



Graph-learning Agents

记忆组织方向

2025年8月4日

马斌

目录

1

论文调研概况

2

论文的理解与分析

3

论文精读

4

总结

4

引用目录



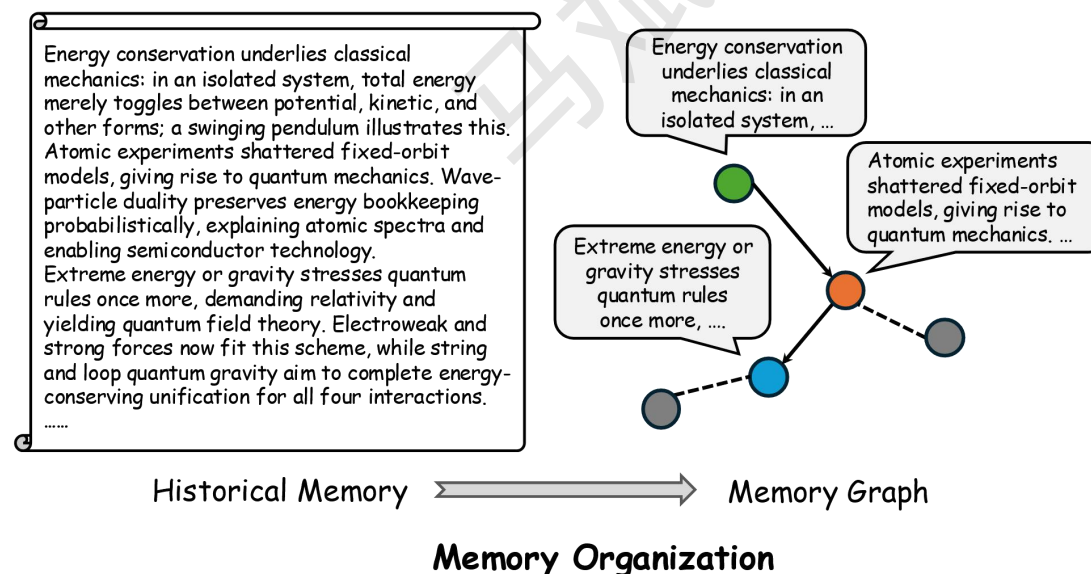
1 论文调研概况



■ 调研背景:

记忆(Memory) 是一项至关重要的能力，它允许智能体存储和回忆过去的经验或相关知识，它允许代理积累经验，从而促进更明智和适当的行动。

而 **基于图的记忆组织(Graph-based Memory Organization)**，增强了代理识别非显式的模式和连接的能力，从而改进了其在动态环境中的执行能力。





1 论文调研概况



■ 调研目录:

- [AriGraph \[74\]](#) --- 三元组的结构化知识图谱
- [IKG \[75\]](#) --- 场景_1: 工业多模态数据
- [Graphusion \[76\]](#) --- 场景_2: 教育
- [MemGraph \[77\]](#) --- 场景_3: 专利分析
- [Mind Map \[78\]](#) --- 场景_4: 网络搜索
- [StructuralMemory \[79\]](#) --- 多种场景下的四种记忆结构: 块、知识三元组、原子事实和摘要
- [DAMCS \[80\]](#) --- 分层知识图谱, 允许代理记录低级细节和高级摘要
- [GraphRAG \[81\]](#) --- 混合知识图谱, 对实体级和社区级节点进行编码
- [KG-Retriever \[82\]](#) --- 分层知识图谱, 节点由实体级和文档级



2 论文的理解与分析

● 论文: [AriGraph \[74\]](#)

■ 背景:

常用RAG -> 非结构化问题 -> KG解决

AriGraph: 更好整合包含结构化、非结构化的记忆

■ 概述:

