# LEGO-GraphRAG: Modularizing Graph-based Retrieval-Augmented Generation for Design Space Exploration

Yukun Cao, Zengyi Gao, Zhiyang Li, Xike Xie, S. Kevin Zhou, Jianliang Xu
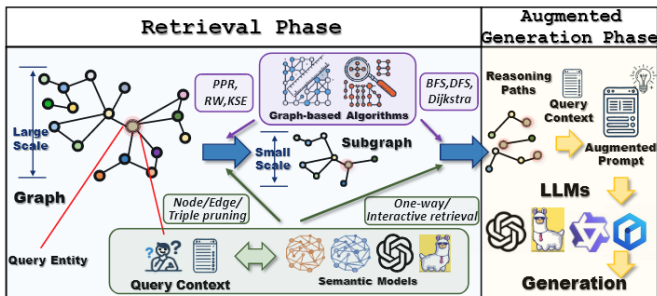
November 16, 2025

# Outline

# GraphRAG

- Retrieval-Augmented Generation (RAG) enhances LLMs with external knowledge.
- GraphRAG improves RAG by leveraging graph structures: entities, relationships, communities.

# Challenges

- Lack of Unified Standards: different graph algorithms and semantic models.
- Difficulty in Modular Optimization: retrieval process is often written as a monolithic block.
- Lack of Testing Platform: There is no public framework that can easily generate instances and perform large-scale comparisons.

# Contributions

- LEGO-GraphRAG framework:
  - dividing the retrieval phase into two flexible modules: subgraph-extraction and path-retrieval.
  - classifies the techniques into structure-based and semantic-augmented methods for each module.
- supports implementing all existing GraphRAG instances.
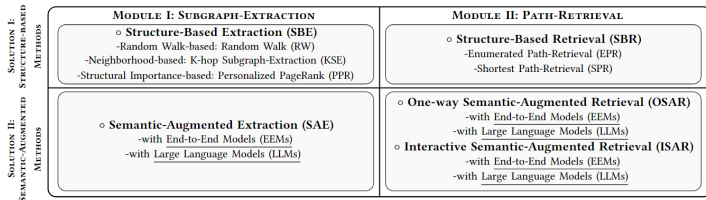
# LEGO-GraphRAG Architecture



Figure: LEGO-GraphRAG Framework

- Two phase: Subgraph-Extraction; Path-Retrieval;
- Two method: Structure-based; Semantic-augmented;

# Stucture-Based Extraction(SBE)



**MODULE I: SUBGRAPH-EXTRACTION**

○ **Structure-Based Extraction (SBE)**
-Random Walk-based: Random Walk (RW)
-Neighborhood-based: K-hop Subgraph-Extraction (KSE)
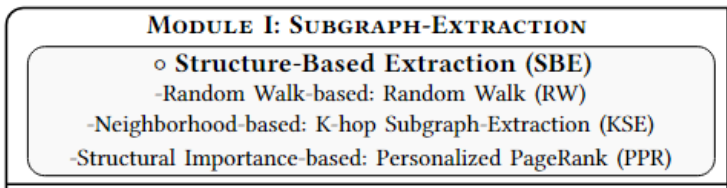-Structural Importance-based: Personalized PageRank (PPR)

Figure: Stucture-Based Extraction

- Random Walk: randomly selecting edges and nodes at each step.
- K-hop Subgraph-Extraction: distance $d(v_i^q, v_j) \leq K$.
- Personalized PageRank: assign an importance score to each node.

# Semantic-Augmented Extraction (SAE)



○ **Semantic-Augmented Extraction (SAE)**
  -with End-to-End Models (EEMs)
  -with Large Language Models (LLMs)

Figure: Semantic-Augmented Extraction

- With EEMs: compute semantic relevance.
  - pre-filtering: use SBE extracts a smaller subgraph.
  - subgraph pruning: node pruning, edge pruning, and triple pruning.
- With LLMs: evaluate and filter the semantic relevance of the initially extracted subgraph.
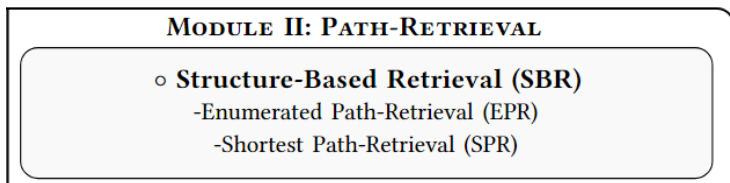
# Structure-Based Retrieval (SBR)

Figure: Structure-Based Retrieval

- Enumerated Path Retrieval: enumerates all possible paths from an entity $v_i^{(q)} \in \epsilon_q$.
- Shortest Path Retrieval: identifies all shortest paths from an entity $v_i^{(q)} \in \epsilon_q$.

