 67.27% | 68.26% |
| *node_out_degree* | 75.83% | 82.17% | 80.43% | 79.17% | 76.25% | 67.14% | 54.78% | 70.48% | 77.50% | 76.96% | 72.27% | 70.87% |
| *node_in_value* | 75.42% | 70.43% | 71.74% | 67.92% | 73.75% | 69.05% | 58.26% | 68.10% | 77.50% | 69.13% | 67.27% | 67.83% |
| *node_out_value* | 77.08% | 74.78% | 71.74% | 71.67% | 78.33% | 55.71% | 55.65% | 58.57% | 79.17% | 70.00% | 68.64% | 65.65% |
| *node_special_info_a* | 100.00% | 23.08% | 14.66% | 11.67% | 99.17% | 47.17% | 47.86% | 44.34% | 100.00% | 45.69% | 46.43% | 49.57% |
| *node_special_info_t* | 71.67% | 81.42% | 64.91% | 79.17% | 78.33% | 76.92% | 72.57% | 72.12% | 83.90% | 63.16% | 64.81% | 68.14% |

suspicious transactions. Although these interpretations are not always fully accurate, they offer valuable insights for understanding underlying motivations and guiding further security analysis.

- *Graph-level Understanding.* LLMs can effectively extract node attributes and synthesize meaningful overviews of transaction graphs, enabling the identification of complex patterns and relational structures. Compared with traditional algorithms, which often struggle with high-dimensional and context-rich data, LLMs demonstrate a stronger capacity to process raw graph representations and capture intricate dependencies and anomalies.
- *Support for Cybercrime and Security Analysis.* LLMs can assist in identifying anomalous transaction patterns, inferring the potential motivations behind suspicious activities, and providing interpretable reasoning processes. These capabilities are particularly valuable in security-critical contexts such as anti-money laundering, fraud detection, and darknet market tracing, offering law enforcement and researchers additional clues and explainable evidence.

**Remaining Challenges.** We also outline three key challenges in applying LLMs to cryptocurrency transaction graphs.

- *Token Limits.* They restrict the amount of graph data that can be processed at once, hindering the analysis of large cryptocurrency transaction graphs. With only limited portions available, LLMs often lack sufficient context to deliver accurate insights, making it difficult to capture the complexity and nuances of the data.
- *Reference Graph Selection.* Choosing representative labeled samples is challenging, as different subgraphs emphasize different features and may introduce bias. Such choices directly affect classification outcomes and their explanations. The problem is further compounded by token limits, which restrict the number of reference samples that can be included, especially for raw graph data, thereby constraining the model's ability to generalize.
- *Explanation Accuracy.* Improving the reliability of LLM-generated explanations remains a major challenge. Current outputs may contain incomplete reasoning or subtle inaccuracies, which limits their utility for high-stakes applications such as illicit activity detection. Developing

more rigorous prompting strategies or integration with external knowledge sources could enhance the fidelity of these explanations. Achieving this would not only improve interpretability but also strengthen the overall effectiveness of cryptocurrency transaction analysis.

### C. Factors Influencing LLM Performance

The performance of LLMs on transaction graphs is shaped by multiple factors, with the first being the model bottleneck. Although major updates and optimizations can enhance the capability of LLMs in certain aspects, they may not resolve the inherent limitations of LLMs in transaction graph understanding. For example, when using the same graph feature data for classification, GPT-4 and GPT-4o show only slight improvements over GPT-3.5. At present, a more effective way to boost performance may involve increasing the number of labeled samples and enriching the feature set in graph feature-based classifications, which also require fewer tokens compared to processing raw graphs.

Another factor lies in the size and complexity of transaction graphs. Smaller graphs tend to yield better results on foundational metrics, while larger graphs may reduce effectiveness. At the contextual interpretation level, we observed that even when using the same data, accuracy may vary slightly, by around 5%. Moreover, due to the application of LLM4TG, there is a slight loss of temporal information, which could affect results when raw graphs are used.

Finally, the representation of data and the choice of features also play a crucial role. The type of input, such as graph features versus raw graphs, and the availability of labeled samples directly influence performance. Graph feature-based data is generally more efficient in terms of LLM token consumption, and its classification effectiveness can be further improved through the inclusion of additional features and labels.

### VII. CONCLUSION

This work evaluates LLMs' capabilities in analyzing Bitcoin transaction graphs. We introduced a three-level framework along with two key innovations: LLM4TG format to enhance readability and reduce graph sizes, and CETraS algorithm to

optimize graph simplification. Our experiments demonstrate satisfactory accuracy in foundational metrics, effectiveness in obtaining useful overview characteristics, and solid top-3 accuracy in classification tasks. These findings highlight the significant potential of LLMs and establish a foundation for their broader application in cryptocurrency analysis.