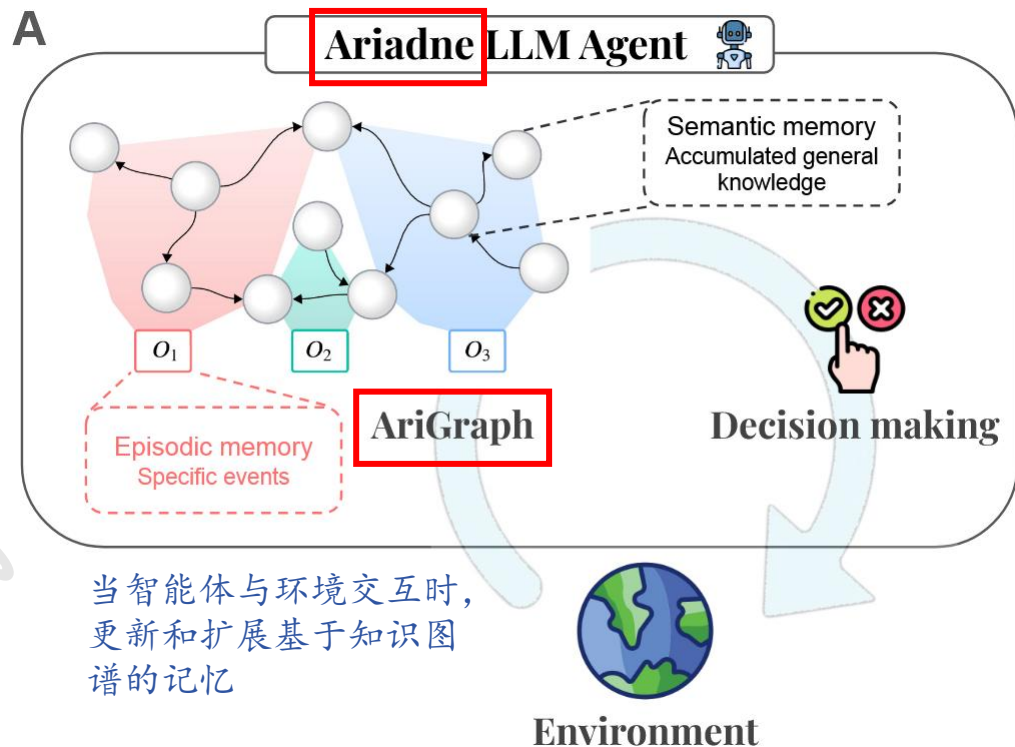
语义记忆知识网络 -> 知识图谱

情景记忆 -> 情景边(记录时序关系)

-> 情景节点 {包含改时刻所有新语义边}

■ 实验:

	TextWorld 环境			NetHack 环境			多跳问答
	Treasure Hunt	Treasure Hunt Hard	Treasure Hunt Hardest	Cooking	Cooking Hard	Cooking Hardest	Cleaning
Full History	0.47	-	-	0.18	-	-	0.05
Summary	0.33	0.17	-	0.52	0.21	-	0.35
RAG	0.33	0.17	-	0.36	0.17	-	0.39
Reflexion	0.93	-	-	1.0	-	-	0.27
Simulacra	0.4	-	-	0.3	-	-	0.7
AriGraph	1.0	1.0	1.0	1.0	1.0	0.65	0.79
AriGraph w/o exploration	0.87	-	-	0.87	-	-	0.76
AriGraph w/o episodic	1.0	0.67	-	0.64	0.45	-	0.92
AriGraph LLaMA-3-70B	0.47	-	-	0.67	-	-	0.5
Human Top-3	1.0	-	-	1.0	-	-	1.0
Human All	0.96	-	-	0.32	-	-	0.59



● 认知科学知识补充:

语义记忆包含有关世界的事实知识;

情景记忆则涉及个人经历, 其中通常包含更丰富、更详细的信息。

语义知识建立在情景记忆的基础上

-> 允许集成各种记忆



2 论文的理解与分析

- 论文: [IKG \[75\]](#), [Graphusion \[76\]](#), [MemGraph \[77\]](#), [Mind Map \[78\]](#)

- [IKG \[75\]](#)

工业知识图谱 (IKG)-资源自配置, 多智能体强化学习 (MARL)-过程自优化。

评估: 多UR5机器人协同到达任务的任务完成率、碰撞率、路径优化效率的真实案例。

- [Graphusion \[76\]](#)

通过种子实体引导、三元组提取、融合实现从自由文本构建全局视角的KG

评估: 数据集: ACL论文摘要 (2017-2023)、TutorQA (1200条教育QA)

对比模型: Zero-shot LLM、RAG、Graphusion

评估指标: 人工评分 (实体/关系质量)、链路预测F1、TutorQA任务准确率

- [MemGraph \[77\]](#)