# One-way/Interactive Semantic-Augmented Retrieval (OSAR/ISAR)



> ○ **One-way Semantic-Augmented Retrieval (OSAR)**
> -with End-to-End Models (EEMs)
> -with Large Language Models (LLMs)
> ○ **Interactive Semantic-Augmented Retrieval (ISAR)**
> -with End-to-End Models (EEMs)
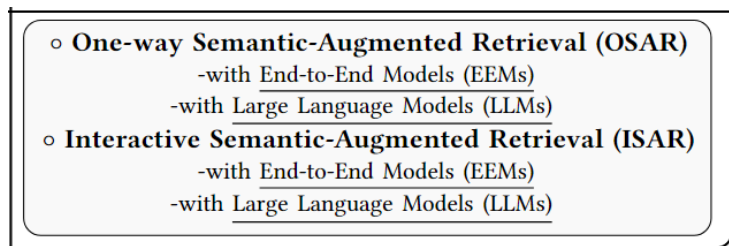> -with Large Language Models (LLMs)

Figure: One-way/Interactive Semantic-Augmented Retrieval

- One-way Semantic-Augmented Retrieval: evaluate and select the $N_p$ most relevant paths.
- Interactive Semantic-Augmented Retrieval: at each step, potential path extensions are evaluated for semantic relevance.

Table 2: Five Groups of Instances under the LEGO-GraphRAG Framework

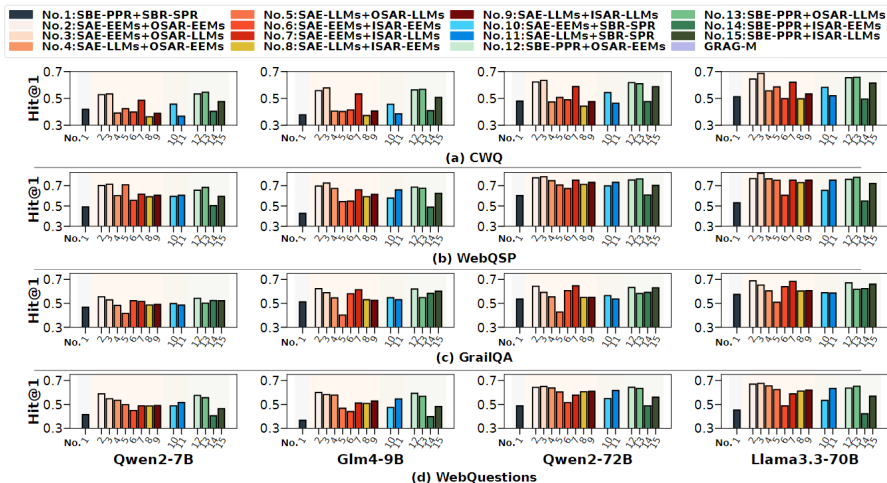| Group | Subgraph-Extraction | Path-Retrieval | Implemented Instances |
|---|---|---|---|
| Structure-based Methods on Both Modules (Group (I): *SBE & SBR*) | SBE (-RW/KSE/PPR) | SBR (-EPR/SPR) | Our Instance: No.1 |
| Semantic-Augmented Methods on Both Modules (Group (II): *SAE & I/OSAR*) | SAE (-EEMs/LLMs) | OSAR (-EEMs/LLMs) | Our Instances: No.2, 3, 4, 5 |
| | SAE (-EEMs/LLMs) | ISAR (-EEMs/LLMs) | GCR (arXiv24) [72] Our Instances: No.6, 7, 8, 9 |
| Semantic-Augmented Methods on Subgraph-Extraction (Group (III): *SAE & SBR*) | SAE (-EEMs/LLMs) | SBR (-EPR/SPR) | RoG (ICLR24) [71] GSR (EMNLP24) [45] Our Instances: No.10, 11 |
| Semantic-Augmented Methods on Path-Retrieval (Group (IV): *SBE & I/OSAR*) | SBE (-RW/KSE/PPR) | OSAR (-EEMs/LLMs) | Our Instances: No.12, 13 |
| | SBE (-RW/KSE/PPR) | ISAR (-EEMs/LLMs) | StructGPT (EMNLP23) [50] Our Instances: No.14, 15 |
| Without Subgraph-Extraction Modules (Group (V): *SBR or I/OSAR*) | None | SBR (-EPR/SPR) | - |
| | None | OSAR (-EEMs/LLMs) | KELP (ACL24) [65] |
| | None | ISAR (-EEMs/LLMs) | ToG (ICLR24) [98], DoG (arXiv24) [74] |