In addressing our research questions, we find that (RQ1) the LLM4TG format, combined with CETraS, reduces token redundancy and preserves graph structure, enabling efficient analysis of large Bitcoin graphs; (RQ2) our three-level framework effectively measures LLMs' capacity to understand transaction graphs, capturing both local node details and broader behavioral patterns; and (RQ3) while engineered features simplify analysis through pre-computed metrics, raw graph data allows deeper insights and higher potential accuracy, but also exposes larger performance disparities across models.

## ACKNOWLEDGMENT

## REFERENCES

[1] Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, et al. A survey on evaluation of large language models. *ACM Transactions on Intelligent Systems and Technology (TIST)*, 15(3):1–45, 2024.

[2] Bonan Min, Hayley Ross, Elior Sulem, Amir Pouran Ben Veyseh, Thien Huu Nguyen, Oscar Sainz, Eneko Agirre, Ilana Heintz, and Dan Roth. Recent advances in natural language processing via large pre-trained language models: A survey. *ACM Computing Surveys (CSUR)*, 56(2):1–40, 2023.

[3] Robert Tinn, Hao Cheng, Yu Gu, Naoto Usuyama, Xiaodong Liu, Tristan Naumann, Jianfeng Gao, and Hoifung Poon. Fine-tuning large neural language models for biomedical natural language processing. *Patterns*, 4(4), 2023.

[4] Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. *Advances in Neural Information Processing Systems (NeurIPS)*, 36, 2024.

[5] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed El-hoseiny. MiniGPT-4: Enhancing vision-language understanding with advanced large language models. In *International Conference on Learning Representations (ICLR)*, 2024.

[6] Jingxuan He and Martin Vechev. Large language models for code: Security hardening and adversarial testing. In *ACM SIGSAC Conference on Computer and Communications Security (CCS)*, pages 1865–1879, 2023.

[7] Gustavo Sandoval, Hammond Pearce, Teo Nys, Ramesh Karri, Siddharth Garg, and Brendan Dolan-Gavitt. Lost at C: A user study on the security implications of large language model code assistants. In *USENIX Security Symposium (USENIX Sec)*, pages 2205–2222, 2023.

[8] Karan Singhal, Shekoofeh Azizi, Tao Tu, S Sara Mahdavi, Jason Wei, Hyung Won Chung, Nathan Scales, Ajay Tanwani, Heather Cole-Lewis, Stephen Pfohl, et al. Large language models encode clinical knowledge. *Nature*, 620(7972):172–180, 2023.

[9] Jiayan Guo, Lun Du, and Hengyu Liu. GPT4Graph: Can large language models understand graph structured data? an empirical evaluation and benchmarking. *arXiv preprint arXiv:2305.15066*, 2023.

[10] Yuntong Hu, Zheng Zhang, and Liang Zhao. Beyond text: A deep dive into large language models' ability on understanding graph data. *arXiv preprint arXiv:2310.04944*, 2023.

[11] Heng Wang, Shangbin Feng, Tianxing He, Zhaoxuan Tan, Xiaochuang Han, and Yulia Tsvetkov. Can language models solve graph problems in natural language? In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 36, pages 30840–30861, 2023.

[12] Zeyang Zhang, Xin Wang, Ziwei Zhang, Haoyang Li, Yijian Qin, and Wenwu Zhu. LLM4DyG: Can large language models solve spatial-temporal problems on dynamic graphs? In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, page 4350–4361, 2024.

[13] Yuexin Xiang, Yuchen Lei, Ding Bao, Tiantian Li, Qingqing Yang, Wenmao Liu, Wei Ren, and Kim-Kwang Raymond Choo. BABD: A Bitcoin address behavior dataset for pattern analysis. *IEEE Transactions on Information Forensics and Security (TIFS)*, 19:2171–2185, 2024.

[14] Qin Wang, Guangsheng Yu, and Shiping Chen. Cryptocurrency in the aftermath: Unveiling the impact of the SVB collapse. *IEEE Transactions on Computational Social Systems (TCSS)*, 2024.

[15] Jintao Huang, Ningyu He, Kai Ma, Jiang Xiao, and Haoyu Wang. Miracle or mirage? a measurement study of NFT rug pulls. *Proceedings of the ACM on Measurement and Analysis of Computing Systems (SIGMETRICS)*, 7(3):1–25, 2023.

[16] Qin Wang, Guangsheng Yu, Yilin Sai, Caijun Sun, Lam Duc Nguyen, Sherry Xu, and Shiping Chen. An empirical study on snapshot DAOs. *arXiv preprint arXiv:2211.15993*, 2022.

