Are LLMs Truly Graph-Savvy? A Comprehensive Evaluation of Graph Generation

Anonymous ACL Submission

Abstract

While large language models (LLMs) have demonstrated impressive capabilities across diverse tasks, their ability to generate valid graph structures remains underexplored. We evaluate fifteen state-of-the-art LLMs on five specialized graph generation tasks spanning delivery networks, social networks, quantum circuits, genedisease networks, and transportation systems. We also test the LLMs using 3 different prompt types: direct, iterative feedback, and programaugmented. Models supported with explicit reasoning modules (o3-mini-high, o1, Claude 3.7 Sonnet, DeepSeek-R1) solve more than twice as many tasks as their general-purpose peers, independent of parameter count. Error forensics reveals two recurring failure modes: Llama-family models often violate basic structural constraints, whereas Claude models respect topology but mismanage higher-order logical rules. Allowing models to iteratively refine their answers yields uneven gains, underscoring fundamental differences in error-correction capacity. This work demonstrates that graph competence stems from specialized architectural design rather than scale, establishing a framework for developing truly graph-savvy language models. Results and verification scripts available at https://github.com/anonymized-forthe-blind-review.

1 Introduction

002

006

007

011

017

019

027

Large Language Models (LLMs) have revolutionized natural language processing by achieving stateof-the-art performance on a diverse range of tasks, from translation and summarization to complex reasoning (Brown et al., 2020). Despite these impressive advancements in text generation, their ability to handle structured data, particularly graphs, remains an emerging area of research. Graphs, which consist of nodes (representing entities) and edges (representing relationships), are fundamental to a wide spectrum of applications including social network analysis, biological systems modeling, and transportation planning. However, while LLMs demonstrate remarkable fluency in natural language, their performance in generating and reasoning about graph structures is often hindered by a persistent challenge: *hallucination*. In many cases, LLMs produce graph outputs that are syntactically plausible yet factually or structurally incorrect (Merrer and Tredan, 2024). While these failures are well documented on individual graph benchmarks, no broad, cross-domain evaluation has yet been performed. 043

045

047

049

051

054

055

057

059

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

077

078

079

To address these gaps, our contribution is three-fold:

(i) We introduce a novel evaluation framework comprising five specialized graph problems designed to challenge and assess LLMs' structural reasoning capabilities: (1) a Time-Dependent Delivery Network with complex spatiotemporal constraints; (2) a Directed Social Network with hierarchical influence relationships; (3) a Quantum Circuit Design requiring an understanding of quantum gate operations; (4) a Gene-Disease Association Network modeling bipartite relationships; and (5) an Optimal Transportation Network with robust connectivity requirements. These problems intentionally extend beyond conventional datasets to mitigate the effects of memorization, identified as confounding factors in the evaluation of LLM performance. Additionally, since these problems are open-ended, they can have multiple solutions which can prove challenging to some LLMs with limited reasoning abilities.

(ii) We conduct a comprehensive evaluation using fifteen state-of-the-art language models spanning multiple architectural families and parameter scales. This selection enables us to conduct thorough comparisons across different architectures, which previous taxonomies by Ren et al. (2024) indicate are crucial for understanding the specific limitations of models in graph processing.

120

121

122

123

124

125

127

128

129

130 131

132

133

135

084

(iii) We systematically investigate three prompting paradigms: direct prompting, iterative feedback, and program-augmented prompting. Building upon the reasoning frameworks of the study, we examine whether these prompting approaches can effectively address the hallucination challenges documented by Tonmoy et al. (2024) and improve structural fidelity in the graph output.

Before detailing our evaluation, we review prior attempts to evaluate LLM graph skills. Early efforts to explore the graph capabilities of LLMs have yielded promising but mixed results. For instance, Yao et al. (2024) introduced LLM4GraphGen, which systematically evaluates the ability of LLMs to generate graphs based on structural rules and distributions. Their findings suggest that while models like GPT-4 exhibit some capacity for rulebased and distribution-based graph generation, conventional prompting methods (e.g., few-shot or chain-of-thought) do not consistently improve performance. In parallel, Wang et al. (2023) proposed the NLGraph benchmark, a set of graph reasoning tasks that ranges from basic connectivity checks to complex algorithmic challenges such as maximum flow and bipartite graph matching. Their study showed that while LLMs demonstrate preliminary reasoning abilities, their performance deteriorates as task complexity increases, and standard prompting strategies often fail to enhance results. Notably, both studies highlight that LLMs have difficulty generalizing beyond examples they have seen. This raises concerns about whether they genuinely learn graph structures or simply rely on memorization, and shows the need for more robust evaluations that go beyond standard datasets and assess LLMs' ability to construct and reason about unseen graphs.

Advances in reasoning frameworks further illustrate both the potential and limitations of LLMs for graph-related tasks. The graph chain-of-thought (Graph-CoT) framework of Jin et al. (2024) promotes iterative reasoning by structuring LLM reasoning paths through explicit graph structures and demonstrating improved performance in complex graph-related inference tasks. Similarly, the Graph of Thoughts (GoT) framework introduced by Besta et al. (2024) models reasoning as a graph rather than a traditional tree, allowing LLMs to explore non-linear reasoning paths that better capture dependencies in structured data. Although these methods significantly improve reasoning accuracy, they do not fully address graph generation. Additionally, approaches such as the GCoder by Zhang et al. (2024) have explored integrating LLM with codebased methodologies to solve generalized graph problems, and have demonstrated substantial improvements over traditional natural language reasoning paradigms. Meanwhile, broader investigations into hallucination mitigation, such as the comprehensive survey by Tonmoy et al. (2024), underscore the need for more robust evaluation protocols that explicitly detect and quantify structural inconsistencies in graph outputs. These collective efforts indicate that while LLMs are becoming increasingly capable of handling graph-based reasoning, their ability to reliably generate novel, structurally valid graphs remains an open challenge requiring further study. 136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

Lastly, Merrer and Tredan (2024) examined how LLMs generate known graphs such as Zachary's Karate Club and Les Misérables. However, their approach is limited in scope as it relies on a small set of benchmark graphs, many of which are widely available in public datasets and may have been seen during model training. Furthermore, their evaluation is based on single-prompt interactions without testing the robustness of model responses across multiple attempts or under varied prompt conditions. This narrow evaluation methodology fails to capture the broader generalization and reasoning abilities of LLMs in generating unseen graph structures, leaving critical questions unanswered regarding their ability to construct complex, structured graphs beyond memorization.

Through (i) crafting five diverse, unconstrained graph tasks, (ii) benchmarking fifteen distinct LLM architectures, and (iii) evaluating three prompting strategies, we offer a comprehensive evaluation of LLM graph-generation capabilities. Our results quantify current performance boundaries with statistical rigor and establish a reusable framework for assessing and improving structural fidelity in LLM outputs. Through our unique approach of targeting structural reasoning rather than memorization, we directly address the gap identified by recent surveys (Yu et al., 2025; Li et al., 2024), and take a step toward building graph-savvy language models that generate and reason about complex networks with higher fidelity and consistency.

2 Methodology

In this section, we describe the procedures used to design our five specialized graph-generation tasks, the verification pipeline for evaluating generated 183 solutions, and the experimental setup employed to assess model performance.