- Defined five types of GraphRAG instance groups.

## Table 3: Existing Instances vs. LEGO-GraphRAG Instances

| GraphRAG Instances | WebQSP | | CWQ | |
|---|---|---|---|---|
| | Hits@1 | Recall | Hits@1 | Recall |
| RoG [71] (RoG planning w/ChatGPT) | 81.51 | 71.60 | 52.68 | 48.51 |
| LEGO-RoG (RoG planning w/ChatGPT) | 82.79 | 64.41 | 56.06 | 49.76 |
| KELP [65] (one-hop w/gpt-4o-mini) | 31.06 | - | 14.16 | - |
| LEGO-KELP (one-hop w/gpt-4o-mini) | 77.36 | 63.99 | 48.65 | 43.88 |
| ToG [98] (w/Llama3-8B) | 59.76 | 43.05 | 36.97 | 32.69 |
| LEGO-ToG (w/Llama3-8B) | 66.44 | 44.77 | 40.26 | 33.63 |

- exhibit performance comparable to their original counterparts.
- KELP implementation outperforms the original.

# Reasoning Performance



Legend:
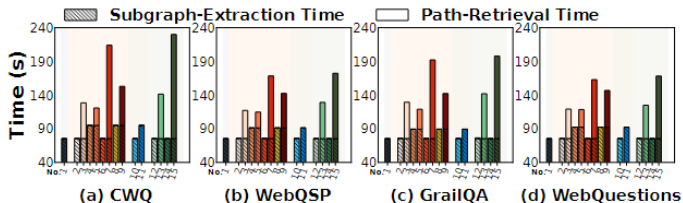- No.1:SBE-PPR+SBR-SPR
- No.2:SAE-EEMs+OSAR-EEMs
- No.3:SAE-EEMs+OSAR-LLMs
- No.4:SAE-LLMs+OSAR-EEMs
- No.5:SAE-LLMs+OSAR-LLMs
- No.6:SAE-EEMs+ISAR-EEMs
- No.7:SAE-EEMs+ISAR-LLMs
- No.8:SAE-LLMs+ISAR-EEMs
- No.9:SAE-LLMs+ISAR-LLMs
- No.10:SAE-EEMs+SBR-SPR
- No.11:SAE-LLMs+SBR-SPR
- No.12:SBE-PPR+OSAR-EEMs
- No.13:SBE-PPR+OSAR-LLMs
- No.14:SBE-PPR+ISAR-EEMs
- No.15:SBE-PPR+ISAR-LLMs
- GRAG-M