[17] Ting Chen, Zihao Li, Yuxiao Zhu, Jiachi Chen, Xiapu Luo, John Chi-Shing Lui, Xiaodong Lin, and Xiaosong Zhang. Understanding Ethereum via graph analysis. *ACM Transactions on Internet Technology (TOIT)*, 20(2):1–32, 2020.

[18] Bingyu Gao, Haoyu Wang, Pengcheng Xia, Siwei Wu, Yajin Zhou, Xiapu Luo, and Gareth Tyson. Tracking counterfeit cryptocurrency end-to-end. *Proceedings of the ACM on Measurement and Analysis of Computing Systems (SIGMETRICS)*, 4(3):1–28, 2020.

[19] Natkamon Tovanich, Nicolas Soulié, Nicolas Heulot, and Petra Isenberg. An empirical analysis of pool hopping behavior in the Bitcoin blockchain. In *IEEE International Conference on Blockchain and Cryptocurrency (ICBC)*, 2021.

[20] Binbing Hou and Feng Chen. A study on nine years of Bitcoin transactions: Understanding real-world behaviors of bitcoin miners and users. In *IEEE International Conference on Distributed Computing Systems (ICDCS)*, 2020.

[21] Xigao Li, Anurag Yepuri, and Nick Nikiforakis. Double and nothing: Understanding and detecting cryptocurrency giveaway scams. In *Network and Distributed Systems Security (NDSS) Symposium*, 2023.

[22] Bowen He, Yuan Chen, Zhuo Chen, Xiaohui Hu, Yufeng Hu, Lei Wu, Rui Chang, Haoyu Wang, and Yajin Zhou. TxPhishScope: Towards detecting and understanding transaction-based phishing on Ethereum. In *ACM SIGSAC Conference on Computer and Communications Security (CCS)*, pages 120–134, 2023.

[23] Gibran Gomez, Pedro Moreno-Sanchez, and Juan Caballero. Watch your back: Identifying cybercrime financial relationships in Bitcoin through back-and-forth exploration. In *ACM SIGSAC Conference on Computer and Communications Security (CCS)*, 2022.

[24] Jiajing Wu, Dan Lin, Qishuang Fu, Shuo Yang, Ting Chen, Zibin Zheng, and Bowen Song. Toward understanding asset flows in crypto money laundering through the lenses of ethereum heists. *IEEE Transactions on Information Forensics and Security (TIFS)*, 2024.

[25] Jialan Chen, Dan Lin, and Jiajing Wu. Do cryptocurrency exchanges fake trading volumes? an empirical analysis of wash trading based on data mining. *Physica A: Statistical Mechanics and its Applications*, 586:126405, 2022.

[26] Pranav Nerurkar, Dhiren Patel, Yann Busnel, et al. Dissecting Bitcoin blockchain: Empirical analysis of Bitcoin network (2009–2020). *Journal of Network and Computer Applications*, 2021.

[27] Luca Serena, Stefano Ferretti, and Gabriele D'Angelo. Cryptocurrencies activity as a complex network: Analysis of transactions graphs. *Peer-to-Peer Networking and Applications*, 2022.

[28] Bishenghui Tao, Hong-Ning Dai, Jiajing Wu, Ivan Wang-Hei Ho, Zibin Zheng, and Chak Fong Cheang. Complex network analysis of the Bitcoin transaction network. *IEEE Transactions on Circuits and Systems II: Express Briefs*, 2021.

[29] Dongchao Guo, Jiaqing Dong, and Kai Wang. Graph structure and statistical properties of Ethereum transaction relationships. *Information Sciences*, 2019.

[30] Dan Lin, Jiajing Wu, Yunmei Yu, Qishuang Fu, Zibin Zheng, and Changlin Yang. DenseFlow: Spotting cryptocurrency money laundering in Ethereum transaction graphs. In *Proceedings of the ACM on Web Conference (WWW)*, 2024.

[31] Deepesh Chaudhari, Rachit Agarwal, and Sandeep Kumar Shukla. Towards malicious address identification in Bitcoin. In *IEEE International Conference on Blockchain (Blockchain)*, 2021.

[32] M Mazhar Rathore, Sushil Chaurasia, and Dhirendra Shukla. Mixers detection in Bitcoin network: a step towards detecting money laundering

in crypto-currencies. In *IEEE International Conference on Big Data (BigData)*, 2022.

[33] Anton Wahrstätter, Jorão Gomes, Sajjad Khan, and Davor Svetinovic. Improving cryptocurrency crime detection: Coinjoin community detection approach. *IEEE Transactions on Dependable and Secure Computing (TDSC)*, 2023.

