Graph-O-Planner: Injecting Graph Neural Tool Embeddings into LLMs for Efficient and Accurate Task Execution

Anonymous ACL submission

Abstract

001

003

800

011

012

014

017

018

027

037

Recent advancements in Large Language Models (LLMs) have enabled the development of AI agents capable of multi-step reasoning. However, deploying these agents in real-world applications requires planners that adapt to domainspecific tools and workflows, where traditional prompting frameworks often struggle to accurately represent available functional dependencies. To address this gap, we propose Graph-O-Planner, a novel graph-learning method that explicitly encodes tool relationships and execution sequences into LLM planning. Our approach constructs graph embeddings of available tools, enabling agents to dynamically map dependencies while minimizing context window overload. Evaluations across multiple benchmarks, including UltraTool and Task Bench, demonstrate that Graph-O-Planner achieves up to 68% higher and 60% higher performance with our approach, compared to state-of-the-art hybrid graph+LLM based planners and LLM-finetuned planners respectively, while significantly reducing any hallucinations in LLM generation. The method's tool knowledge compression further reduces inference latency by 20%, validating its effectiveness in resource-constrained environments and making it more compatible for real-life practical deployment. We release our code here¹.

1 Introduction

The advent of Large Language Model (LLM)powered agents marks a paradigm shift in artificial intelligence, with transformative potential for real-world applications ranging from autonomous robotics to precision medicine. Early implementations like HuggingGPT (Shen et al., 2023a) demonstrate problem-solving flexibility in controlled benchmarks, and agents such as Voyager (Wang et al., 2023a) showcase emergent strategic reasoning in gaming environments. Effective planning modules with precise tool alignment are essential for developing practical AI agentic systems in both consumer and industrial applications. Recent advances leverage prompting strategies to decompose complex tasks: Wei et al. (2022); Yu et al. (2025) pioneered Chain-of-Thought(CoT) reasoning through sequential step generation, while Wang et al. (2023b) introduced plan-and-solve prompting for systematic task decomposition. Yao et al. (2024) later expanded these concepts with tree-based reasoning architectures. A parallel research trajectory has focused on translating these reasoning structures into executable tool operations (Schick et al., 2023; Shen et al., 2023b; Singh et al., 2023; Song et al., 2023). 042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

077

081

These methods, however, are purely promptbased and use LLMs deployed on the cloud. This centralization exposes sensitive data to privacy risks, incurs network-dependent latency, and complicates regulatory compliance. Another focus has been locally deployable solutions (Erdogan et al., 2024; Wu et al., 2024) where the model is finetuned for planning and tool-calling. Locally deploying preserves data privacy and minimizes inference delay, yet small to mid-sized models introduce their own constraints. We describe these constraints:

- 1. *Challenge 1.* Tool Hallucination/ Tool Grounding significantly drives down task-sequence determination, especially with a bigger set of tools/sub-tasks (CodeLlama13B and GPT3.5 Turbo see 60% hallucination in edge prediction with a set of 260 sub-tasks) (Wu et al., 2024).
- 2. *Challenge 2.* Fixed Context: The limited context window of LLMs prevents them from handling a large number of tools, often forcing tools to truncate from the context.
- 3. *Challenge 3.* Token Overhead and Latency Issues: Tool usage in current systems involves repeatedly injecting detailed descriptions of

¹https://anonymous.4open.science/r/Graph-O-Planner-16B3

- 100 101

099

- 102
- 104
- 105 106

107

109

- 110 111 112
- 113

114 115

116 117

118

119

120

122

123

124

125

126

127

128

130

2

2.1 GNN-based Learning

Related Work

lines by 35%.

Recent works employ task graphs to model interlinked sub-tasks and align LLMs with tool-relevant information. GNNs demonstrate strong capabilities for complex decision-making (Khalil et al., 2017; Xu et al., 2019; Dudzik and Veličković, 2022). Graph-based QA systems have successfully

tools into prompts. This token increase leads

to quadratic increase in memory requirements

due to attention, impacting latency during in-

4. Challenge 4. Shallow GNN+LLM Fusion:

Existing GNN+LLM hybrid lacks deep struc-

tural grounding, limiting planning quality and

We draw upon these insights to propose a frame-

work that effectively uses GNN for tool information

injection during task-planning. Building upon the

initial work done by Wu et al. (2024), we propose

Graph-O-Planner, that represents tools as nodes

in a dynamic graph, leveraging Graph Neural Net-

works (GNNs) to capture functional dependencies

and enable scalable tool representation. By fusing

GNN with Large Language Models (LLMs), we

create a more grounded and hallucination-resistant

decision-making process while reducing latency

Our main contributions are summarized as:

• To the best of our knowledge, Graph-O-

Planner is the first attempt to deeply integrate

multilevel interaction of a GNN with LLMs

for task planning. This setup uses the lan-

guage understanding skills of LLMs in con-

junction with the effective information propa-

• We show the efficacy of adding GNN based

interaction by comparing against LLM-only

models. Graph-O-Planner improves upon the

finetuned LLM performance by nearly 60%

while beating hybrid LLM+graph-based base-

• By incorporating dependency-aware tool

graph reduces tool hallucination metrics by

13%, ensuring that model generates correct

tool names instead of similar looking tool

names. It further reduces the context tokens

which is a major bottleneck for local deploy-

ment and speeds up inference by 1.25x.

gation capability of GNNs.

ference

scalability.

and memory requirements.

leveraged external KGs for factual queries through LLMs, using relevant subgraph retrieval (Luo et al., 2023; Zhang et al., 2023, 2024; Yao et al., 2022). GNNs have also been shown to enhance LLMs' ability to model textual relational structures. While transformers (Guo et al., 2025; Chung et al., 2022; Dubey et al., 2024; Yang et al., 2024) dominate sequential processing, they falter with long-range dependencies like tool relationships. GNNs overcome this by representing text as graphs (Huang et al., 2019; Pham et al., 2023; Zhu et al., 2021), improving LLMs' understanding of local and global patterns (Wu et al., 2024, 2021). 131

132

133

134

135

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

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

2.2 Tool Graph-based Planning

Taking inspiration from KG based learning(Liu et al., 2021; Wang et al., 2021; Ye et al., 2022; Tena Cucala et al., 2022; Chen et al., 2020), training GNN using task graphs has become a powerful tool for task planning, enabling the modeling of complex dependencies between tasks, resources, and constraints. They have been applied across diverse domains, including MoE task planning (Zhou et al., 2022; Cai et al., 2024; Li et al., 2025) and multi-agent coordination (Wu et al., 2023; Chan et al., 2023; Talebirad and Nadiri, 2023; Nascimento et al., 2023) by facilitating decentralized decision making.

3 Preliminaries

In this section, we describe tool graphs, which we define as dynamically changing graphs. Subsequent paragraphs provide a detailed explanation of tool graphs, including their structure, tool descriptions & input/output formats.

Tool Graph: Let G = (V, E, A, T, X) where, V is a set of tool nodes, E corresponds to edges between node embedding $(v_i \rightarrow v_j)$ if output of V_i can be fed to V_j . A denotes the edge weight matrix between pair of nodes, such that $A[i, j] \in (0, 1]$, if $v_i, v_j \in V$ and $(v_i, v_j) = e_{ij} \in E$, and 0 otherwise. Tool information is defined as $T = \{n, d, i, o\}_{(k=1)}^{|V|}$, where n is k^{th} tool name, d is k^{th} tool description, i and o corresponds to k^{th} input and output format of tool respectively. Thus, X = Emb(T), where Emb is the embedding function. $X = \{x_i\}_{(k=1)}^{|V|}$ contains feature embedding of tool's information for each $v_i \in V$.

Planning task definition: A task query Q can be decomposed into sub-tasks $S = s_1, s_2...s_n$, such that each sub-task s_i can be completed by an unique tool v_i . The objective is to construct



Figure 1: Overall schematic block diagram of Graph-O-Planner including data flow, key component and internal interactions. Yellow layers in the diagram represent trainable parameters, while purple indicates frozen parameters.



Figure 2: Creation of dependency-aware tool embeddings from GNN

a Directed Acyclic Graph (G) that represents the
sequence of processes to solve query Q. This can
be formally represented as:

$$DAG = (v_1 \rightarrow v_2 \rightarrow \ldots \rightarrow v_P)$$

showcasing the sequence of tools that needs to be executed in order to solve the query Q. Each tool v_i is connected to the next tool v_{i+1} through a dependency edge.