186

187

191

192

193

194

195

197

198

199

200

204

209

210

211

212

214

215

216

217

218

219

221

223

224

228

234

We evaluate the ability of Large Language Models (LLMs) to generate valid graphs using five tasks that each emphasize a distinct set of structural and logical challenges. These tasks are inspired by classical problem domains, including combinatorial optimization, network analysis, and biological systems modeling. Detailed prompts and constraints can be seen in the Appendix.

Time-Dependent Delivery Network: This scenario requires scheduling deliveries across multiple locations using a fleet of vehicles. Constraints include vehicle and storage capacities, dynamically adjusted travel times, and delivery time windows. It is similar to a time-windowed Vehicle Routing Problem (VRP) (Toth and Vigo, 2001) often encountered in logistics and supply-chain management, where resource utilization and schedule feasibility are essential.

Directed Social Network with Influence Relationships: We construct a social network in which users (categorized by trust scores) exert directed influence over others. The graph must remain acyclic while respecting category-based constraints (e.g., celebrities requiring sufficient outgoing edges). This setup reflects common problems in social network analysis, trust-based recommendation systems, and hierarchical structures where influence needs to be rigorously defined and free of feedback loops.

Quantum Circuit: This task involves organizing qubits, gates (single- and multi-qubit), and measurement operations under strict limitations on gate adjacency, temporal layering, and measurement rules. It mirrors quantum circuit scheduling challenges (Romero-Alvarez et al., 2024), where quantum gates must be placed in a Directed Acyclic Graph (DAG)-like structure, to ensure no conflicting operations and respect hardware constraints (such as non-adjacent CNOT requirements).

Gene-Disease Association Network: A bipartite graph is formed between genes and diseases, with each node set governed by specific degree constraints and edges indicating association strengths in the range [0.0,1.0]. In particular, our design draws inspiration from recent findings on the bipartite structure of vertebrate centromeres (Sacristan et al., 2024). This problem is an example of biological networks (e.g., gene-regulatory or gene-disease association mappings) that capture the confidence of links between genetic factors and clinical conditions. The valid bipartite structure and bounded association strengths are essential for realistic biological modeling.

Optimal Transportation Network: In this problem, LLMs need to develop a strongly connected, cost-effective, and resilient network of cities (nodes) and directed roads (edges). Important constraints include limits on road length and cost to ensure accessibility for the population. Additionally, the design should incorporate redundancy through multiple edges to enhance resilience. This problem is similar to multi-constraint transportation (Li et al., 2023) or flow networks, with a particular focus on two-edge robustness and minimizing path lengths to ensure that the network remains reliable and efficient under stress.

We evaluate a set of fifteen state-of-the-art LLMs, spanning multiple architectures and parameter sizes. These include GPT-40 (January 29 version), GPT-40-mini, o1, and o3-mini-high by OpenAI (2024a,b,c); Claude 3.5 Sonnet, Claude 3.5 Haiku, and Claude 3.7 Sonnet (with extended thinking) by Anthropic (2024a,b,c); Gemini 2.0 Pro and Gemini 2.0 Flash by Google (2024a,b); Llama 3.1 (8B), Llama 3.1 (405B), and Llama 3.2 (3B) by Meta AI (2024a,b); DeepSeek-V3 and DeepSeek-R1 by DeepSeek AI (2025, 2024); and Grok-V3 by xAI (2025). Models from Llama family are run in Ollama (2025), allowing direct control over parameter settings and token decoding, while the remaining models are accessed through their respective chat-based interfaces following each provider's recommended prompt-completion protocol.

We explore three prompting paradigms:

- *Direct Prompting:* The model receives a single, comprehensive prompt containing the entire task description, without additional feedback during generation.
- *Iterative Prompting:* After the initial direct prompt, if the model's output is unsatisfactory, it receives the verification script output as feedback. This feedback helps to refine the subsequent response, allowing for a multi-step corrective process.
- Program-Augmented Prompting: In the initial prompt, we include both the task description and the verification script. The model is encouraged to refer to this script during the generation process to self-assess and ensure

384

335

336

that the output meets the specified structural requirements.

For each of the five tasks, we generate solutions using every model and prompting style combination across five independent runs. This approach is necessary because large language models (LLMs) are inherently non-deterministic, meaning they can produce different responses to the same prompt due to the stochastic elements in their decoding processes. Conducting multiple independent runs allows us to capture this variability.

287

290

291

292

302

307 308

310

311

312

313

314

315

316

317

320

321

324

325

326

We save each generated output in a JSON file, which includes the graph definition (such as nodes and edges) and any numerical attributes (like costs and trust scores). After saving the output, we use a task-specific verification script to validate the generated graph. This script parses the JSON file into the required Python data structures and checks each constraint. During this process, any errors or constraints that are not met in the output are recorded in a separate JSON file. This file summarizes which constraints were satisfied and explicitly lists any errors made by the model.

We then aggregate these files across the five runs, and look at following metrics:

- **Overall Pass Rate:** The fraction of outputs that satisfy all constraints for a given (model, prompt style) pair.
- Error Breakdown: The frequency of constraint failures in structural vs. logical vs. attribute categories.
- Average Constraint Passing: The average count of successfully met constraints, offering more granularity than a strict pass/fail.

Finally, we compile all verification reports to create a per-run summary of pass/fail outcomes. Another report aggregates the results at the model and prompting method level, computing average pass rates and error counts across the five runs.

3 Results

Our evaluation reveals variations in graph generation capabilities among state-of-the-art language models, providing empirical evidence on the extent to which LLMs are genuinely graph-savvy. The results show critical insights into architectural differences, the efficacy of different prompting strategies, and the distinctive challenges posed by structured graph problems.

3.1 Performance Stratification Across Model Architectures

As shown in Figure 1(c), we observe a pronounced stratification in performance across model families, with specialized reasoning models demonstrating markedly superior capabilities. o3-mini-high and o1 (OpenAI's reasoning-focused models released in January 2025 and December 2024, respectively) achieved exceptional performance with average pass rates of 82.7% and 78.7%, substantially outperforming the cross-model average of 34.0%. Claude 3.7 Sonnet, Anthropic's hybrid reasoning model released in February 2025, followed with a 69.3% success rate, while DeepSeek-R1, another reasoning-specialized architecture, achieved a 48.0% pass rate.

This performance distribution aligns with our hypothesis that graph generation requires sophisticated structural reasoning beyond basic pattern recognition. Notably, the four models explicitly designed with enhanced reasoning capabilities (o3mini-high, o1, Claude 3.7 Sonnet, and DeepSeek-R1) occupy four of the top five positions in overall performance, suggesting that architectural innovations specifically targeting complex reasoning transfer effectively to graph-related tasks.

In contrast, smaller parameter-count models and those without explicit reasoning enhancements struggled significantly. Llama 3.1 (8B) and Llama 3.2 (3B) achieved only 1.3% pass rates, while Chat-GPT 4o-mini reached just 14.7%, indicating fundamental limitations in graph representation abilities. This pattern supports our premise that graph generation constitutes a distinctive challenge requiring specialized architectural capabilities rather than merely scaling parameters.