利用实体和本体层级关系优化专利语义理解与检索。

评估: 数据集: PatentMatch (1000条专利匹配问题)

对比模型: 领域微调LLM (MoZi、PatentGPT)、RAG

消融实验: 验证Z_IR (实体增强检索) 和Z_Gen (本体增强生成) 的性能提升效果。

- [Mind Map \[78\]](#)

评估: 数据集: Humanity's Last Exam、GPQA、GAIA、FreshWiki

消融实验: 验证工具组合 (Web-Search + Code + Mind-Map最优) 及Web-Search组件有效性

案例研究: 医疗决策问题中工具协同应用

不同层度的应用场景:

场景_1: 工业多模态数据

场景_2: 教育

场景_3: 专利分析

场景_4: 网络搜索

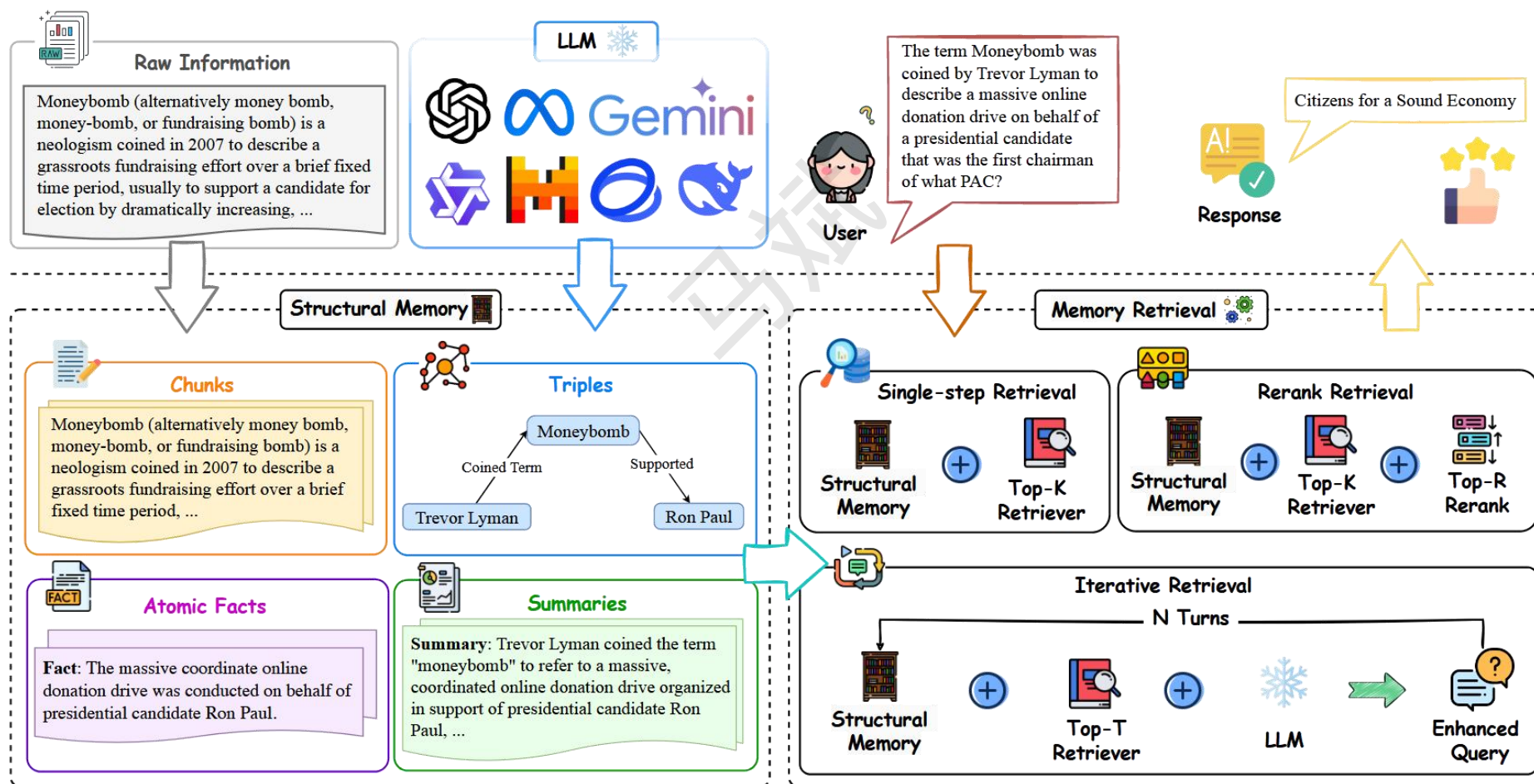


2 论文的理解与分析



● 论文: [StructuralMemory \[79\]](#)