4 Graph-O-Planner

185

We introduce Graph-O-Planner, a novel graph-189 infused learning framework for task planning, to effectively align task steps to the available tools. 191 Figure 1 illustrates the pipeline of our proposed 192 approach. We first convert the target dataset into an aligned tool-graph, incorporating the tool name, 195 description, required inputs and generated output (Section 4.1). This is then passed to a GNN to cre-196 ate its aligned node embeddings & edge scores at 197 each layer, as defined in Section 3. Finally, a layerwise knowledge injection method is utilized to in-199

ject knowledge from GNN layers to corresponding 200 LLM layers (Section 4.3). This allows the LLM 201 to effectively map the subtask to the correct tool 202 sequence. The rough decomposed plan (attained 203 from any global LLM) is processed sequentially 204 through the LLM layers with injected information 205 from the GNN to generate a sequence of tools as a 206 DAG to be executed. With Graph-O-Planner, both 207 LLMs and GNNs are trained in alignment with in-208 jected knowledge, to allow the model to build an 209 understanding of how a subtask relates to a particu-210 lar sequence of nodes in the available tool graph. 211

4.1 Tool Graph Creation

We first encode dataset's Tool information using213a text encoder. We chose ModernBERT Large214(Warner et al., 2024) which outperforms other215encoder-based language models on Natural Language Understanding (NLU) and retrieval tasks.217Using the encoder model we generate a tool's em-218

239

240

241

242

243

244

246

247

248

253

254

260

262

264

bedding representation x.

$$x_{i} = Emb(ToolName, ToolDesc., ToolInputs, ToolOutput)$$
(1)

where x_i is i^{th} tool embedding and Emb is embedding model. For a given set of task steps, some 223 nodes and edges are more semantically relevant than others. To effectively model this by leveraging core semantic information, the edge connections between the tools are scored. Alignment with the requirements of the decomposed task steps is 227 obtained by using a bilinear layer to estimate the relevance score for each edge given the embedding of the sub-task. Finally, these embeddings are normalized to increase computational efficiency. Motivated by Jang et al. (2017), we used Gumbel softmax approach to model the output as soft labels with a stop gradient mechanism to address the problem of gradient propagation of hard labels dur-235 ing backward pass. The node and the scored edge 236 embeddings together comprise the required Tool Graph (defined in Section 3).

4.2 Dependency-Aware Tool Graph

The scored tool graph is then passed through a graph network, to obtain its graphical embedding representation at every layer. This is pictorially shown in Figure 2. We use Graph Attention (GAT) layers for encoding the tool info graph representations $X = \{x_1x_2, \ldots, x_n\}$, via iterative convolutional operations between neighboring nodes of the graph network.

1. Edge Encoding

Given a graph, G = (V, E, A, T, X) (refer section 3), with node features $h_i^l \in R^d$ (initially $h_i^0 = x_i, x_i \in X$), and auxiliary node features $\phi \in R^{|K| \times d}$, our goal is to learn node representations that captures structural neighborhood patterns and edge semantic relationship. For each edge $e_{ij} \in E$, we obtain its encoding ϵ_{ij} as:

$$\epsilon_{ij} = f_{edge}(\phi_{e_{ij}} \oplus \tau_{k_i} \oplus \tau_{k_j}) \qquad (2)$$

where f_{edge} is a multi-layer MLP, $\phi_{e_{ij}}$ represents the Gumbel Softmax of edge e_{ij} , and τ_{k_i} and τ_{k_j} represents the Gumbel Softmax of node k_i and k_j respectively, with \oplus as the concatenate function.

2. Multi-Head Heterogeneous Attention

Each edge's representation ϵ_{ij} is then uti-

lized to transform node representations using a multi-head heterogeneous attention, M. Specifically, for each graph node, we obtain the normalized attention of the pre-head computations. The node embedding \tilde{h}_i is obtained by concatenating (\oplus) node features and its features of its neighboring edges N(i), to better capture neighborhood patterns every iteration:

$$\tilde{h_i} = h_i \oplus \bigoplus_{j \in N(i)} \epsilon_{ij} \tag{3}$$

265

266

267

269

270

271

272

273

274

275

276

277

278 279

281

282

283

284

289

291

292

293

294

295

297

298

299

300

301

302

303

305

306

307

308

The pre-heads of Key, Query and Value for the k-th node are thus computed as:

$$Q_i^k = W_Q^k \tilde{h_i} \tag{4}$$

$$K_i^k = W_K^k \tilde{h_i} \tag{5}$$

$$Q_i^k = W_V^k \tilde{h_i} \tag{6}$$

which is then normalized based on the degree of the node:

$$\alpha_{ij}^k = \frac{d_j}{Z_i} * \exp \frac{\langle q_j^k, k_i^k \rangle}{\sqrt{d}} \tag{7}$$

where $Z_i = \sum_{j \in N(i)} d_j * \exp \frac{\langle q_j^k, k_i^k \rangle}{\sqrt{d}}$ and $d_i = |N(i)|$ is the out degree of node j.

3. Message Aggregation

The multi head attention output are aggregated through a two stage process to synthesize neighborhood information. First for each attention head k, messages from neighboring nodes are weighted by their normalized coefficients α_{ij}^k producing head-specific representations

$$m_i^k = \sum_{j \in N(i)} (\alpha_{ij}^k v_{ij}^k) \tag{8}$$

These head embeddings are then concatenated across all k heads and linearly projected to the original dimension d using learnable weights W_o to obtain M - edge aware attention. This ensures structured fusion of heterogeneous relation patterns. The hybrid approach retains structural information while allowing the model to learn incremental feature updates, balancing neighborhood influence with nodespecific characteristics.

Finally, we use a three phase computation over the edge-attention M to obtain the embeddings for every time-stamp:

$$H^{l+1} = MLP(Agg(M(H^l, E, \epsilon)))$$
(9)

We use an additional node, master node which connects to all the nodes in the graph to pool the repre-310 sentations of all the nodes and obtain representation 311 of the graph. 312

4.3 Deep Fusion of GNN with LLM 313

314

317

319

325

326

334

336

341

347

Finally, we inject the GNN knowledge into LLM layers for effective tool-aligned subtask creation. 315 For a chosen subset of LLM layers, the standard multi-head self-attention and feed-forward layer is extended to fuse the modalities between text and graph domain. The pooled representation of the graph embeddings i.e. representation of the master node H from equation 9 is fused by concatenating (\oplus) with LLM decoder module to obtain an intermediate representation I_c formulated as follows:

$$I_{c} = \{\theta_{key} \oplus \psi_{key}, \theta_{query} \oplus \psi_{query}, \\ \theta_{value} \oplus \psi_{value}\}$$
(10)

where ψ_i are aligned values obtained from GNN after passing through two Feed Forward Layers. The injection pipeline is aligned with the LoRA layers of the LLM to avoid re-training of the complete LLM. As the knowledge injection methodology only needs aligned layer output fusion, it is independent of the architecture of the LLM in use.

5 Experimental Setup and Results

In this section, we describe the training setup and datasets used. We also enumerate a detailed set of experiments to evaluate the performance of Graph-O-Planner against state-of-the-art graphbased baselines and LLMs finetuned on four opensource datasets. For our primary pipeline, we choose Flan-T5-XL (3B) as the LLM with a LoRA of rank 8 and alpha 16. All models are trained using Adam Optimizer, with the batch size of 32 and learning rate of 1e - 4 and 3e - 4 for the LLM and GNN respectively. All models are trained using A-6000 GPU in a Pytorch framework and CUDA 12.6. We provide specific hardware and software versions information in appendix A.2

5.1 Datasets

We train and evaluate our model on four open source datasets, UltraTool (Huang et al., 2024), HuggingFace, Multimedia and TaskBench-Daily Life (Shen et al., 2023b). Each dataset is converted 351 to a Tool Graph as described in Section 4.1. Detailed dataset description has been provided in the Appendix A.6. 354



Figure 3: Statistics of tool graph and samples of datasets.

Results & Analysis 5.2

We present our model performance across accuracy, hallucination and latency metrics, detailed in the section below. In all experiments, our base pipeline uses FlanT5-XL as LLM and a GAT GNN. We compare our model performance against various existing SOTA models. We also demonstrate the impact of a graph based knowledge injection by comparing the performance against traditional LLM only approaches, shallow GNN interaction methods and reasoning only based settings. We report tool and sequence performances along with hallucination and model latency for all four target datasets in below sections.

355

356

357

358

359

360

361

362

363

364

365

366

369

370

371

372

373

374

375

376

377

379

380

381

384

385

386

387

389

5.2.1 **Tool and Sequence Detection**

The primary requirement of a tool aligned planner is to ensure that the model is able to correctly pick from the provided set of available tools. A planner that provides a "somewhat-correct" output is not scalable in a real-life application. We thus measure node prediction F1-score as a primary metric of performance evaluation. Node-F1 score is estimated as the correctness of predicted tool nodes required to complete a given task.

Another crucial performance metric is to ensure that the correct nodes are detected in the correct sequence. Thus, the edges between various task nodes and their relative sequence is of utmost importance. Accordingly, we design Edge-F1 score which compares the predicted links with the ground truth edges, using the tool network topology populated adjacency matrices. We present the edge and node F1 algorithm in Appendix A.9.