3.2 Problem-Specific Performance

The performance gradient across tasks remained consistent across model families: **the Time-Dependent Delivery Network** presented the greatest challenge (with error counts averaging 18-49 for most models under direct prompting), followed by **the Gene-Disease Association Network** (10-38 errors). This hierarchy persisted despite iterative feedback, suggesting fundamental differences in task complexity rather than mere prompting limitations. The consistency of this pattern indicates that temporal reasoning with multiple interacting constraints presents a qualitatively different challenge compared to static structural properties.



Figure 1: **Performance analysis of LLMs on graph generation tasks.** Figure panels summarize key trends across fifteen LLMs and five problem domains. (a) Pass rates per model and task reveal that only a few models consistently satisfy all constraints across problems, with stronger results under iterative prompting. (b) Error heatmaps show which models struggle with which types of graphs. (c) Average pass rates across all tasks highlight the performance stratification between reasoning-enhanced and general-purpose models. (d) Performance deltas from iterative feedback quantify each model's ability to self-correct, with Grok-v3 showing the largest improvement.

Error analysis reveals that failures in the Directed Social Network stemmed primarily from specific constraint violations rather than general structural confusion. LLMs mainly struggled with the acyclicity constraint, which involves avoiding prohibited influence cycles, as well as with category-based rules. Claude family showed minimal errors, averaging between 0 and 1, while others, like ChatGPT 40, produced between 6.6 and 7.8 errors under direct prompting, particularly regarding celebrity outgoing edge requirements. Furthermore, specialized reasoning models exhibited a better ability to uphold global structural properties like acyclicity. The deliberately introduced gap in trust score categorization (50-70) shows a consistent tendency across models to hallucinate classifications for these ambiguous values rather than adhering strictly to provided rules. This clas-

390

400

401

402

sification completion bias persisted across multiple prompt iterations especially for simpler models, suggesting an intrinsic tendency to complete perceived patterns rather than strictly adhering to explicit constraints. This is a concerning finding for domain applications requiring rigid adherence to rules.

The Gene-Disease Association task shows another structural pattern. Traditional LLMs struggled specifically with maintaining bipartite integrity (creating forbidden gene-gene or diseasedisease connections) and balancing degree constraints simultaneously. Llama 3.1 (405B) generated 35.4 errors on average under direct prompting, with approximately 70% related to bipartite violations and degree constraint failures. Even with iterative feedback, these models continued to generate structurally invalid networks, suggesting a 403

404

523

473

474

fundamental difficulty in conceptualizing strict categorical separation between node types. In contrast, reasoning-specialized models primarily made errors in strength attribute assignments while maintaining valid bipartite structures.

421

422

423

494

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454 455

456

457

458

459

460

461

462

463

464

465

466

467 468

469

470

471

472

For the Quantum Circuit task, lowerperforming models like Llama 3.1 (8B) and DeepSeek-V3 (which recorded 7.4 errors under direct prompting) primarily struggled with gate adjacency requirements and constraints related to layered operations. This led to the creation of technically invalid quantum circuits. In contrast, errors from Claude and OpenAI models focused more on gate optimization and final state compliance. These were more subtle violations that resulted in operationally valid but suboptimal circuits. This pattern suggests a hierarchy in understanding quantum circuits, where basic structural validity must be established before addressing optimization capabilities. The tendency to selectively violate constraints indicates that domain-specific requirements may be overshadowed by more familiar structural patterns, which raises concerns for specialized domain applications.

The Optimal Transportation Network task revealed a distinctive error pattern focusing on costdistance consistency and accessibility requirements. Even models with high overall pass rates struggled with balancing mutually constraining objectives: Llama models generated 27.4-38.8 errors under direct prompting, primarily violating strategic road placement constraints while maintaining valid connectivity. In contrast, reasoning models made significantly fewer errors (0-1.4) and effectively balanced multiple competing constraints. This suggests that multi-objective optimization in graphs represents a distinctive capability of reasoningenhanced architectures that general-purpose models have not yet mastered.

The most pronounced error pattern emerged in **the Time-Dependent Delivery Network task**, where even high-performing models exhibited cascading failure modes. Error analysis reveals that violations typically began with time window inconsistencies that propagated to vehicle capacity and storage compliance failures. Claude 3.7 Sonnet's unusually high error count (49.0) under direct prompting stems primarily from creating temporally impossible delivery sequences that subsequently violated multiple dependent constraints. This suggests that temporal reasoning in graphs triggers a distinctive failure mode where local inconsistencies propagate through interconnected constraint networks.

Furthermore, across multiple problems, we observed that models frequently generated locally valid edges (satisfying pairwise constraints) that violated global structural properties such as acyclic or strong connectivity. This pattern suggests a limitation in maintaining coherent global graph properties while simultaneously satisfying local edge constraints. This finding has significant implications for applications requiring global structural guarantees.

These detailed error patterns across problem domains collectively indicate that graph hallucination is not a uniform phenomenon but manifests differently depending on the structural properties required. Reasoning-enhanced models demonstrate superior constraint reconciliation abilities, particularly for maintaining global structural properties while satisfying local edge constraints, which is a critical capability for real-world graph applications.

3.3 Constraint Satisfaction by Category

Figure 2(e) demonstrates that reasoning-enhanced models (o3-mini-high, o1, Claude 3.7 Sonnet, and DeepSeek-R1) consistently passed 10-12 structural constraints regardless of prompting strategy. This suggests that structural reasoning capabilities emerge as inherent properties of these architectures rather than being prompt-dependent.

Figure 2(f) reveals greater variability in logical constraint satisfaction, with iterative feedback substantially improving performance across most models (e.g., Grok-v3 improving from 11.6 to 14.0). This differential responsiveness suggests that logical constraints, which often require multi-step reasoning about consequences, benefit most from decomposed reasoning in iterative feedback loops, aligning with prior findings on step-by-step reasoning (Jin et al., 2024).

Figure 2(g) reveals that attribute constraints pose a relatively manageable challenge for most models, with top-performing reasoning models like Claude 3.7 Sonnet, o1, and o3-mini-high consistently achieving perfect or near-perfect scores of 9.0 passed constraints. Even models with moderate overall performance generally exhibited strong attribute constraint satisfaction, suggesting that handling spatial, quantitative, and categorical graph properties represents a more tractable aspect of graph generation compared to structural or logical constraints for current LLM architectures.



Figure 2: **Constraint satisfaction and error analysis.** Breakdown of model performance across constraint types and error categories. (e-f-g) show the average number of structural, logical, and attribute constraints passed per model and prompting strategy. Reasoning-enhanced models (e.g., o1, o3-mini-high, Claude 3.7 Sonnet) consistently score higher, especially on logical constraints. (h) displays average error types by model, revealing that Llama models tend to accumulate structural errors, while Claude models exhibit a higher proportion of logical errors. This analysis shows consistent error signatures across architectures and shows that constraint handling is both task-specific and model-dependent.

3.4 The Efficacy of Prompting Paradigms

524

525

529

533

534

535

537

541

As quantified in Figure 1(d), the improvement from direct prompting to iterative feedback varied dramatically across model families. Grok-v3 exhibited a striking 48% absolute increase, while reasoningspecialized models showed more modest gains (16-28%), suggesting these models possess inherent graph reasoning capabilities less dependent on external guidance. Llama models showed minimal improvement (<5%) despite their poor baseline performance, indicating fundamental architectural limitations that feedback alone cannot overcome.