- 研究LLM代理中, 4种单一记忆结构(chunks、知识三元组、原子事实、摘要)及混合结构, 3种记忆检索方法(单步检索、重排序、迭代检索), 在4类任务(多跳QA、单跳QA、对话理解、阅读理解)的6个数据集上实验的对性能的影响





2 论文的理解与分析

- 论文: [DAMCS \[80\]](#), [GraphRAG \[81\]](#)

- [DAMCS \[80\]](#)

去中心化多智能体协作框架。

评估: Multi-agent Crafter (MAC) 环境下协作收集钻石, 完成任务的平均步数。

- [GraphRAG \[81\]](#)

微软研究院

两阶段图索引构建: 使用LLM

1. 从源文档中提取实体, 构建实体知识图谱;
2. 为所有密切相关的实体组预生成社区摘要。

评估: 数据集采用Podcast transcripts、News articles (约100万toke)



2 论文的理解与分析



● 论文: [KG-Retriever \[82\]](#)

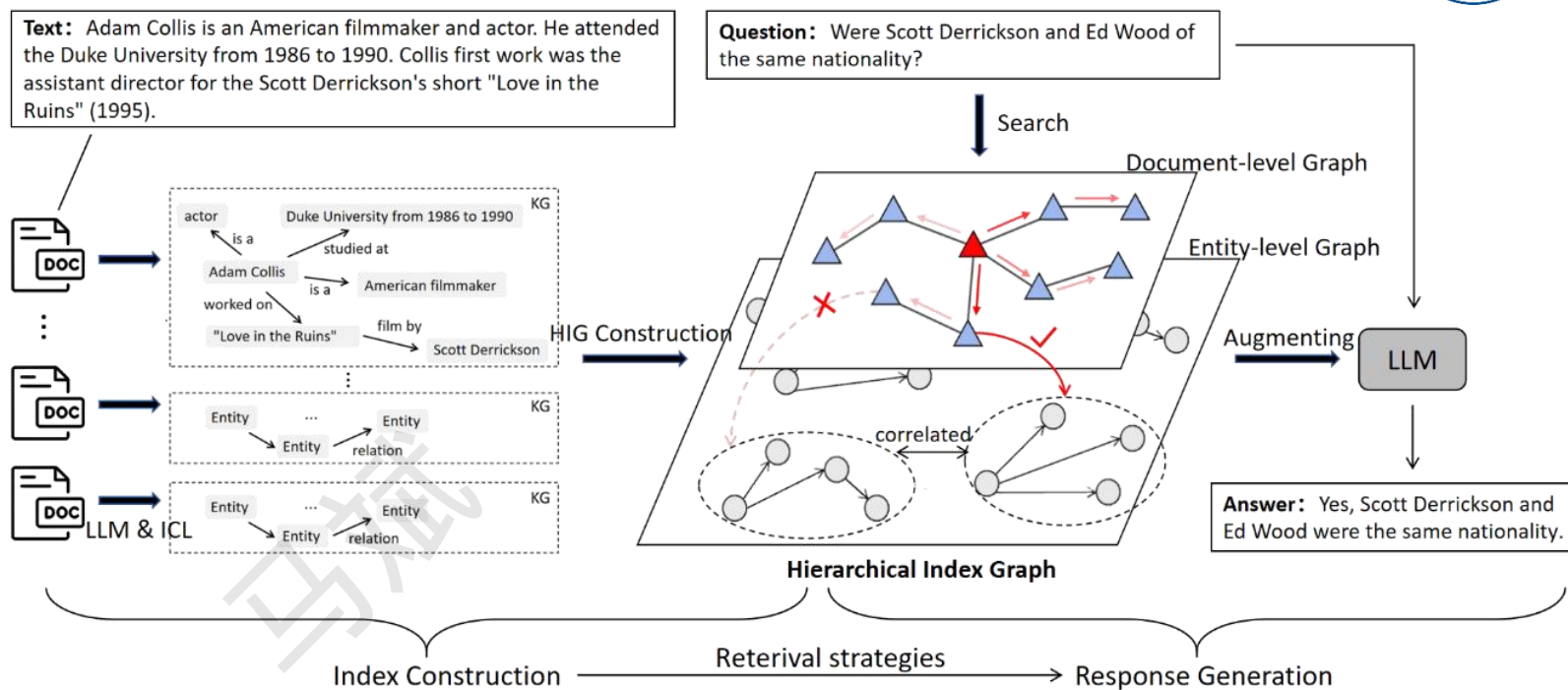
■ 概述: 基于KG的分层知识检索器。

□ 知识图谱层

文档 -LLM→ 知识三元组

□ 协作文档层

基于语义相似性分层,
将相关文档相连。



Datasets	HotpotQA		MuSiQue		2WikiMultiHopQA		CRUD-QA1			CRUD-QA2		
	EM	Time	EM	Time	EM	Time	BLEU	Rouge-L	Time	BLEU	Rouge-L	Time
Naive LLM	0.102	0.69s	0.040	0.31s	0.170	0.28s	0.073	0.24	1.75s	0.079	0.237	0.95s
COT + LLM	0.172	2.11s	0.060	1.32s	0.180	1.45s	0.035	0.146	2.11s	0.030	0.151	4.06s