For F1 scores, the set of predicted nodes/edges and set of ground truth nodes/edges is used for each

Method/Dataset	Huggingface		Ultratool		Multimedia			Dailylife				
	Acc.	Edge-F1	Node-F1	Acc.	Edge-F1	Node-F1	Acc.	Edge-F1	Node-F1	Acc.	Edge-F1	Node-F1
Mistral7b + GraphToken*	20.08	32.55	62.15	64.57	68.57	85.63	35.06	74.57	63.71	69.42	73.57	92.50
e5-335M + GraphSAGE*	24.74	40.60	66.52	37.56	42.36	71.23	62.37	70.25	88.86	86.57	85.80	97.42
e5-335M + GCN*	16.36	30.23	60.06	46.23	32.54	63.54	61.25	50.76	73.34	75.49	65.49	86.39
e5-335M + GIN*	23.31	39.37	65.77	57.84	37.25	69.25	62.55	69.84	88.74	84.49	82.65	93.36
e5-335M + GAT*	23.39	40.66	66.44	60.25	40.23	80.23	63.69	70.24	88.90	89.91	89.32	97.21
e5-335M + GraphTransformer*	24.54	41.64	66.92	46.42	57.68	76.38	NA	70.24	88.90	NA	85.80	97.43
Qwen 2.5 Coder 3B [‡]	81.32	68.23	<u>87.77</u>	74.80	78.57	91.87	21.48	26.96	56.94	89.86	91.83	99.24
Deepseek R1 1.5B ‡	79.28	80.22	84.33	85.28	89.27	87.33	82.25	83.24	81.27	87.23	88.23	86.24
Flan T5 XL‡	49.0	63.48	74.11	73.8	73.87	83.87	43.40	58.25	58.25	63.87	48.77	66.45
RAG + Qwen 2.5 Coder 3B (not fintuned) ψ	32.8	42.86	40.09	56.8	46.84	38.36	30	32.89	36.65	30.13	32.56	36.36
RAG + Qwen 2.5 Coder 3B (finetuned) ψ	75.00	76.33	78.26	<u>95.56</u>	<u>95.87</u>	<u>96.56</u>	75.6	76.82	78.76	<u>96.8</u>	96.00	96.41
Graph-O-Planner(Ours)	82.6	89.73	97.36	97.39	97.73	99.19	92.0	93.55	99.13	96.81	<u>95.59</u>	<u>99.82</u>

Table 1: Accuracy, Edge-F1 and Node-F1 scores across all four datasets. All result are in (%) with the best bold and runner-up underlined. * indicates hybrid LLM+graph based approaches and that the numbers are sourced from Wu et al. (2024). \ddagger indicates LLM-only models finetuned for each dataset. ψ indicated Retrieval Augmented Generation (RAG) based results. The information about the baseline models are described in Appendix A.7.

sample *i* among *N* total samples.

$$Precision = \frac{1}{N} \sum_{i=1}^{N} \frac{|Predicted_i \cap Ground Truth_i|}{|Predicted_i|}$$
(11)

$$\operatorname{Recall} = \frac{1}{N} \sum_{i=1}^{N} \frac{|\operatorname{Predicted}_{i} \cap \operatorname{Ground} \operatorname{Truth}_{i}|}{|\operatorname{Ground} \operatorname{Truth}_{i}|}$$
(12)

$$F1 \text{ Score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$
(13)

We also evaluate the success rate at the task level using the Accuracy (Acc) metric. The Accuracy metric is defined as:

Accuracy =
$$\frac{1}{N} \sum_{i=1}^{N} I(F1_i = 1)$$
 (14)

Here, $I(F1_i = 1)$ is an indicator function that assigns a value of 1 if the F1 score for a given task *i* is perfect (i.e., all nodes are correctly predicted), and 0 otherwise. This binary evaluation allows us to assess the model's ability to accurately predict all nodes in a task.

Table 1 shows our performance across four different datasets, which included a variety of task types and complexities. This suggests that our method is flexible and can be applied to different problems and datasets. Specifically, across the more voluminous **ultratool dataset** (with 260+ tools), Graph-O-Planner shows a **33% improvement** over older graph interaction methods and 8% improvement against Deepseek, 54%**improvement** againt retrieval+prompt based settings and **6.4% accuracy improvement** against retrieval+finetune. even when the pipeline using a lesser competent llm (Flan-t5). This improvement is even more significant when considering the more convoluted **huggingface dataset**, with Graph-O-Planner seeing nearly **68% improvement** over previous models as presented in Table 1.

5.2.2 Tool Hallucination Reduction

We also demonstrate the efficacy of integrating GNNs in reducing hallucination in LLMs. We use two hallucination metrics: micro hallucination and macro hallucination. These metrics are designed to quantify the extent of hallucination in the predicted sets of nodes compared to the ground truth sets.

Let N be the total number of samples, P_i be the predicted set of nodes for the i^{th} sample, and let V be the set of valid nodes.

Micro hallucination calculates the fraction of predicted nodes that are absent in the ground truth, averaged over all samples, represented as:

Micro Hallucination
$$= \frac{1}{N} \sum_{i=1}^{N} \frac{|P_i \setminus V|}{|P_i|}$$
 (15)

where $|P_i \setminus V|$ represents the number of nodes in P_i that are not in V, essentially the number of hallucinated nodes in the prediction.

Macro hallucination checks if any of the predicted nodes are absent from the ground truth and assigns 1 if at least one node is absent, 0 otherwise, and then averages over all samples:

Macro Hallucination
$$= \frac{1}{N} \sum_{i=1}^{N} I(P_i \setminus V \neq \emptyset)$$
(16)

where $I(P_i \setminus V \neq \emptyset)$ equals 1 if there are any nodes in P_i not in V (i.e., $P_i \setminus V$ is not empty), and 0 otherwise.

As shown in figure 4a and 4b, our proposed GNN-based approach achieves a substantial reduction in prediction hallucination, with a 13% decrease in incorrect edge predictions.

392

390

50

396

397

.....

400 401

402

405

406

407

408

409

410

411

412

413 414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

444 445

446 447

448

449



(a) Micro hallucination in Flan T5 XL vs Graph-O-Planner (b) Macro hallucination in Flan T5 XL vs Graph-O-Planner

Figure 4: Tool Prediction Hallucination in Flan T5 XL vs Graph-O-Planner

The results suggest that the GNN's ability to model complex structural relationships between tasks is instrumental in mitigating hallucination. By representing task sequences as graphs and leveraging the strengths of GNNs, we can better capture the nuances of task dependencies and generate more accurate and contextually relevant responses.

5.2.3 Model Latency

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479 480

481

482

483

484

A complementary benefit of our proposed approach is the reduction in input context size during both training and inference as presented in more detail in Appendix A.5. In the current literature, training Large Language Model (LLM) planners pass all the tool information, including name, description, input/output format, directly to the prompt resulting into longer context with a large number of tools, as seen in the Ultratool dataset. Even with context length of 8192, we observed a spill-over of input tokens. Since we inject tool knowledge as tool embeddings directly into the Graph Neural Network (GNN) layers, while the LLM focuses on the input query and the steps needed to execute the task. This makes it a more scalable and reliable solution for handling complex task sequences and a large number of tools.

This also significantly reduces inference time latency, as shown in Figure 5, making it more suitable for real-time applications due to the reduced input size, computational requirements, and focused input context. Our method achieves faster inference than encoder-decoder, reasoning, and promptbased LLMs. Compact inputs and rapid generation enable more efficient, reliable outputs with fewer hallucinations.





5.3 Ablation Study

1. Effect of Number of GNN Layers. We fur-486 ther conducted a study to evaluate the impact of 487 varying GNN layer depths (2, 5, 10) across tool 488 graph datasets of different structures and reported 489 the result in Table 3. We observe that for smaller 490 tool graphs such as huggingface, multimedia and 491 dailylife, lesser layers of GNN suffice. Increasing 492 the number of layers can cause a negative effect 493 due to over-smoothing effect. For UltraTool with 494 260 tools, 5 GNN layers perform the best. For the 495 Ultratool graph increasing depth from 2 to 5 lay-496 ers yielded statistically significant improvements 497 in edge prediction (edge-F1: $95.79\% \rightarrow 97.73\%$), 498 accompanied by a marginal yet consistent gain in 499 node prediction (node-F1: $98.83\% \rightarrow 99.19\%$). 500 Performance degradation at 10 layers (edge-F1: 501 97.26%) aligns with established phenomena of 502 over-smoothing in excessively deep architectures. 503 Conversely, datasets with dense connections but less number of nodes such as dailylife and hug-505

Method/Dataset	Huggingface		Ultratool		Multimedia		Dailylife	
	Edge-F1	Node-F1	Edge-F1	Node-F1	Edge-F1	Node-F1	Edge-F1	Node-F1
Message Passing	88.25	96.21	97.56	98.88	93.55	97.11	95.06	97.51
GCN	90.15	97.50	97.73	<u>98.99</u>	92.80	<u>99.05</u>	94.96	98.85
SGC	89.73	97.23	97.20	<u>98.99</u>	91.43	99.00	<u>95.43</u>	99.82
GAT	<u>89.73</u>	<u>97.36</u>	97.73	99.19	91.49	99.13	95.59	99.82

Table 2: Effect of different GNN networks used in Graph-O-Planner, evaluated on Node -f1 and Edge-f1. All result are in (%) with the best bold and runner-up underlined.