Contrary to our hypothesis, program-augmented prompting, which provided explicit verification code, did not consistently outperform iterative feedback and sometimes produced worse results than direct prompting. This finding challenges assumptions about LLMs' ability to leverage programmatic verification during generation and suggests limitations in code comprehension or selfmonitoring capabilities. The pattern aligns with Zhang et al. (2024)'s findings that code-based methodologies require tight integration with model architecture rather than simply being provided as context.

3.5 Error Patterns

Figure 2(h) shows distinctive error patterns across model families that illuminate the nature of graph hallucination:

We identified two predominant error patterns: (1) models with high structural but low logical errors (Llama family), suggesting fundamental difficulty with graph topology; and (2) models with low structural but moderate logical errors (Claude family), indicating stronger topological understand-

558

542

543

544

545

546

547

548

549

- 5
- 5
- 565 566
- 567
- 5

5

572 573

task types.

4

Discussion

language models:

reasoning are necessary.

575

577 578

579

5

5 5

585 586

588

590

592 593 594

5

596

598 599

> 6 6

6

606

607

tion.

The variable efficacy of prompting strategies

ing but challenges with constraint reasoning. These

distinct profiles suggest different mechanisms un-

derlying graph hallucination across architectures.

remarkably balanced and minimal error profiles across all categories, while Llama models exhib-

ited compounded failures across structural, logi-

cal, and attribute dimensions. Anthropic models

showed moderate but balanced error distributions,

suggesting a more comprehensive but imperfect

graph understanding. These distinctive signatures

indicate that architectural design decisions create

consistent patterns in graph processing capabilities

that transcend individual prompting strategies or

Our thorough evaluation of fifteen advanced lan-

guage models (LLMs) across five different graph

generation tasks provides an insightful answer to

the question: "Are LLMs truly graph-savvy?"

The evidence indicates that the ability to gener-

ate graphs is not consistently found across models.

Instead, it is closely linked to the architectural de-

sign of the models, especially those enhancements

Our findings have several important theoretical

The consistent superiority of reasoning-

implications for the development of graph-capable

enhanced models (o3-mini-high, o1, Claude

3.7 Sonnet, DeepSeek-R1) over larger but

graph reasoning requires specialized architectural

capabilities rather than merely scaling parameters

or training data. This contradicts the notion that

larger models will naturally develop sophisticated

graph reasoning, suggesting instead that architec-

tural innovations specifically targeting complex

lem types challenge the notion of general graph

reasoning capabilities. Models that excelled at op-

timal transportation networks often struggled with

time-dependent delivery networks, suggesting that

LLMs develop domain-specific structural compe-

tencies that transfer imperfectly across problem

domains. This domain-specificity has implications

for applications requiring cross-domain generaliza-

The pronounced performance gaps across prob-

general-purpose architectures indicates

that focus on improving reasoning capabilities.

OpenAI's models (01, 03-mini-high) displayed

across model families indicates that prompting can enhance but not fundamentally transform an architecture's graph processing capabilities, challenging perspectives that view prompting as a substitute for architectural innovation. This suggests that prompting should be viewed as complementary to, rather than a replacement for, architectural improvements. 609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

Despite our comprehensive evaluation, several limitations should be acknowledged. First, our iterative feedback paradigm utilized only a single round of feedback, potentially limiting the improvements possible through iterative correction. Future work could explore multi-step interactive protocols that better leverage the potential of decomposed reasoning to address complex graph constraints. Second, while our five graph problems span diverse domains, they represent only a subset of possible graph structures and constraint types. Expanding the evaluation to include additional problem domains such as knowledge graphs, molecule generation, and program synthesis graphs would provide a more comprehensive assessment of LLMs' graph capabilities. Third, our evaluation focused primarily on constraint satisfaction rather than generative creativity or optimization quality. Future work could explore how models balance adherence to constraints with the generation of novel or optimal graph structures, particularly in open-ended design tasks. Finally, the black-box nature of many commercial LLMs limits our ability to analyze the underlying mechanisms responsible for performance differences. Future research could benefit from more transparent model architectures that enable detailed analysis of how graph structures are represented and manipulated internally. These limitations suggest several promising directions for future research. The development of specialized finetuning approaches for graph-related tasks could address the observed domain transfer limitations. Hybrid architectures that combine LLMs with graph neural networks or constraint satisfaction solvers might use the complementary strengths of different approaches. In conclusion, our findings demonstrate that while recent architectural advances have significantly improved graph generation capabilities, LLMs' graph-savviness remains highly dependent on architectural design, with specialized reasoning capabilities playing a crucial role. Future advances will likely come from architectural innovations specifically targeting structured reasoning rather than simply scaling existing models or refining prompting strategies.

that

References

661

670

671

673

674

694

703

704

706

708

710

711

712

713

714

- Anthropic. 2024a. Claude 3.5 haiku. https://www. anthropic.com/claude/haiku. Large Language Model.
- Anthropic. 2024b. Claude 3.5 sonnet. https:// www.anthropic.com/news/claude-3-5-sonnet. Large Language Model.
- Anthropic. 2024c. Claude 3.7 sonnet. https:// www.anthropic.com/news/claude-3-7-sonnet. Large Language Model.
- Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoefler. 2024. Graph of thoughts: Solving elaborate problems with large language models. Proceedings of the AAAI Conference on Artificial Intelligence, 38(16):17682–17690.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In Advances in Neural Information Processing Systems, volume 33, pages 1877–1901. Curran Associates, Inc.
 - DeepSeek AI. 2024. Deepseek-v3. https:// api-docs.deepseek.com/news/news1226. Large Language Model.
 - DeepSeek AI. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *Preprint*, arXiv:2501.12948.
 - Google. 2024a. Gemini 2.0 flash. https://deepmind. google/technologies/gemini/flash/. Large Language Model.
 - Google. 2024b. Gemini 2.0 pro. https://deepmind. google/technologies/gemini/pro/. Large Language Model.
 - Bowen Jin, Chulin Xie, Jiawei Zhang, Kashob Kumar Roy, Yu Zhang, Zheng Li, Ruirui Li, Xianfeng Tang, Suhang Wang, Yu Meng, and Jiawei Han. 2024.
 Graph chain-of-thought: Augmenting large language models by reasoning on graphs. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 163–184, Bangkok, Thailand. Association for Computational Linguistics.
- Xinguang Li, Jun Zhan, Fuquan Pan, Tong Lv, and Shen Wang. 2023. A multi-objective optimization model of urban passenger transportation structure