[34] Mark Weber, Giacomo Domeniconi, Jie Chen, Daniel Karl I Weidele, Claudio Bellei, Tom Robinson, and Charles E Leiserson. Anti-money laundering in Bitcoin: Experimenting with graph convolutional networks for financial forensics. *arXiv preprint arXiv:1908.02591*, 2019.

[35] Yuexin Xiang, Tiantian Li, and Yuquan Li. Leveraging subgraph structure for exploration and analysis of Bitcoin address. In *IEEE International Conference on Big Data (BigData)*, pages 1957–1962, 2022.

[36] Yu Gai, Liyi Zhou, Kaihua Qin, Dawn Song, and Arthur Gervais. Blockchain large language models. *arXiv preprint arXiv:2304.12749*, 2023.

[37] Yuhan Li, Zhixun Li, Peisong Wang, Jia Li, Xiangguo Sun, Hong Cheng, and Jeffrey Xu Yu. A survey of graph meets large language model: Progress and future directions. In *International Joint Conference on Artificial Intelligence (IJCAI)*, 2024.

[38] Guangsheng Yu, Qin Wang, Tanzeela Altaf, Xu Wang, Xiwei Xu, and Shiping Chen. Predicting nft classification with GNN: A recommender system for web3 assets. In *IEEE International Conference on Blockchain and Cryptocurrency (ICBC)*, 2023.

[39] Jianheng Tang, Qifan Zhang, Yuhan Li, Nuo Chen, and Jia Li. Grapharena: Evaluating and exploring large language models on graph computation. In *The Thirteenth International Conference on Learning Representations (ICLR)*, 2025.

[40] Yuan Sui, Mengyu Zhou, Mingjie Zhou, Shi Han, and Dongmei Zhang. Table meets LLM: Can large language models understand structured table data? a benchmark and empirical study. In *ACM International Conference on Web Search and Data Mining (WSDM)*, pages 645–654, 2024.

[41] Jinhao Jiang, Kun Zhou, zican Dong, KeMing Ye, Xin Zhao, and Ji-Rong Wen. StructGPT: A general framework for large language model to reason over structured data. In *The Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2023.

[42] Debarati Das, Ishaan Gupta, Jaideep Srivastava, and Dongyeop Kang. Which modality should I use - text, motif, or image? : Understanding graphs with large language models. In *Findings of the Association for Computational Linguistics (NAACL)*, pages 503–519, 2024.

[43] Zhikai Chen, Haitao Mao, Hang Li, Wei Jin, Hongzhi Wen, Xiaochi Wei, Shuaiqiang Wang, Dawei Yin, Wenqi Fan, Hui Liu, et al. Exploring the potential of large language models (LLMs) in learning on graphs. *ACM SIGKDD Explorations Newsletter*, 25(2):42–61, 2024.

[44] Zhikai Chen, Haitao Mao, Hongzhi Wen, Haoyu Han, Wei Jin, Haiyang Zhang, Hui Liu, and Jiliang Tang. Label-free node classification on graphs with large language models (LLMs). In *The Twelfth International Conference on Learning Representations (ICLR)*, 2024.

[45] Kai Sun, Yifan Ethan Xu, Hanwen Zha, Yue Liu, and Xin Luna Dong. Head-to-tail: How knowledgeable are large language models (LLMs)? a.k.a. will LLMs replace knowledge graphs? In *Proceedings of the NAACL-HLT*, 2024.

[46] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

[47] Bowen Jin, Gang Liu, Chi Han, Meng Jiang, Heng Ji, and Jiawei Han. Large language models on graphs: A comprehensive survey. *IEEE Transactions on Knowledge and Data Engineering (TKDE)*, 2024.

[48] Mahdi Hajiabadi, Jasbir Singh, Venkatesh Srinivasan, and Alex Thomo. Graph summarization with controlled utility loss. In *ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD)*, pages 536–546, 2021.

[49] Shinhwan Kang, Kyuhan Lee, and Kijung Shin. Personalized graph summarization: formulation, scalable algorithms, and applications. In *IEEE International Conference on Data Engineering (ICDE)*, pages 2319–2332, 2022.

[50] Kyuhan Lee, Hyeonsoo Jo, Jihoon Ko, Sungsu Lim, and Kijung Shin. Ssumm: Sparse summarization of massive graphs. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)*, pages 144–154, 2020.

[51] Meiquan Lai, Yaqi Huang, Zhidan Liu, and Kaishun Wu. An optimized lossless graph summarization for large-scale graphs. In *IEEE International Conference on Parallel and Distributed Systems (ICPADS)*, pages 355–362. IEEE, 2023.

[52] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 35, pages 24824–24837, 2022.