Layers/Dataset	Huggingface		Ultratool		Multimedia		Dailylife	
	Edge-F1	Node-F1	Edge-F1	Node-F1	Edge-F1	Node-F1	Edge-F1	Node-F1
2	89.73	97.12	95.79	98.83	93.52	98.87	95.59	99.84
5	88.74	97.36	97.73	99.19	93.18	99.13	92.78	99.71
10	88.8	98.96	97.26	98.72	93.55	99.04	94.54	99.82

Table 3: Performance comparison across datasets by GNN layer depth. For each dataset best performing epoch in 15 epochs is reported in above table. Best numbers were highlighted in bold.

gingface exhibited peak edge-F1 scores at 2 layers (95.59% and 89.73%, respectively), with deeper configurations inducing performance declines (dailylife edge-F1: 92.78% at 5 layers). Node prediction metrics remained stable across depths for these datasets (node-f1 $\Delta < 0.2\%$), suggesting limited utility of extended neighborhood aggregation in locally cohesive graphs.

506

507

508

509

510

511

512

514

515

516

517

518

519

520

521

522

523

524

525

527

528

529

531

533

535

536

539

2. Comparison of different GNN approaches. We experiment with different GNN approaches, namely Message Passing Neural Networks (MPNNs), Simplified Graph Convolution (SGC), standard Graph Convolutional Networks (GCN) and Graph Attention Networks (GAT). The experimental results presented (Table 2) demonstrate that GAT outperforms other GNNs for Ultratool, Multimedia and Dailylife and second best performance on Ultratool.

3. Other Design Choices - In addition to the type and size of GNN used, Graph-O-Planner reuses some of modules and techniques that have been proven as state-of-art in earlier works. As discussed earlier, we use frozen ModernBERT to encode the task graph information. It provides strong semantic representations while maintaining a lightweight and efficient architecture in comparison to previous work such as BERT and DeBERTa (discussed in Warner et al. (2024)). Our decision to freeze the encoder ensures that the entire learning capacity is focused on the GNN and LLM layers, keeping the pipeline modular and resource-efficient. Additionally, instead of normal softmax, we choose to use Gumbel Softmax. The recent work of (Zhang et al., 2024) validates the increased efficacy of using gumbel softmax for graphical architectures. As these choices are derivative usages of proven SOTA techniques, we do not include detailed insights into variations caused by these design choices. We provide more details of multi-layer and multi-level interactions in Graph-O-Planner in Appendix A.8. 540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

6 Conclusion

In this work, we propose Graph-O-Planner, a graphbased task selection method for generalized agent planning. Traditional prompt based methods of LLM based agent creation are hindered by concerns related to ever increasing tool context length, hallucinations and inductive biases. We propose using a GNN based network to effectively embed the information of the available tools, and use a knowledge-injection methodology in Graph-O-Planner to empower the LLM to map the sub-tasks to the appropriate tool sequence. Our method enables a more modular and flexible architecture by decoupling tool knowledge from the input prompt and injecting it into GNN layers, allowing for seamless integration of new tools and task sequences complementing the LLM for better performance as presented in Appendix A.3. We evaluate our model across 4 open-source datasets, comparing with multiple existing SOTA methodologies. As noted in results, we beat existing benchmarks by significant levels, enforcing the efficacy of the proposed model. The impact of tool information compression is also seen in inference latency, with a 1.25x improvement in inference speed. To the best of our knowledge, proposed work is the first of its kind, exploring a deeply integrated GNN-LLM framework for effective task planning.

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

625

626

574 Limitations

Despite encouraging performance, this work is only the beginning of exploring the GNN-LLM interac-576 tion in-depth. We also want to extend the pipeline to make the planning truly generalizable across any unseen tool-graph, task type, ambiguous data and erroreneous formats/descriptions. In real-life application scenarios, the available tools and user preferences will be constantly evolving and varied from 582 person-to-person. A truly intelligent agent should be able to effectively generalize across all such interactions without the need of any fine-tuning or 585 adaption. We aim to look into more details on these 586 in future works.

References

589

590

591

592

593

594

595

597

598

599

607

610

611

612

613

614

615

616

617

618

619

620

- Weilin Cai, Juyong Jiang, Fan Wang, Jing Tang, Sunghun Kim, and Jiayi Huang. 2024. A survey on mixture of experts. *arXiv preprint arXiv:2407.06204*.
- Chi-Min Chan, Weize Chen, Yusheng Su, Jianxuan Yu, Wei Xue, Shanghang Zhang, Jie Fu, and Zhiyuan Liu. 2023. Chateval: Towards better llm-based evaluators through multi-agent debate. *arXiv preprint arXiv:2308.07201*.
- Xiaojun Chen, Shengbin Jia, and Yang Xiang. 2020. A review: Knowledge reasoning over knowledge graph. *Expert systems with applications*, 141:112948.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. 2022. Scaling instruction-finetuned language models. arXiv preprint.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Andrew J Dudzik and Petar Veličković. 2022. Graph neural networks are dynamic programmers. Advances in neural information processing systems, 35:20635–20647.
- Lutfi Eren Erdogan, Nicholas Lee, Siddharth Jha, Sehoon Kim, Ryan Tabrizi, Suhong Moon, Coleman Hooper, Gopala Anumanchipalli, Kurt Keutzer, and Amir Gholami. 2024. Tinyagent: Function calling at the edge. *Preprint*, arXiv:2409.00608.

- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Lianzhe Huang, Dehong Ma, Sujian Li, Xiaodong Zhang, and Houfeng Wang. 2019. Text level graph neural network for text classification. *arXiv preprint arXiv:1910.02356*.
- Shijue Huang, Wanjun Zhong, Jianqiao Lu, Qi Zhu, Jiahui Gao, Weiwen Liu, Yutai Hou, Xingshan Zeng, Yasheng Wang, Lifeng Shang, Xin Jiang, Ruifeng Xu, and Qun Liu. 2024. Planning, creation, usage: Benchmarking llms for comprehensive tool utilization in real-world complex scenarios. *Preprint*, arXiv:2401.17167.
- Eric Jang, Shixiang Gu, and Ben Poole. 2017. Categorical reparameterization with gumbel-softmax. *Preprint*, arXiv:1611.01144.
- Elias Khalil, Hanjun Dai, Yuyu Zhang, Bistra Dilkina, and Le Song. 2017. Learning combinatorial optimization algorithms over graphs. *Advances in neural information processing systems*, 30.
- Jiachen Li, Xinyao Wang, Sijie Zhu, Chia-Wen Kuo, Lu Xu, Fan Chen, Jitesh Jain, Humphrey Shi, and Longyin Wen. 2025. Cumo: Scaling multimodal Ilm with co-upcycled mixture-of-experts. *Advances in Neural Information Processing Systems*, 37:131224– 131246.
- Shuwen Liu, Bernardo Grau, Ian Horrocks, and Egor Kostylev. 2021. Indigo: Gnn-based inductive knowledge graph completion using pair-wise encoding. *Advances in Neural Information Processing Systems*, 34:2034–2045.
- Linhao Luo, Yuan-Fang Li, Gholamreza Haffari, and Shirui Pan. 2023. Reasoning on graphs: Faithful and interpretable large language model reasoning. *arXiv preprint arXiv:2310.01061*.
- Nathalia Nascimento, Paulo Alencar, and Donald Cowan. 2023. Self-adaptive large language model (llm)-based multiagent systems. In 2023 IEEE International Conference on Autonomic Computing and Self-Organizing Systems Companion (ACSOS-C), pages 104–109. IEEE.
- Phu Pham, Loan TT Nguyen, Witold Pedrycz, and Bay Vo. 2023. Deep learning, graph-based text representation and classification: a survey, perspectives and challenges. *Artificial Intelligence Review*, 56(6):4893–4927.
- Timo Schick, Jane Dwivedi-Yu, Roberto Dessì, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. 2023. Toolformer: Language models can teach themselves to use tools. *Advances in Neural Information Processing Systems*, 36:68539–68551.