| under low-carbon orientation considering participat- | 715 |
|--|------------------|
| ing subjects. <i>Environmental Science and Pollution</i> | 716 |
| <i>Research</i> , 30(54):115839–115854. | 717 |
| Yuhan Li, Zhixun Li, Peisong Wang, Jia Li, Xiangguo | 718 |
| Sun, Hong Cheng, and Jeffrey Xu Yu. 2024. A survey | 719 |
| of graph meets large language model: Progress and | 720 |
| future directions. <i>Preprint</i> , arXiv:2311.12399. | 721 |
| Erwan Le Merrer and Gilles Tredan. 2024. Llms hallu- | 722 |
| cinate graphs too: a structural perspective. <i>Preprint</i> , | 723 |
| arXiv:2409.00159. | 724 |
| Meta AI. 2024a. Llama 3.1. https://ai.meta.com/ | 725 |
| blog/meta-llama-3-1/. Large Language Model. | 726 |
| Meta AI. 2024b. Llama 3.2. | 727 |
| https://ai.meta.com/blog/ | 728 |
| llama-3-2-connect-2024-vision-edge-mobile-dev | ices <i>†</i> 29 |
| Large Language Model. | 730 |
| Ollama. 2025. ollama/ollama: Get up and running | 731 |
| with llama 3.3, deepseek-r1, phi-4, gemma 3, mistral | 732 |
| small 3.1 and other large language models. https: | 733 |
| //github.com/ollama/ollama. Version v0.7.0, ac- | 734 |
| cessed 2025-05-16. | 735 |
| OpenAI. 2024a. Gpt-4o. https://openai.com/ | 736 |
| index/hello-gpt-4o/. Large Language Model. | 737 |
| OpenAI. 2024b. o1. https://openai.com/o1/. | 738 |
| Large Language Model. | 739 |
| OpenAI. 2024c. o3-mini-high. https://openai.com/ | 740 |
| index/openai-o3-mini/. Large Language Model. | 741 |
| Xubin Ren, Jiabin Tang, Dawei Yin, Nitesh Chawla, | 742 |
| and Chao Huang. 2024. A survey of large language | 743 |
| models for graphs. In <i>Proceedings of the 30th ACM</i> | 744 |
| <i>SIGKDD Conference on Knowledge Discovery and</i> | 745 |
| <i>Data Mining</i> , KDD '24, page 6616–6626, New York, | 746 |
| NY, USA. Association for Computing Machinery. | 747 |
| Javier Romero-Alvarez, Jaime Alvarado-Valiente, Jorge | 748 |
| Casco-Seco, Enrique Moguel, Jose Garcia-Alonso, | 749 |
| and Juan M. Murillo. 2024. Scheduling Process of | 750 |
| Quantum Circuits to Optimize Tasks Execution on | 751 |
| Quantum Computers . In 2024 IEEE International | 752 |
| Conference on Quantum Computing and Engineering | 753 |
| (QCE), pages 182–186, Los Alamitos, CA, USA. | 754 |
| IEEE Computer Society. | 755 |
| Carlos Sacristan, Kumiko Samejima, Lorena Andrade | 756 |
| Ruiz, Moonmoon Deb, Maaike L.A. Lambers, Adam | 757 |
| Buckle, Chris A. Brackley, Daniel Robertson, Tet- | 758 |
| suya Hori, Shaun Webb, Robert Kiewisz, Tristan Be- | 759 |
| pler, Eloïse van Kwawegen, Patrik Risteski, Kruno | 760 |
| Vukušić, Iva M. Tolić, Thomas Müller-Reichert, Tat- | 761 |
| suo Fukagawa, Nick Gilbert, Davide Marenduzzo, | 762 |
| William C. Earnshaw, and Geert J.P.L. Kops. 2024. | 763 |
| Vertebrate centromeres in mitosis are functionally | 764 |
| bipartite structures stabilized by cohesin. <i>Cell</i> , | 765 |

766

187(12):3006-3023.e26.

- 767 772 774 775 776 777 778 779 780 781 783 784 789
- 790 791

794 795

796

799

803

810 811

812

- 813
- 814

- S. M Towhidul Islam Tonmoy, S M Mehedi Zaman, Vinija Jain, Anku Rani, Vipula Rawte, Aman Chadha, and Amitava Das. 2024. A comprehensive survey of hallucination mitigation techniques in large language models. Preprint, arXiv:2401.01313.
- Paolo Toth and Daniele Vigo, editors. 2001. The vehicle routing problem. Society for Industrial and Applied Mathematics, USA.
- Heng Wang, Shangbin Feng, Tianxing He, Zhaoxuan Tan, Xiaochuang Han, and Yulia Tsvetkov. 2023. Can language models solve graph problems in natural language? In Thirty-seventh Conference on Neural Information Processing Systems.
- xAI. 2025. Grok-v3. https://x.ai/blog/grok-3. Large Language Model.
- Yang Yao, Xin Wang, Zeyang Zhang, Yijian Qin, Ziwei Zhang, Xu Chu, Yuekui Yang, Wenwu Zhu, and Hong Mei. 2024. Exploring the potential of large language models in graph generation. Preprint, arXiv:2403.14358.
- Shuo Yu, Yingbo Wang, Ruolin Li, Guchun Liu, Yanming Shen, Shaoxiong Ji, Bowen Li, Fengling Han, Xiuzhen Zhang, and Feng Xia. 2025. Graph2text or graph2token: A perspective of large language models for graph learning. Preprint, arXiv:2501.01124.
- Qifan Zhang, Xiaobin Hong, Jianheng Tang, Nuo Chen, Yuhan Li, Wenzhong Li, Jing Tang, and Jia Li. 2024. Gcoder: Improving large language model for generalized graph problem solving. Preprint, arXiv:2410.19084.

Appendix: Graph Generation Problem Statements

This appendix contains the detailed problem statements for the five graph generation tasks used in our evaluation framework.

A.1 Time-Dependent Delivery Network

Problem Description:

Create a delivery network that schedules deliveries across multiple locations using a fleet of vehicles. The network must account for vehicle capacities, location storage capacities, delivery time windows, dynamic travel times, and vehicle speeds to ensure efficient and timely deliveries.

- **Constraints:**
- 1. Locations:
 - Total Locations: 15, labeled from L0 to L14.
 - Attributes:

| – Storage Capacity: Each location | 815 |
|---|-----|
| has a storage capacity specified in | 816 |
| kilograms (kg). Example: L0 has a | 817 |
| capacity of 500 kg. | 818 |
| - Time Window: Each location has a | 819 |
| delivery time window represented as | 820 |
| a list of two integers [start_hour, | 821 |
| end_hour] in 24-hour format. Ex- | 822 |
| ample: L3 has a time window of [9, | 823 |
| 11] corresponding to 09:00-11:00. | 824 |
| 2. Vehicles: | 825 |
| • Total Vehicles: 7, labeled from V1 to V7. | 826 |
| Attributes: | 827 |
| - Capacity: Each vehicle has a spe- | 828 |
| cific capacity in kilograms (kg). Ex- | 829 |
| ample: V1 has a capacity of 100 kg. | 830 |
| - Speed: Each vehicle has a defined | 831 |
| speed in kilometers per hour (km/h). | 832 |
| Example: V1 travels at 60 km/h. | 833 |
| 3. Edges (Routes): | 834 |
| • Definition: Represents travel paths be- | 835 |
| tween two distinct locations. | 836 |
| Attributes: | 837 |
| - From: The starting location ID (e.g., | 838 |
| L0). | 839 |
| - To: The destination location ID (e.g., | 840 |
| L1). | 841 |
| - Base Travel Time: The fundamental | 842 |
| travel time for the route in minutes. | 843 |
| - Hourly Adjustments: A dictionary | 844 |
| where keys are time ranges in the | 845 |
| format "HH-HH" (24-hour format) | 846 |
| and values are additional travel time | 847 |
| in minutes applicable during those | 848 |
| hours. Example: {"8-10": 15} | 849 |
| adds 15 minutes to the base travel | 850 |
| time between 08:00-10:00. | 851 |
| - Maximum Weight Limit: The max- | 852 |
| imum weight a vehicle can carry on | 853 |
| that route in kilograms (kg). | 854 |
| 4. Operational Constraints: | 855 |
| • Storage Capacity Compliance: The | 856 |
| sum of incoming goods to any location | 857 |