(a) CWQ

(b) WebQSP

(c) GrailQA

(d) WebQuestions

Qwen2-7B    Glm4-9B    Qwen2-72B    Llama3.3-70B

Figure 5: Runtime of SE and PR Modules for GraphRAG Instances

- Subgraphextraction is the primary bottleneck in GraphRAG runtime.
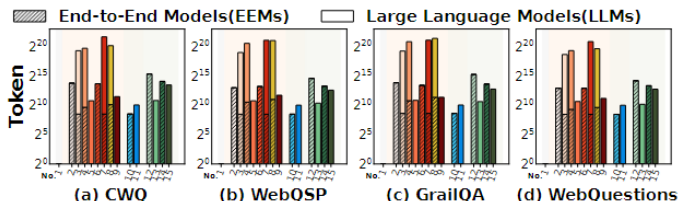- ISAR-LLMs methods in the path-retrieval may result in unacceptably low query efficiency.
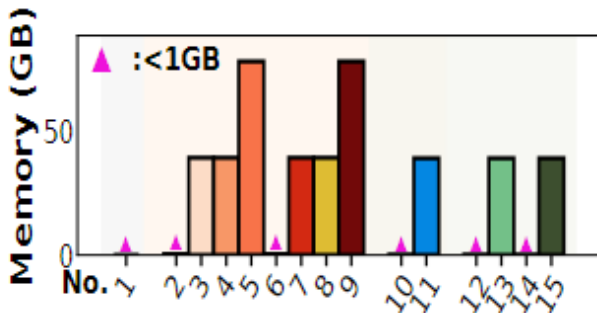
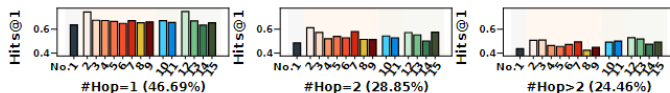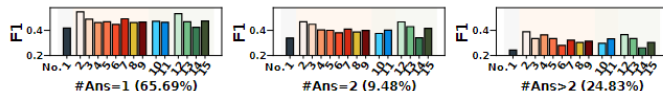**Figure 6:** Token Costs for EEMs and LLMs in GraphRAG Instances
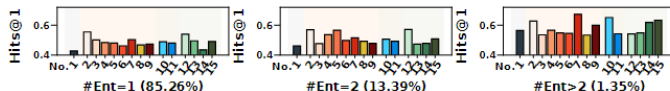
**(a) Queries with Varied Number of Reasoning Hops**



**(b) Queries with Varied Number of Answers**



**(c) Queries with Varied Number of Entities**

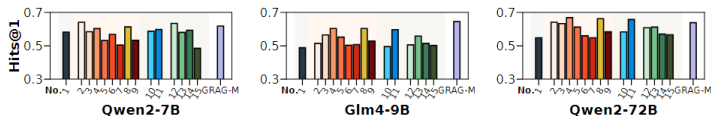**Figure 7: Instance Performance w.r.t. Queries**

Figure 3: Results of LEGO-GraphRAG Instances and GraphRAG-M Instance on MetaQA Dataset

- implement an instance (GRAG-M): performing community detection on the graph; precomputing textual summaries for each community;

**Table 6:** Computational acceleration of PPR

| Approximate PPR (Freebase) | | | Distributed PPR | |
|---|---|---|---|---|
| **Method** | **Recall** | **Ave. Time (s)** | **Method** | **Speedup** |
| RBS [107] | 0.31 | **0.51** | HGPA [31] | 3.4–4.1× |
| Fora [110] | 0.18 | 7.61 | PAFO [109] | 58.7× |
| TopPPR [113] | 0.44 | 42.08 | Delta-Push [39] | **123–162×** |
| **Standard PPR** | **0.96** | 75.29 | **Standard PPR** | 1× (baseline) |

**Table 7:** Performance and Cost Analysis of the SE Acceleration Methods on Freebase (About 100M Nodes, 300M Edges)

| Method | Pre-time | Online-time | Recall | F1 | WCC |
|---|---|---|---|---|---|
| Precomputation | 19373s | 19.02s | 0.52 | 0.0021 | **1** |
| Vector Database | 14570s | **0.85s** | 0.65 | 0.0023 | 12.86 |
| **Standard PPR** | **0s** | 75.29s | **0.96** | **0.0038** | **1** |

**Table 8: Performance of strategies designed to mitigate the limitations of LLM-generated outputs**

| Method | CWQ | WebQSP | GrailQA | WebQuestions |
|---|---|---|---|---|
| SBE+OSAR-LLMs | 0.304 | 0.401 | 0.401 | 0.328 |
| SBE+OSAR-LLMs (M-LLMs) | 0.349 | **0.438** | 0.381 | 0.348 |
| SBE+OSAR-LLMs (EEMs-S) | **0.381** | 0.427 | **0.476** | **0.353** |