681

- 722 723 725 726 727 728 730 731 732
- 733 734 735

- Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. 2023a. Hugginggpt: Solving ai tasks with chatgpt and its friends in hugging face. Preprint, arXiv:2303.17580.
- Yongliang Shen, Kaitao Song, Xu Tan, Wenqi Zhang, Kan Ren, Siyu Yuan, Weiming Lu, Dongsheng Li, and Yueting Zhuang. 2023b. Taskbench: Benchmarking large language models for task automation. arXiv preprint arXiv:2311.18760.
- Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. 2023. Progprompt: Generating situated robot task plans using large language models. In 2023 IEEE International Conference on Robotics and Automation (ICRA), pages 11523-11530. IEEE.
- Yifan Song, Weimin Xiong, Dawei Zhu, Wenhao Wu, Han Qian, Mingbo Song, Hailiang Huang, Cheng Li, Ke Wang, Rong Yao, et al. 2023. Restgpt: Connecting large language models with real-world restful apis. arXiv preprint arXiv:2306.06624.
- Yashar Talebirad and Amirhossein Nadiri. 2023. Multiagent collaboration: Harnessing the power of intelligent llm agents. arXiv preprint arXiv:2306.03314.
- DJ Tena Cucala, B Cuenca Grau, Egor V Kostylev, and Boris Motik. 2022. Explainable gnn-based models over knowledge graphs.
- Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. 2023a. Voyager: An openended embodied agent with large language models. Preprint, arXiv:2305.16291.
- Lei Wang, Wanyu Xu, Yihuai Lan, Zhigiang Hu, Yunshi Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. 2023b. Plan-and-solve prompting: Improving zeroshot chain-of-thought reasoning by large language models. arXiv preprint arXiv:2305.04091.
- Yu Wang, Zhiwei Liu, Ziwei Fan, Lichao Sun, and Philip S Yu. 2021. Dskreg: Differentiable sampling on knowledge graph for recommendation with relational gnn. In Proceedings of the 30th ACM International Conference on Information & Knowledge Management, pages 3513-3517.
- Benjamin Warner, Antoine Chaffin, Benjamin Clavié, Orion Weller, Oskar Hallström, Said Taghadouini, Alexis Gallagher, Raja Biswas, Faisal Ladhak, Tom Aarsen, Nathan Cooper, Griffin Adams, Jeremy Howard, and Iacopo Poli. 2024. Smarter, better, faster, longer: A modern bidirectional encoder for fast, memory efficient, and long context finetuning and inference. Preprint, arXiv:2412.13663.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837.

Lingfei Wu, Yu Chen, Kai Shen, Xiaojie Guo, Hanning Gao, Shucheng Li, Jian Pei, and Bo Long. 2021. Graph neural networks for natural language processing: A survey. CoRR, abs/2106.06090.

736

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

775

778

781

782

783

784

785

787

- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. 2023. Autogen: Enabling next-gen llm applications via multiagent conversation framework. arXiv preprint arXiv:2308.08155.
- Xixi Wu, Yifei Shen, Caihua Shan, Kaitao Song, Siwei Wang, Bohang Zhang, Jiarui Feng, Hong Cheng, Wei Chen, Yun Xiong, et al. 2024. Can graph learning improve task planning? arXiv preprint arXiv:2405.19119.
- Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2024. Efficient streaming language models with attention sinks. Preprint, arXiv:2309.17453.
- Keyulu Xu, Jingling Li, Mozhi Zhang, Simon S Du, Ken-ichi Kawarabayashi, and Stefanie Jegelka. 2019. What can neural networks reason about? arXiv preprint arXiv:1905.13211.
- Peng Xu, Xinchi Chen, Xiaofei Ma, Zhiheng Huang, and Bing Xiang. 2021. Contrastive document representation learning with graph attention networks. In Findings of the Association for Computational Linguistics: EMNLP 2021, pages 3874-3884, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Daviheng Liu, Fei Huang, Haoran Wei, et al. 2024. Qwen2. 5 technical report. arXiv preprint arXiv:2412.15115.
- Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2024. Tree of thoughts: Deliberate problem solving with large language models. Advances in Neural Information Processing Systems, 36.
- Yunzhi Yao, Shaohan Huang, Li Dong, Furu Wei, Huajun Chen, and Ningyu Zhang. 2022. Kformer: Knowledge injection in transformer feed-forward layers. In CCF International Conference on Natural Language Processing and Chinese Computing, pages 131–143. Springer.
- Zi Ye, Yogan Jaya Kumar, Goh Ong Sing, Fengyan Song, and Junsong Wang. 2022. A comprehensive survey of graph neural networks for knowledge graphs. IEEE Access, 10:75729-75741.
- Junwei Yu, Yepeng Ding, and Hiroyuki Sato. 2025. Dyntaskmas: A dynamic task graph-driven framework for asynchronous and parallel llm-based multiagent systems. Preprint, arXiv:2503.07675.

- Qinggang Zhang, Junnan Dong, Hao Chen, Xiao Huang, Daochen Zha, and Zailiang Yu. 2023. Knowgpt: Black-box knowledge injection for large language models. *arXiv preprint arXiv:2312.06185*.
- Yu Zhang, Kehai Chen, Xuefeng Bai, Quanjiang Guo, Min Zhang, et al. 2024. Question-guided knowledge graph re-scoring and injection for knowledge graph question answering. *arXiv preprint arXiv:2410.01401*.
- Yanqi Zhou, Tao Lei, Hanxiao Liu, Nan Du, Yanping Huang, Vincent Zhao, Andrew M Dai, Quoc V Le, James Laudon, et al. 2022. Mixture-of-experts with expert choice routing. *Advances in Neural Information Processing Systems*, 35:7103–7114.
- Jason Zhu, Yanling Cui, Yuming Liu, Hao Sun, Xue Li, Markus Pelger, Tianqi Yang, Liangjie Zhang, Ruofei Zhang, and Huasha Zhao. 2021. Textgnn: Improving text encoder via graph neural network in sponsored search. In *Proceedings of the Web Conference 2021*, pages 2848–2857.

A Appendix

A.1 Notations

790

796

803

810

811

812

813

814

815

816

817

818

819

821

822

823

825

826

827

830

832

834

835

838

The symbolic notations used in the paper are summarized in Table 4.

A.2 Implementation Details

Detailed information on the experimental setup section 5.

Hardware. All the models are trained using Py-Torch 2.3.1 framework in Python 3.11 conda environment. Ubuntu server equipped with four 48GB Nvidia-RTX A6000 with driver version 560.35.03 and CUDA 12.6 are utilized to perform the study.

Hybrid GNN+LLM Baselines: We used the repo by Wu et al. (2024) to reproduce the results.

LLM-Only: We used the Transformer library's Trainer with LoRA similar to Graph-O-Planner to train the piepeline.

RAG Baselines We utilize open-source ChromaDB framework which uses a hierarchical configuration architecture, employing the all-MiniLM-L6-v2 transformer model as its default text embedding engine, which produces 384-dimensional semantic vectors optimized for cosine similarity computations through PyTorch-based inference. Input text undergoes preprocessing via whitespace normalization and lexical truncation at 256 tokens prior to vectorization, with document batches processed at 32-sample granularity to balance memory efficiency and throughput. The system implements an in-memory HNSW (Hierarchical Navigable Small World) index with empirically tuned parameters (ef_construction=200, M=16) to optimize approximate nearest neighbor search latency. Graph-O-Planner: GNN. The dimensions of GNN module is converted from 1024 embedding size to 200, which can be modified as per user's choice and a variable number of layer of GNN modules between 5 and 7 both inclusive with a dropout of 0.18 applied within each consecutive layer. Training. All models are trained using Adam optimizer. Learning rate of LLM module and GNN module is kept 1e-4 and 3e-4 respectively, with batch size of 32. We choose FLAN-T5-XL(3B) model as LLM for Low-Rank Adaption (LoRA) training with rank 8 and alpha 16. The maximum token length across tokenizer is kept variable as per the requirement of dataset.

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

A.3 Additional aid to LLM with Graph-O-Planner

In the figure 6, 7, 8, 9, 10, 11 shown Dev, test and loss for only LLM approach in comparison with Graph-O-Planner approach. From figure 6 and 7, it is observed that Graph-O-Planner approach with the help of tool embeddings from GNN through multilevel interaction learn meaningful insights about the tool graph and surpass 80% edgef1 in merely 5 epochs and settles at >90% after 10 epochs. While in T5 only training approach we find that the overall edge-f1 cannot surpass 60% even after 30 epochs as shown in figure 8 and 11.

We also observed a much reliable training with Graph-O-Planner approach. As seen from Figure 10, the loss curve much cohesively justifies the overall loss when compared with the improvement seen from test eval curve. While for T5 only approch in Figure 11 we can observe that the loss drops drastically till 40th epoch, but soon reaches a stagnant curve, however as can be observed from 9, the test edge-f1 has a lot of scope for improvement. From these result we can come to conclusion that Knowledge fusion between LLM and GNN can lead to benefits listed below:

• Unified Task Perspective. In our approach, the LLM can be directly leveraged to produce outputs for multiple tasks. For varying tasks, it can either operate in a masked mode using precomputed embeddings—eliminating the need for re-computing task graph. This underscores our core contribution: the LLM+GNN functions as a flexible, "plug-and-play" mod-

Notations	Definition
G, V, E	Tool graph with set of nodes V and edges E
T_i, A_i	Features embeddings i-th tool in graph G, A represents adjacency matrix
T	Tool information in graph G
Q_i	i-th query of dataset
$S_i = s_1, \dots s_n$	Subtasks of query Q_i
Emb(.)	Embedding function representing tool info in graph space for GNN training
$h^{(l)}$	GNN node representation at step l
M(.)	Edge aware attention function
$1_{\phi(e_{ij})}$	Encoded edge features obtained after applying Gumbel softmax transform
$1_{\tau(v_i)}$	Encoded node features obtained after applying Gumbel softmax transform
(e_{ij})	Edge representation obtained after concatenation of encoded edge and node features
q_j^k, k_j^k, v_{ij}^k	Query, Key and Value projections of k-th GNN node
α_{ij}^k	Normalized node attention based on out degree of k-th GNN node
N(i)	Out degree of i-th node
m_i^k	Head Specific attention of k-th GNN node after message passing to neighbors
h'_i	Edge representation of GNN node after message passing
I_c	Intermediate LLM+GNN interaction layer
$ heta_{key}, heta_{query}, heta_{value}$	Key, Query and Value obtained from LLM decoder
$\phi_{key}, \phi_{query}, \phi_{value}$	Key, Query and Value obtained from GNN decoder

Table 4: Notation table in Graph-O-Planner

ule, significantly improving efficiency and performance over conventional large language model (LLM)-only approaches by storing precomputed task graph embeddings and lower context length requirement.