858

859

860

861

must not exceed its storage capacity. • Vehicle Capacity Compliance: No vehicle should exceed its capacity on any edge it traverses.

| 862 | • Time Window Compliance: Departures |
|-----|---|
| 863 | and arrivals must respect the time win- |
| 864 | dows of locations. Specifically: |
| 865 | - Departure Time: Must be within |
| 866 | the from location's time window. |
| 867 | - Arrival Time: Must be within the |
| 868 | to location's time window. |
| 869 | - Loading Time: Assume a fixed load- |
| 870 | ing time of 10 minutes at each loca- |
| 871 | tion, which must be accounted for |
| 872 | when scheduling departures. |
| 972 | Required Autput Format: |
| 874 | <pre><format></format></pre> |
| 014 | |
| 875 | A.2 Directed Social Network with Influence |
| 876 | Relationships |
| 877 | Problem Description: |
| 878 | Create a social network graph representing influ- |
| 879 | ence relationships among users. Each user has |
| 880 | specific attributes, and influence connections must |
| 881 | adhere to defined constraints to maintain the in- |
| 882 | tegrity and intended structure of the network. |
| 883 | Constraints: |
| 884 | 1. Users: |
| 885 | • Total of 20 users labeled from U0 to U19. |
| 886 | Each user has a "trust_score" ranging |
| 887 | from 0 to 100. |
| 888 | • Each user belongs to a "category" based |
| 889 | on their trust score: |
| 890 | - "celebrity" (trust_score ≥ 80) |
| 891 | - "expert" ($70 \le \text{trust_score} < 80$) |
| 892 | – "regular" (trust_score < 50) |
| 893 | 2. Edges (Influence Relationships): |
| 894 | • Directed edges where $Ux \rightarrow Uy$ indi- |
| 895 | cates that Ux influences Uv. |
| 896 | • No self-loops: A user cannot influence |
| 897 | themselves. |
| 898 | Category Constraints: |
| 800 | - Celebrities: Must have at least 5 out- |
| 900 | going edges |
| 901 | - Regular Users: Cannot influence ex- |
| 902 | perts. |
| 003 | Granh Structure |
| 503 | The graph must be savelie (no surles |
| 904 | - The graph must be acyclic (no cycles |
| 900 | in infuence relationships). |
| 906 | Required Output Format: |
| | |

A.3 Quantum Circuit Design

Problem Description:

Design a quantum circuit consisting of multiple qubits and quantum gates. The circuit must adhere to specific constraints to ensure proper gate operations, circuit efficiency, and overall functionality. The design should incorporate structural elements like depth and a Directed Acyclic Graph (DAG) while simplifying some of the gate-related rules to enhance accessibility. 908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

Constraints:

| 1. | Qubits | |
|----|--------|--|
| | · · | |

- Total Qubits: 10, labeled from Q0 to Q9.
- Initialization: All qubits must start in the $|0\rangle$ state.

2. Gates:

- Types of Gates to Include:
 - Single-Qubit Gates: Hadamard (H), Pauli-X (X), Pauli-Z (Z)
 Multi-Qubit Gates: Controlled
 - NOT (CNOT), SWAP – Measurement: Measure (Measure)

• Gate Operations:

- Each gate operates on specific qubits at designated times.
- **CNOT Gates:** Must operate on qubits that are not adjacent (e.g., Q0 and Q2 are valid; Q0 and Q1 are invalid).
- SWAP Gates: Must operate between pairs of qubits that have identical gate sequences up to that point.
- Measurements: Each qubit can be measured only once and must be the last operation on that qubit.

• Gate Restrictions:

- Gate Frequency: No single-qubit gate can be applied more than twice consecutively on the same qubit.

3. Circuit Structure:

- The circuit must be a Directed Acyclic Graph (DAG); no repeated times for the same qubit.
- Layered Operations: Gates at the same time step must operate on disjoint sets of qubits (i.e., no two gates at the same time can act on the same qubit).

| • Depth Constraint: The total number of | 3. Degree Constraints: | 1002 |
|--|---|--|
| time steps (layers) must not exceed 30. | • Genes: | 1003 |
| 4. Operational Constraints: | - Each gene must be associated with at | 1004 |
| • Circuit Reversibility: Measurements | least 2 and at most 5 diseases. | 1005 |
| must be the final operations on their re- | • Diseases: | 1006 |
| spective qubits to maintain circuit re- | - Each disease must be associated with | 1007 |
| versibility. | at least 3 and at most 10 genes. | 1008 |
| • Gate Optimization: The circuit should minimize the total number of gates while | 4. Structural Constraints: | 1009 |
| satisfying all other constraints. | • The network must be bipartite; no edges | 1010 |
| • Final State: After all operations, all | should connect nodes within the same | 1011 |
| qubits must either be measured or re- turned to the $ 0\rangle$ state | set (i.e., no gene-gene or disease-disease | 1012 |
| turned to the $ 0\rangle$ state. | associations). | 1013 |
| Required Output Format: | • There should be no duplicate edges (i.e., | 1014 |
| <formal></formal> | each gene-uisease pan is unique). | 1015 |
| A.4 Gene-Disease Association Network | Required Output Format: | 1016 |
| Problem Description: | <format></format> | 1017 |
| Create a bipartite network that models the asso- ciations between genes and diseases. This net- | A.5 Optimal Transportation Network | 1018 |
| work will represent which genes are associated | Problem Description: | 1019 |
| with which diseases, capturing the strength of each | Design an optimal transportation network repre- | 1020 |
| association. The network should adhere to defined | sented as a directed graph where nodes represent | 1021 |
| constraints to ensure biological relevance and struc- | cities and edges represent one-way roads. The net- | 1022 |
| tural integrity. | work must satisfy constraints to ensure efficiency, | 1023 |
| ~ | connectivity robustness and cost-ettectiveness | 102/ |
| Constraints: | connectivity, robustness, and cost-effectiveness. Constraints: | 1024 1025 |
| Constraints: 1. Nodes: | connectivity, robustness, and cost-effectiveness. Constraints: | 1024 1025 |
| Constraints: 1. Nodes: • Genes: | connectivity, robustness, and cost-effectiveness. Constraints: 1. Nodes (Cities): | 1024 1025 1026 |
| Constraints: 1. Nodes: • Genes: – Total of 20 genes labeled from G0 to | connectivity, robustness, and cost-effectiveness. Constraints: 1. Nodes (Cities): • Total: 8, labeled from C0 to C7. | 1024 1025 1026 1027 |
| Constraints: 1. Nodes: • Genes: – Total of 20 genes labeled from G0 to G19. | connectivity, robustness, and cost-effectiveness. Constraints: Nodes (Cities): Total: 8, labeled from C0 to C7. Attributes: | 1024 1025 1026 1027 1028 |
| Constraints: 1. Nodes: • Genes: – Total of 20 genes labeled from G0 to G19. – Each gene has a "name" and a "func- | connectivity, robustness, and cost-effectiveness. Constraints: Nodes (Cities): Total: 8, labeled from C0 to C7. Attributes: Population: Number of inhabitants | 1024 1025 1026 1027 1028 1029 |
| Constraints: 1. Nodes: • Genes: – Total of 20 genes labeled from G0 to G19. – Each gene has a "name" and a "func- tion". | connectivity, robustness, and cost-effectiveness. Constraints: Nodes (Cities): Total: 8, labeled from C0 to C7. Attributes: Population: Number of inhabitants in each city. | 1024 1025 1026 1027 1028 1029 1030 |
| Constraints: 1. Nodes: Genes: Total of 20 genes labeled from G0 to G19. Each gene has a "name" and a "function". Diseases: Total of 20 diseases labeled from D0 | connectivity, robustness, and cost-effectiveness. Constraints: Nodes (Cities): Total: 8, labeled from C0 to C7. Attributes: Population: Number of inhabitants in each city. C0: 1,000 | 1024 1025 1026 1027 1028 1029 1030 1031 |
| Constraints: 1. Nodes: • Genes: – Total of 20 genes labeled from G0 to G19. – Each gene has a "name" and a "func- tion". • Diseases: – Total of 20 diseases labeled from D0 to D19 | connectivity, robustness, and cost-effectiveness. Constraints: Nodes (Cities): Total: 8, labeled from C0 to C7. Attributes: Population: Number of inhabitants in each city. C0: 1,000 C1: 500 | 1024 1025 1026 1027 1028 1029 1030 1031 1032 |
| Constraints: 1. Nodes: Genes: Total of 20 genes labeled from G0 to G19. Each gene has a "name" and a "function". Diseases: Total of 20 diseases labeled from D0 to D19. Each disease has a "name" and | connectivity, robustness, and cost-effectiveness. Constraints: Nodes (Cities): Total: 8, labeled from C0 to C7. Attributes: Population: Number of inhabitants in each city. C0: 1,000 C1: 500 C2: 750 | 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 |
| Constraints: 1. Nodes: Genes: Total of 20 genes labeled from G0 to G19. Each gene has a "name" and a "function". Diseases: Total of 20 diseases labeled from D0 to D19. Each disease has a "name" and a "severity level" (e.g., "Low". | connectivity, robustness, and cost-effectiveness. Constraints: Nodes (Cities): Total: 8, labeled from C0 to C7. Attributes: Population: Number of inhabitants in each city. C0: 1,000 C1: 500 C2: 750 C3: 600 | 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 |
| Constraints: 1. Nodes: Genes: Total of 20 genes labeled from G0 to G19. Each gene has a "name" and a "function". Diseases: Total of 20 diseases labeled from D0 to D19. Each disease has a "name" and a "severity_level" (e.g., "Low", "Medium", "High"). | connectivity, robustness, and cost-effectiveness. Constraints: Nodes (Cities): Total: 8, labeled from C0 to C7. Attributes: Population: Number of inhabitants in each city. C0: 1,000 C1: 500 C2: 750 C3: 600 C4: 900 | 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 |
| Constraints: 1. Nodes: Genes: Total of 20 genes labeled from G0 to G19. Each gene has a "name" and a "function". Diseases: Total of 20 diseases labeled from D0 to D19. Each disease has a "name" and a "severity_level" (e.g., "Low", "Medium", "High"). | connectivity, robustness, and cost-effectiveness. Constraints: Nodes (Cities): Total: 8, labeled from C0 to C7. Attributes: Population: Number of inhabitants in each city. C0: 1,000 C1: 500 C2: 750 C3: 600 C4: 900 C5: 400 | 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036 |
| Constraints: 1. Nodes: Genes: Total of 20 genes labeled from G0 to G19. Each gene has a "name" and a "function". Diseases: Total of 20 diseases labeled from D0 to D19. Each disease has a "name" and a "severity_level" (e.g., "Low", "Medium", "High"). 2. Edges (Associations): | connectivity, robustness, and cost-effectiveness. Constraints: Nodes (Cities): Total: 8, labeled from C0 to C7. Attributes: Population: Number of inhabitants in each city. C0: 1,000 C1: 500 C2: 750 C3: 600 C4: 900 C5: 400 C6: 800 | 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036 1037 |
| Constraints: 1. Nodes: Genes: Total of 20 genes labeled from G0 to G19. Each gene has a "name" and a "function". Diseases: Total of 20 diseases labeled from D0 to D19. Each disease has a "name" and a "severity_level" (e.g., "Low", "Medium", "High"). 2. Edges (Associations): Represents the association between a | connectivity, robustness, and cost-effectiveness. Constraints: Nodes (Cities): Total: 8, labeled from C0 to C7. Attributes: Population: Number of inhabitants in each city. C0: 1,000 C1: 500 C2: 750 C3: 600 C4: 900 C5: 400 C6: 800 C7: 650 | 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036 1037 1038 |
| Constraints: 1. Nodes: Genes: Total of 20 genes labeled from G0 to G19. Each gene has a "name" and a "function". Diseases: Total of 20 diseases labeled from D0 to D19. Each disease has a "name" and a "severity_level" (e.g., "Low", "Medium", "High"). 2. Edges (Associations): Represents the association between a gene and a disease. Bineratite Constraints Association periods | connectivity, robustness, and cost-effectiveness. Constraints: Nodes (Cities): Total: 8, labeled from C0 to C7. Attributes: Population: Number of inhabitants in each city. C0: 1,000 C1: 500 C2: 750 C3: 600 C4: 900 C5: 400 C5: 400 C6: 800 C7: 650 2. Edges (Roads): | 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036 1037 1038 1039 |
| Constraints: 1. Nodes: Genes: Total of 20 genes labeled from G0 to G19. Each gene has a "name" and a "function". Diseases: Total of 20 diseases labeled from D0 to D19. Each disease has a "name" and a "severity_level" (e.g., "Low", "Medium", "High"). 2. Edges (Associations): Represents the association between a gene and a disease. Bipartite Constraint: Associations can only exist between genes and diseases | connectivity, robustness, and cost-effectiveness. Constraints: Nodes (Cities): Total: 8, labeled from C0 to C7. Attributes: Population: Number of inhabitants in each city. C0: 1,000 C1: 500 C2: 750 C3: 600 C4: 900 C5: 400 C6: 800 C7: 650 2. Edges (Roads): Definition: Represents a one-way road | 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036 1037 1038 1039 1040 |
| Constraints: 1. Nodes: Genes: Total of 20 genes labeled from G0 to G19. Each gene has a "name" and a "function". Diseases: Total of 20 diseases labeled from D0 to D19. Each disease has a "name" and a "severity_level" (e.g., "Low", "Medium", "High"). 2. Edges (Associations): Represents the association between a gene and a disease. Bipartite Constraint: Associations can only exist between genes and diseases, not within the same set | connectivity, robustness, and cost-effectiveness. Constraints: Nodes (Cities): Total: 8, labeled from C0 to C7. Attributes: Population: Number of inhabitants in each city. C0: 1,000 C1: 500 C2: 750 C3: 600 C4: 900 C5: 400 C5: 400 C6: 800 C7: 650 2. Edges (Roads): Definition: Represents a one-way road from one city to another. | 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036 1037 1038 1039 1040 1041 |
| Constraints: 1. Nodes: Genes: Total of 20 genes labeled from G0 to G19. Each gene has a "name" and a "function". Diseases: Total of 20 diseases labeled from D0 to D19. Each disease has a "name" and a "severity_level" (e.g., "Low", "Medium", "High"). 2. Edges (Associations): Represents the association between a gene and a disease. Bipartite Constraint: Associations can only exist between genes and diseases, not within the same set. Association Strength: Each association | connectivity, robustness, and cost-effectiveness. Constraints: Nodes (Cities): Total: 8, labeled from C0 to C7. Attributes: Population: Number of inhabitants in each city. C0: 1,000 C1: 500 C2: 750 C3: 600 C4: 900 C5: 400 C6: 800 C7: 650 2. Edges (Roads): Definition: Represents a one-way road from one city to another. Attributes: | 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036 1037 1038 1039 1040 1041 1042 |
| Constraints: 1. Nodes: Genes: Total of 20 genes labeled from G0 to G19. Each gene has a "name" and a "function". Diseases: Total of 20 diseases labeled from D0 to D19. Each disease has a "name" and a "severity_level" (e.g., "Low", "Medium", "High"). Edges (Associations): Represents the association between a gene and a disease. Bipartite Constraint: Associations can only exist between genes and diseases, not within the same set. Association Strength: Each association has a "strength" value ranging from 0.0 | connectivity, robustness, and cost-effectiveness. Constraints: Nodes (Cities): Total: 8, labeled from C0 to C7. Attributes: Population: Number of inhabitants in each city. C0: 1,000 C1: 500 C2: 750 C3: 600 C4: 900 C5: 400 C5: 400 C6: 800 C7: 650 2. Edges (Roads): Definition: Represents a one-way road from one city to another. Attributes: Distance: Length of the road in kilo- | 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036 1037 1038 1039 1040 1041 1042 1043 |
| Constraints: 1. Nodes: Genes: Total of 20 genes labeled from G0 to G19. Each gene has a "name" and a "function". Diseases: Total of 20 diseases labeled from D0 to D19. Each disease has a "name" and a "severity_level" (e.g., "Low", "Medium", "High"). 2. Edges (Associations): Represents the association between a gene and a disease. Bipartite Constraint: Associations can only exist between genes and diseases, not within the same set. Association Strength: Each association has a "strength" value ranging from 0.0 to 1.0, indicating the confidence of the | connectivity, robustness, and cost-effectiveness. Constraints: Nodes (Cities): Total: 8, labeled from C0 to C7. Attributes: Population: Number of inhabitants in each city. C0: 1,000 C1: 500 C2: 750 C3: 600 C4: 900 C5: 400 C5: 400 C6: 800 C7: 650 2. Edges (Roads): Definition: Represents a one-way road from one city to another. Attributes: Distance: Length of the road in kilometers (km). (Each road must be ≤ | 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036 1037 1038 1039 1040 1041 1042 1043 1044 |
| Constraints: 1. Nodes: Genes: Total of 20 genes labeled from G0 to G19. Each gene has a "name" and a "function". Diseases: Total of 20 diseases labeled from D0 to D19. Each disease has a "name" and a "severity_level" (e.g., "Low", "Medium", "High"). 2. Edges (Associations): Represents the association between a gene and a disease. Bipartite Constraint: Associations can only exist between genes and diseases, not within the same set. Association Strength: Each association has a "strength" value ranging from 0.0 to 1.0, indicating the confidence of the association. | connectivity, robustness, and cost-effectiveness. Constraints: Nodes (Cities): Total: 8, labeled from C0 to C7. Attributes: Population: Number of inhabitants in each city. C0: 1,000 C1: 500 C2: 750 C3: 600 C2: 750 C3: 600 C4: 900 C5: 400 C6: 800 C7: 650 2. Edges (Roads): Definition: Represents a one-way road from one city to another. Attributes: Distance: Length of the road in kilometers (km). (Each road must be ≤ 300 km.) | 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036 1037 1038 1039 1040 1041 1042 1043 1044 1045 |