890

892

893

894

895

896

900

901

902

903

904

905

906

907

908

909

910 911

912

913

914

915

Integrated Fusion Strategy. Our fusion strategy facilitates concurrent information propagation from tool graphs (hidden embeddings) and (task-specific output vectors) to the query input. This enables structured knowledge injection and task-specific adaptation, making our LLM + Graph Network (GNN) paradigm superior to LLM-only models, which often struggle with structural reasoning and compositional generalization. By leveraging graph representations, our approach effectively captures relational dependencies, improving both adaptability and interpretability across tasks.

A.4 Necessity of GNN integration

Large Language Model (LLM)-only planners require all tool information—including tool descriptions and capabilities to be included within the input prompt at both training and inference time. This quickly becomes infeasible due to the limited context window of current LLMs (**Context Overflow Problem**). When the tool-related content exceeds the model's context length, critical parts of the input are truncated, leading to incomplete information available for planning. 916

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

Even though models with longer context windows are emerging, they come with practical tradeoffs:

- Latency and Memory Bottlenecks: The latency increases nearly quadratically as attention has $(O(n^2))$ time complexity. So, when the prompt size increases linearly, the latency increases quadratically.
- Windowed Attention Limitations: Many long-context LLMs rely on windowed attention mechanisms, which suffer from "window sink" issues that prevent the model from capturing long-range dependencies across windows, as demonstrated in earlier works (Xu et al., 2021; Xiao et al., 2024).

To address these limitations, Graph Neural Networks (GNNs) encodes and pools tool information into a fixed-size vector representation, independent of the number of tools. GNNs are well-suited to model complex and structured data, capturing longrange and interdependent relationships between tools efficiently.

943

947

951

952

954

955

957

961

962

963

964

966

968

969

970

971

974

975

976

979

A.5 How Graph-O-Planner overcomes context overflow and latency problem

The integration of Large Language Models (LLMs) with Graph Neural Networks (GNNs) presents a compelling advancement over LLM-only approaches for task planning, particularly in scenarios involving an extensive set of tools with complex specifications. Traditional LLM-based methods rely on tokenization to encode tool-related information, which inherently limits scalability due to increasing sequence lengths and associated computational costs. In our experiments, we observed that models such as Qwen struggle when provided with a large number of tools and their descriptions in Ultratool dataset. The excessive tokenization required to process tool details not only constrains the model's ability to handle more elaborate queries but also results in increased latency and memory overhead, making real-time task planning inefficient.

Context Overflow mitigation. Our proposed Graph-O-Planner framework mitigates these limitations by encoding tool information as embeddings within a graph structure, rather than representing them as lengthy text sequences. By leveraging GNNs to store and propagate tool-specific embeddings, we significantly reduce tokenization overhead, enabling the LLM to allocate more of its token budget toward processing complex queries rather than repetitive tool descriptions. This structured approach enhances efficiency by shifting the burden of tool representation from token-based encoding to a graph-based framework, leading to more scalable and interpretable reasoning over available tools. Furthermore, the graph structure inherently captures relational dependencies between tools, facilitating a more structured understanding of tool applicability and interoperability.

Reduced Model Latency. Beyond tokenization efficiency, the incorporation of GNNs also enhances inference speed by caching for repetitive Task info embeddings on first fly. In LLM-only approaches, each query requires reprocessing tool descriptions, leading to redundant computation. In contrast, our GNN-enhanced model precomputes and stores tool embeddings, allowing for direct retrieval and propagation of relevant tool information without unnecessary recomputation. This not only accelerates inference but also ensures that the model retains a more contextually enriched and persistent representation of tools across different task planning queries. By leveraging message-passing mechanisms within the GNN, our approach ensures efficient information flow, reducing the reliance on autoregressive decoding for tool-related reasoning.

By encoding tool knowledge in a structured graph representation, we achieve a dual advantage: reducing tokenization demands while improving inference efficiency. This allows for handling more complex and multi-step task planning scenarios, where an LLM alone would struggle due to token constraints and redundant processing. Our findings demonstrate that integrating structured graph-based reasoning with LLMs enables more effective tool selection, faster response times, and improved scalability, making it a superior approach for real-world Agentic planning applications.



Figure 6: Dev edge-f1 of Graph-O-Planner



Figure 7: Test edge-f1 of Graph-O-Planner

A.6 Dataset

In this section we will deep dive into dataset mentioned in section 5.1.

Ultratool. It consist of 260 tools with 3527 1010 task and steps samples. On average each sample's 1011

1007

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

1005



Figure 8: Dev edge-f1 of Flan T5 XL



Figure 9: Test edge-f1 of Flan T5 XL



Figure 10: Loss graph of Graph-O-Planner after training for 15 epochs

plan include 2.42 tool callings. All samples within ultratool have at least one tool calling. In particular 64.24% of samples consist of two tool calling and rest consist of multiple tool calling. Each sample contains at least two tool calls.

Huggingface. It consist of 40 tools with 7546 training samples. The tools comprise of hugging face hosted models fine-tuned to perform various



Figure 11: Loss of Flan T5 XL while training for 30 Epochs

downstream tasks. The overall dataset requires 20177 tool callings with an average of 3.28 arguments per tool call. The dataset consists of 40.64% of samples with single tool calling.

Multimedia. This dataset also consist of 40 unique tools with 5584 training samples. The tools comprises of generic multimedia tools like 'Video-to-Audio', 'Audio-Splicer' etc. It consist of 15860 distinct tool calls with 3.49 arguments per tool call. Out of 5584 samples 36.48% of samples requires only single tool calling.

Dailylife. The dataset contains 40 distinct tools with 4320 samples out of which 1258 samples contains single tool calls. In the dataset it requires on average of 3.09 tool calls per sample with average of 4.95 arguments per tool call. The tool consist of general tool present in most virtual assistants like 'book_hotel', 'book_flight' etc.

Next we show sample input data fed from these dataset.

Huggingface

text { 'id ': '57993067 ',	1041
'seed ': 513420,	1042
'n_tools ': 1,	1043
'sampled_nodes ':	1044
[{ 'task ': 'Object Detection ',	1045
'input-type ': ['image'],	1046
'output-type ': ['text']}],	1047
'sampled_links ': [],	1048
'user_request ': "I need	1049
to identifyand label objects	1050
in the provided image	1051
'example.jpg'.",	1052
'task_steps ': [1053
'Step 1: Use Object	1054
Detection to identify	1055
5	

```
objects in the image
1056
                and label them.'
1057
1058
            ],
             'task_nodes ': [{
1059
                'task ': 'Object Detection ',
                'arguments ': [ 'example.jpg ']
1061
1062
                }
             ],
1063
             'task_links ': [],
1064
             'type ': 'single'}
1065
             Multimedia
1066
           { 'id ': '16097613',
1067
             'seed ': 154967,
1068
             'n_tools ': 3,
1069
             'sampled_nodes ':
               [{ 'input-type ': ['audio',
               'text '],
1072
               'output-type': ['audio'],
1073
               'task ': 'Audio Effects '},
1074
              { 'input-type ': ['audio'],
1075
               'output-type ': ['audio'],
               'task ': 'Audio Noise
1077
               Reduction '},
1078
1079
              { 'input-type ': [ 'video '],
               'output-type ': ['audio'],
1080
               'task ': 'Video-to-Audio'}],
1081
             'sampled_links ': [
                { 'source ': 'Audio Noise
1083
                Reduction ',
1084
                 'target ': 'Audio Effects '},
1085
                { 'source ':
                             'Video-to-Audio',
1086
                 'target ': 'Audio Noise
1087
                 Reduction '}],
1088
1089
             'user_request ': 'I have a video
             file example.mp4, and I want to
1090
            extract its audio track, reduce
1091
1092
            background
            noise, and then add a reverb
1093
            effect.
1094
            Please provide the
1095
            processed audio file.',
1096
             'task_steps ': [
                'Extract audio from the given
1098
                 video file',
1099
                'Reduce noise from the
1100
                extracted audio',
1101
1102
                 'Apply audio effects to the
                 noise-reduced
1103
                  audio according to user
1104
                  instructions '],
1105
             'task_nodes ': [
1106
```