| 1046 1047 | - Construction Cost: Cost to build the road in thousand dollars (\$K). |
|--------------|--|
| 1048 3. Ad | ditional Constraints: |
| 1049 (a |) Connectivity: The network must be |
| 1050 | strongly connected, meaning there is a |
| 1051 | directed path from any city to every other |
| 1052 | city. |
| 1053 (b |) Road Capacity: No single road should |
| 1054 | be longer than 300 km . |
| 1055 (c |) Cost Optimization: The total construc- |
| 1056 | tion cost of all roads should not exceed |
| 1057 | \$10,000K. |
| 1058 (d |) Population Accessibility: Each city |
| 1059 | must have at least two incoming roads |
| 1060 | to ensure redundancy and accessibility. |
| 1061 (e |) Strategic Road Placement: Cities C0 |
| 1062 | and C7 are major hubs and must have |
| 1063 | at least three outgoing roads each to |
| 1064 | distribute traffic efficiently. |
| 1065 (f |) Avoiding Redundancy: No two cities |
| 1066 | should have more than one direct road |
| 1067 | connecting them in the same direction. |
| 1068 (g |) Minimizing Total Distance: The sum |
| 1069 | of all road distances should be mini- |
| 1070 | mized to ensure efficient transportation. |
| 1071 (h |) 2-Edge Robustness: The network must |
| 1072 | remain strongly connected if any sin- |
| 1073 | gle road is removed (i.e., there must be |
| 1074 | two edge-disjoint paths between every |
| 1075 | ordered pair of cities). |
| 1076 (1 |) Edge-Disjoint Paths Guarantee: For |
| 1077 | every pair of distinct cities, there must |
| 1078 | dent (adge disjoint) noths connecting |
| 1079 | them |
| 1091 (; |) Balanced Autacing Degrees Except for |
| 1082 | the designated hubs (C0 and C7) the dif- |
| 1083 | ference between the maximum and mini- |
| 1084 | mum number of outgoing roads among |
| 1085 | all cities must not exceed 2 . This pre- |
| 1086 | vents "overloaded" junctions. |
| 1087 (k |) Path Efficiency Constraint: For every |
| 1088 | pair of cities, the shortest route (by total |
| 1089 | distance) should be less than 500 km to |
| 1090 | ensure quick intercity transit. |
| 1091 (1 |) Cost–Distance Consistency: For every |
| 1092 | road, the construction cost (in \$K) must |
| 1093 | be between 1.0 and 1.5 times its dis- |
| 1094 | tance (in km). <i>Example:</i> A road that |
| | |

| is 90 km long must have a cost between | 1095 |
|---|------|
| 90K and 135K . | 1096 |
| (m) Maximum Edge-Hop Constraint: For | 1097 |
| every pair of cities, you need to be able | 1098 |
| to get to every other city in at most 3 | 1099 |
| edges. | 1100 |
| Required Output Format: | 1101 |
| <format></format> | 1102 |