{ 'task ': 'Audio Effects ',	1107
'arguments ':	1108
[' <node-1>', 'reverb']},</node-1>	1109
{ 'task ': 'Audio Noise	1110
Reduction ',	1111
'arguments ': [' <node-2>']},</node-2>	1112
{ 'task ': 'Video-to-Audio',	1113
'arguments ':['example.mp4']}],	1114
'task_links ': [1115
{ 'source ': 'Audio	1116
Noise Reduction ',	1117
'target ': 'Audio Effects'},	1118
{ 'source ': 'Video-to-Audio',	1119
'target': 'Audio Noise	1120
Reduction '}],	1121
'type ': 'chain'}	1122
Dailvlife	1123
{'id ': '13590101',	1124
seed ': 283/17,	1125
'n_tools': 1,	1126
sampled_nodes ': [1127
{ task : play_movie_by_title ,	1128
arguments : [{ name : title ,	1129
type : string ,	1130
desc : Ine title of the	1131
movie to play }]}],	1132
sampled_links : [],	1133
user_request . I want to	1134
'Example Movie'"	1135
Lample Movie ,	1107
nlay movie by title API	1120
with title: 'Example Movie'"	1130
'task nodes '· [1140
{ 'arguments '· [1141
{'name': 'title '	1142
'value ': 'Example Movie'}].	1143
'task ': 'play movie by title '}].	1144
'task links ': [],	1145
'type ': 'single'}	1146
Ultratool	1147
{ 'id ': '3186',	1148
'user_request': 'I need to	1149
cancel the single alarm set	1150
for 8:00 AM today, and change	1151
the daily alarm from 7:00 AM	1152
to 6:30 AM every day.\n',	1153
'task_steps ': [1154
'Step 1 Call clock_alarm_cancel	1155
to cancel the alarm set for	1156

1157	8:00 AM today',
1158	'Step 2 Call clock_alarm_change
1159	to change the daily
1160	alarm from 7:00 AM to 6:30 AM
1161	every day'],
1162	'task_nodes ': [
1163	{ 'task ':
1164	<pre>'clock_alarm_cancel'} ,</pre>
1165	{ 'task ':
1166	<pre>'clock_alarm_change '}],</pre>
1167	'task_links ': [
1168	{ 'source ':
1169	'clock_alarm_cancel ',
1170	'target ':
1171	<pre>'clock_alarm_change '}],</pre>
1172	'n_tools ': 2,
1173	'type ': 'chain'}

A.7 Baselines

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

1190

1191

1192

1193

1194

1195

1196

1197

1198

1199

1201

1203

In this appendix section, we present the details of baselines shown in Table 1.

- **Graph Token.** A method that introduces a global virtual token to GNNs allowing improved global information aggregation and better graph-level representations.
- **GraphSAGE.** A GNN that learns node embeddings by sampling and aggregating information from a nodes' neighborhood, enabling scalable learning on large graphs.
- GCN(Graph Convolutional Network). A fundamental GNN model that extends convolutional operations to graph structure by propagating and aggregating node features using adjacency-based weight metrics
- GAT(Graph Attention Network). A GNN model that incorporates attention mechanism to assign different importance weights to neighboring nodes, improving feature aggregation adaptively.
- **GIN(Graph Isomophism Network).** A powerful GNN variant designed to be as expressive as the Weisfeiler-Lehman graph isomorphism test, using MLP-based neighborhood aggregation.
- Deepseek R1. DeepSeek-R1 is a reasoning model that achieves performance comparable to OpenAI-o1 across math, code, and reasoning tasks, and is open-sourced along with

its distilled dense models to support the re-1204 search community. DeepSeek-R1 is devel-1205 oped through a pipeline that incorporates rein-1206 forcement learning and supervised fine-tuning, 1207 and its reasoning patterns can be distilled 1208 into smaller models, resulting in better per-1209 formance on benchmarks. like Multi head 1210 Latent attention. 1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1248

1249

1250

1251

1252

1253

• Qwen 2.5 Coder. Qwen2.5-Coder is a large language model series with six mainstream model sizes, offering improved code generation, reasoning, and fixing capabilities. It has become the state-of-the-art open-source codeLLM, matching the coding abilities of GPT-40 with enhanced coding capabilities and long-context support up to 128K tokens.

A.8 More Detailed study

1. Effect of multilayer interaction between LLM 1221 and GNN. We investigate the contribution of multi-1222 layer interaction, experiments using single-layer 1223 interaction of SGC, GCN, GAT, and GraphSAGE 1224 without multi-layer feature injection showed that 1225 this component significantly boosts performance 1226 (up to 30%) by aggregating richer contextual infor-1227 mation across layers as observed in Table 5. Our ex-1228 periments demonstrate the critical role of hierarchi-1229 cal feature aggregation through direct comparison 1230 with shallow graph convolution baselines. Single-1231 layer variants (SGC, GCN, GAT, GraphSAGE) 1232 exhibit substantially inferior performance across 1233 all datasets-Huggingface edge-F1 reaches just 1234 43.09% (GraphSAGE) versus our 89.73%, while 1235 Multimedia node-F1 plateaus at 75.51% (Graph-1236 SAGE) versus our 99.13%. The performance dif-1237 ferential is most pronounced in edge prediction 1238 tasks, where our method achieves relative improve-1239 ments of 108.3% (Huggingface edge-F1: 89.73 1240 vs. 43.09) and 44.6% (Dailylife edge-F1: 95.59 1241 vs. 66.57) over the strongest shallow baselines. 1242 Even node classification, traditionally less depth-1243 sensitive, shows absolute gains of 29.85% (Hug-1244 gingface) and 23.62% (Multimedia) compared to 1245 single-layer counterparts, empirically validating 1246 the necessity of cross-layer information fusion. 1247

2. Architectural implication of Multilevel fusion. The stark performance margins (Δ edge-F1 > 46% across all datasets) reveal fundamental limitations of shallow interaction paradigms. While shallow methods like SGC achieve computational efficiency through layer truncation, they

Method/Dataset	Huggingface		Multimed	lia	Dailylife	
	Node-F1	Edge-F1	Node-F1	Edge-F1	Node-F1	Edge-F1
SGC	67.43	42.08	74.07	49.90	87.13	66.49
GCN	66.54	40.74	73.34	50.76	86.39	65.49
GAT	66.77	40.74	73.36	50.20	86.39	65.49
GraphSAGE	68.12	43.09	75.51	52.94	87.51	66.57
Graph-O-Planner (ours)	97.36	89.73	99.13	91.49	99.82	95.59

Table 5: Comparison of shallow interaction between LLM and various GNN settings vs. ours

 $DE(V \cap V)$

DE(V V)

forfeit the ability to capture hierarchical dependen-1254 cies-evidenced by Huggingface's edge prediction 1255 collapse to 42.08% (SGC) versus our 89.73%. Our 1256 multi-layer architecture addresses this through de-1257 liberate feature injection across depths, enabling 1258 progressive refinement of both local and global 1259 1260 graph patterns. This is particularly crucial for complex edge prediction tasks, where shallow mod-1261 els lack the representational capacity to resolve 1262 indirect relationships (e.g., Multimedia edge-F1: 1263 52.94% vs. 91.49%). The consistent outperformance (minimum Δ node-F1: 12.31% in Dailylife) 1265 across diverse graph types further substantiates 1266 multi-layer interaction as a generalizable design 1267 principle rather than a dataset-specific optimiza-1268 tion. 1269

1270

3. Proof of effectiveness of GNN.

In this section, we will theoretically prove that 1271 using Graph-O-Planner can significantly improve 1272 the LLM generation performance. Assume X_l as 1273 input tokens to the LLM and G as input tool graph 1274 features to the GNN, Y represents target output 1275 tool labels. We introduce a dependency function DF(.) that quantifies the dependency between in-1278 put labels X and Y, which reflects the performance of LLM. By introducing tool-graph knowledge into 1279 GNN, we can impactfully improve model perfor-1280 mance in predicting labels Y as $DF(X_l, G, Y) \ge$ $DF(X_l, Y)$. The following outlines the derivation: 1282

$$\begin{aligned} &= \sum_{X_l,G,Y} p(X_l,G,Y) = DF(X_l,Y) \\ &= \sum_{X_l,G,Y} p(X_l,G,Y) \log\left(\frac{p(X_l,G,Y)}{p(X,G)p(Y)}\right) \\ &- \sum_{X_l,Y} p(X_l,Y) \log\left(\frac{p(X_l,Y)}{p(X)p(Y)}\right) \\ &= \sum_{X_l,G,Y} p(X_l,G,Y) \log\left(\frac{p(X_l,G,Y)}{p(X,G)p(Y)}\right) \\ &- \sum_{X_l,G,Y} p(X_l,G,Y) \log\left(\frac{p(X_l,G,Y)}{p(X,G)p(Y)} \cdot \frac{p(X_l)p(Y)}{p(X_l,Y)}\right) \\ &= \sum_{X_l,G,Y} p(X_l,G,Y) \log\left(\frac{p(X_l,G,Y)}{p(G|X)p(Y)p(X_l,Y)}\right) \\ &= \sum_{X_l,G,Y} p(Y,G|X) \log\left(\frac{p(Y,G|X)}{p(G|X)p(Y)p(Y)p(Y|X)}\right) \end{aligned}$$

1284

1285

1286

1287

1288

1289

1290

1291

1292

1293

1294

1295

1296

1297

1298

1299

1300

1301

1302

1303

1304

1305

1306

4. Training time Computation analysis In this section we provide the computation comparison of LLM only v/s Graph-O-Planner approach. During the experiments, the number of trainable parameters remains constant across all the dataset. We observed that for Graph-O-Planner approach the time required for each epoch ranges between 23-32 minutes, while for LLM only approach the time taken to complete one epoch ranges between 150-180 minutes. From these results we infer that our approach is much faster and promising than other SOTA methods.

5. Rationale for usage of ModernBERT as text-embedder for encoding tool info. ModernBERT-Large has been used to generate initial embeddings for each node/tool description, which are then passed as input features to the GNN. The recent work Warner et al. (2024) shows ModernBERT-Large performing better than its predecessors on NLU tasks. It is also finetuned on using triplet networks (i.e. NLI tasks), thus a suitable choice to work with similarity, clustering and retrieval tasks. Morever, it improves upon BERT

1309

1310 1311

1312

1313

1314

1315

1316

1317

1318

1319

1320

1321

1322

1323

1324 1325

1326

1327

1328

1329

1330

1331

1332

1333

1334

1335

1336

1337

1338 1339

1340

1345

1355

and RoBERTa by combining masked and permuted language modelling, capturing global dependencies better other approaches like MiniLM or MPNet.

Algorithms A.9

In this section we provide detailed description of all the major algorithms explained in section 5.2.1.

Algorithm 1: Node F1. The algorithm takes two list as input. Lines 1-2 contains ground truth tools L and predicted tools R which is given to the function in Line 3-19 to calculate node f1. In more detail Lines 4-5 computes the length of lists L and R and stores them in gt_len and $pred_len$ respectively. Lines 8-10 stores unique tool names from ground truth in set *gt_tools* and respectively for predicted tools in *pred_tools* in Lines 11-13.

Finally, the node f1 is calculated in Line 14-17 by taking precision and recall and storing in variables p and r and then computing $node_f1$.

Algorithm 2: Edge F1. The algorithm takes two list. Lines 1-2 contains ground truth tools Land predicted tools R which is given to the function in Line 3-19 to calculate edge f1. Lines 4-5 computes the length of lists L and R and stores them in gt_len and pred_len respectively. Lines 7-13 takes every tool link present in predicted links and checks if the tool is present in ground truth links. If the tools is present, the counter of common links *c_links* is increased by 1. Finally in Line 14, 15 precision and recall for links are computed and then $edge_f1$ is calculated in Line 17.

Finally, the node accuracy is calculated in Line 14-15 as length of intersection set between predicted tool names and ground truth tool names set over length of *qt_tools*.

Algorithm 1:Node-F1
1: $L \leftarrow [l_1, l_2, \dots, l_n] \triangleright$ List of ground truth tool
pairs. $r_i: (tool_{i1}, tool_{i2})$
2: $R \leftarrow [r_1, r_2,, r_n] \triangleright$ List of predicted tool
pairs. $r_i: (tool_{j1}, tool_{j2})$
3: procedure NODE $F1(L, R)$
4: $gt_len \leftarrow length \ ofL$
5: $pred_len \leftarrow length \ of R$
6: $gt_tools \leftarrow \{\}$
7: $pred_tools \leftarrow \{\}$
8: for $i = 1$ to gt_len do
9: $gt_tools = gt_nodes \cup L[i][0] \cup$
$L\left[i ight]\left[1 ight]$
10: end for
11: for $i = 1$ to $pred_len$ do

12:	$pred_tools = gt_nodes \cup R[i][0] \cup$	1356
	R[i][1]	1357

14:
$$c_tools \leftarrow pred_tools \cap gt_tools$$
 1359

1358

1363

1364

1365

1388

15:
$$p \leftarrow \frac{length(c_tools)}{length(pred_tools)}$$
 1360

$$16: \quad r \leftarrow \frac{length(c_tools)}{length(gt_tools)}$$

$$1361$$

17:
$$node_f 1 = \frac{2*p*r}{p+r+\alpha}$$
 1362

Algorithm 2:Edge-F1

1:	$L \leftarrow [l_1, l_2, \dots, l_n] \triangleright$ List of ground truth tool	1366
	pairs. $r_i: (tool_{i1}, tool_{i2})$	1367
2:	$R \leftarrow [r_1, r_2,, r_n] \triangleright$ List of predicted tool	1368
	pairs. $r_i: (tool_{j1}, tool_{j2})$	1369
3:	procedure EDGE $F1(L, R)$	1370
4:	$gt_len \leftarrow length \ of L$	1371
5:	$pred_len \leftarrow length \ of R$	1372
6:	$c_links \leftarrow 0$	1373
7:	for $i = 1$ to $pred_len$ do	1374
8:	for $j = 1$ to gt_len do	1375
9:	if $R[i] == L[j]$ then	1376
10:	$c_links \leftarrow c_links + 1$	1377
11:	end if	1378
12:	end for	1379
13:	end for	1380
14:	$c_tools \leftarrow pred_tools \cap gt_tools$	1381
15:	$p \leftarrow \frac{length(c_tools)}{length(pred_len)}$	1382
16:	$r \leftarrow \frac{length(c_tools)}{length(gt_len)}$	1383
17:	$edge_f1 \leftarrow \frac{c_links}{gt_len}$	1384
18:	return edge_f1	1385
19:	end procedure	1386
A	Algorithm 3: Pseudocode for Graph-O-Planner	1387

1:	Input: Task query Q , tool metadata $T =$	1389
	$\{(n_i, d_i, i_i, o_i)\}_{i=1}^V$	1390
2:	Output: Tool execution DAG $DAG = (v_1 \rightarrow v_2)$	1391
	$v_2 \rightarrow \ldots \rightarrow v_P)$	1392
3:	procedure GRAPHOPLANNER (Q, T)	1393
4:	// Step 1: Tool Graph Construction	1394
5:	for each tool $t_i \in T$ do	1395
6:	Compute embedding: $x_i =$	1396
	$\operatorname{Emb}(n_i, d_i, i_i, o_i)$ (Eq. 1)	1397
7:	end for	1398

1399 8: Form tool graph
$$G = (V, E, A, T, X)$$
 with
1400 nodes V , edges E , and embeddings $X = \{x_i\}$
141 9: $//$ Step 2: Edge Encoding
142 10: for each edge $(v_i, v_j) \in E$ do
143 11: Compute edge representation:
144 $\epsilon_{ij} = f_{edge}(\phi_{e,ij} \oplus \tau_{k_i} \oplus \tau_{k_j})$ (Eq. 2)
145 12: end for
146 13: $//$ Step 3: Attention-Enhanced Graph
147 Convolution
148 for each GNN layer l do
149 15: for each node v_i do
140 16: $h_i^{(l)} = x_i$ if $l = 0$
141 17: Aggregate neighboring edge fea-
141 17: Aggregate neighboring edge fea-
141 18: Compute attention heads:
141 $Q_i^k = W_Q^k \tilde{h}_i, K_i^k = W_K^k \tilde{h}_i, V_i^k = W_V^k \tilde{h}_i$ (Eqs. 4–6)
141 18: Compute attention:
141 18: Compute attention:
141 19: Normalize attention:
141 19: Normalize attention:
141 19: $\alpha_{ij}^k = \frac{d_j}{Z_i} \cdot \exp\left(\frac{(Q_j^k, K_i^k)}{\sqrt{d}}\right)$ (Eq. 7)
141 20: Aggregate messages:
141 $m_i^k = \sum_{j \in N(i)} \sqrt{d}$) (Eq. 8)
142 21: end for
142 22: Fuse heads and update node embed-
142 dings:
142 $H^{(l+1)} = MLP(Agg(M(H^{(l)}, E, \epsilon)))$ (Eq. 9)
142 23: end for
142 24: J Step 4: Inject GNN Embeddings into
143 12: LLM
144 25: For selected LLM layers:
145 $I_c = \{\theta_{key} \oplus \psi_{key}, \theta_{query} \oplus \psi_{query}, \theta_{value} \oplus \psi_{value}\}$
146 (Eq. 10)
147 25: Where ψ , are GNN projections passed
147 17: where ψ . are GNN projections passed
148 28: $//$ Step 5: Decode Tool DAG
149 29: Use LLM (augmented with injected GNN
140 $DAG = (v_1 \rightarrow v_2 \rightarrow ... \rightarrow v_P)$
143 30: end procedure