Graph-guided reasoning	0.197	40.03s	0.070	42.30s	0.210	39.08s	0.311	0.393	34.59s	0.095	0.244	27.18s
BM25	0.236	1.25s	0.070	0.67s	0.250	0.73s	0.209	0.509	1.08s	0.069	0.233	1.37s
DenseRetriever	0.282	0.25s	0.050	0.35s	0.220	0.45s	0.23	0.395	1.06s	0.155	0.259	1.44s
ITRG (5-Iteration)	0.306	10.99s	0.120	9.60s	0.320	6.30s	0.333	0.457	14.7s	0.172	0.277	13.55s
ITER-RETGEN (3-Iteration)	0.323	6.65s	0.170	9.71s	0.290	10.03s	0.24	0.551	6.65s	0.069	0.236	8.95s
KGP (3-Iteration)	0.278	6.34s	0.140	6.98s	0.220	7.03s	0.275	0.376	7.07s	0.155	0.235	7.30s
KG-Retriever	0.328	0.93s	0.210	1.21s	0.350	0.74s	0.449	0.611	0.95s	0.233	0.353	1.46s
KG-Retriever (Attention)	0.322	1.14s	0.200	0.88s	0.340	0.83s	0.458	0.600	0.99s	0.239	0.357	1.63s
KG-Retriever (Multi-Hop)	0.328	1.19s	0.210	1.03s	0.350	0.83s	0.458	0.600	1.15s	0.238	0.354	2.20s

■ 评估:

数据集: 5个公开QA数据集。

基线: 传统RAG、高级RAG(其他研究)、图-RAG(本文)。

指标: EM(二元评估匹配率)、BLEU(生成文本与参考译文的n-gram相似度)、Rouge-L(文本摘要召回率)、Time(生成响应平均用时)。

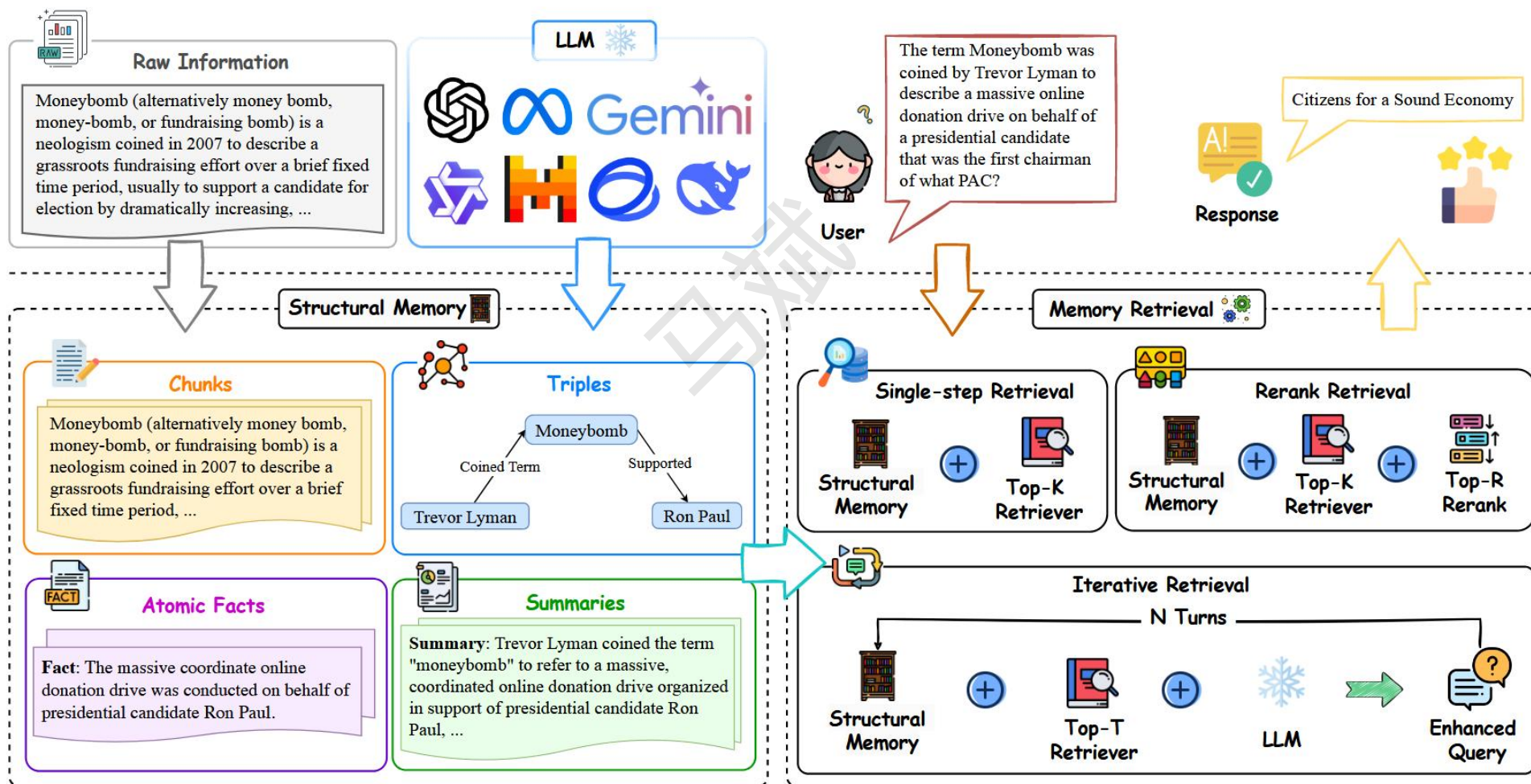


3 论文精读 - 1



● 论文: [StructuralMemory \[79\]](#)

■ 相关背景: 记忆模块(结构记忆生成、记忆检索方法、答案生成)





3 论文精读 - 1



● 论文: [StructuralMemory \[79\]](#)

■ 研究内容

■ 4种记忆结构:

□ 块 Chunks

$$\mathcal{C}_q(\mathcal{D}_q) = \{c_1, c_2, \dots, c_j\}$$

□ 知识三元组 Knowledge triples

$$\mathcal{T}_q = \text{LLM}(\mathcal{D}_q, \mathcal{P}_{\mathcal{T}})$$

□ 原子事实 Atomic facts

$$\mathcal{A}_q = \text{LLM}(\mathcal{D}_q, \mathcal{P}_{\mathcal{A}}).$$

□ 原子事实 Summaries

$$\mathcal{S}_q = \text{LLM}(\mathcal{D}_q, \mathcal{P}_{\mathcal{S}})$$

□ 混合结构 Mixed

$$\mathcal{M}_q^{\text{Mixed}} = \mathcal{C}_q \cup \mathcal{T}_q \cup \mathcal{A}_q \cup \mathcal{S}_q$$

■ 3种记忆检索方法:

□ 单步检索 Single-step retrieval

$$\mathcal{M}_r = \text{Retriever}(q, \mathcal{M}_q, K)$$

□ 重排序 Reranking

$$\mathcal{M}_r = \text{LLM}(q, \mathcal{M}_i, R, \mathcal{P}_R), \text{ where } \mathcal{M}_i = \text{Retriever}(q, \mathcal{M}_q, K)$$

□ 迭代检索 Iterative retrieval

$$q_j = \text{LLM}(\mathcal{M}_j, \mathcal{P}_{\text{Refine}}), \text{ where } \mathcal{M}_j = \text{Retriever}(q_{j-1}, \mathcal{M}_q, T)$$

■ 2种答案生成方法:

$$\mathcal{M}_r = \text{Retriever}(q_N, \mathcal{M}_q, K)$$

□ Memory-Only

直接将检索结果 \mathcal{M}_r 作为上下文

□ Memory-Doc

根据检索结果定位到原始文档 \mathcal{D}_q 作为上下文



3 论文精读 - 1



● 论文: [StructuralMemory \[79\]](#)

■ 实验设计

■ 4类任务:

- 多跳QA
- 单跳QA
- 对话理解
- 阅读理解

■ 6个数据集

Task	Dataset	Avg. # Tokens	# Samples
Multi-hop QA	HotpotQA	1,362	200
Multi-hop QA	2WikiMultihopQA	985	200
Multi-hop QA	MuSiQue	2,558	200
Single-hop QA	NarrativeQA	24,009	200
Dialogue Understanding	LoCoMo	24,375	191
Reading Comprehension	QuALITY	4,696	200



3 论文精读 - 1



● 论文: [StructuralMemory \[79\]](#)

■ 实验结果

Memory Structure	HotPotQA		2WikiMultihopQA		MuSiQue		NarrativeQA		LoCoMo		QuALITY
	EM	F1	EM	F1	EM	F1	EM	F1	EM	F1	ACC
Full Content	55.50	75.77	44.00	54.33	36.00	51.60	7.00	24.99	13.61	41.82	81.50
Single-step Retrieval											
Chunks	<u>61.50</u>	76.93	43.50	59.17	35.50	54.45	13.50	29.78	9.95	40.63	<u>76.00</u>
Triples	59.50	74.09	<u>44.50</u>	<u>60.82</u>	31.00	50.13	11.50	22.04	8.42	41.08	61.50
Atomic Facts	62.50	77.22	39.50	58.63	30.50	51.31	13.50	27.49	9.42	42.92	71.50
Summaries	57.00	74.81	42.00	57.21	<u>34.00</u>	<u>52.83</u>	16.50	32.93	10.99	44.94	<u>76.00</u>
Mixed	60.00	<u>77.10</u>	48.50	65.25	33.00	51.65	<u>14.50</u>	<u>29.86</u>	<u>10.47</u>	<u>44.73</u>	78.00
Reranking											
Chunks	<u>63.00</u>	77.35	<u>45.00</u>	<u>61.31</u>	37.00	55.32	16.00	<u>31.63</u>	<u>9.95</u>	43.47	78.50
Triples	61.00	76.75	43.50	55.43	26.50	42.05	10.00	20.65	8.83	41.82	60.00
Atomic Facts	<u>63.00</u>	<u>78.31</u>	40.50	59.31	28.50	49.95	<u>14.00</u>	28.19	8.90	44.27	67.50
Summaries	61.00	77.80	<u>45.00</u>	61.18	<u>35.50</u>	<u>54.59</u>	16.00	32.26	12.04	44.83	75.00
Mixed	65.00	78.58	45.50	61.77	34.00	52.45	11.98	28.02	9.42	<u>44.51</u>	<u>77.50</u>
Iterative Retrieval											
Chunks	63.00	79.10	46.50	62.13	37.00	56.78	<u>14.50</u>	<u>30.88</u>	<u>10.47</u>	<u>45.14</u>	<u>77.00</u>
Triples	64.00	78.78	<u>47.50</u>	62.06	<u>38.00</u>	55.93	10.50	21.67	9.47	41.41	60.50
Atomic Facts	<u>65.50</u>	<u>81.29</u>	44.00	<u>63.89</u>	34.50	<u>57.55</u>	<u>14.50</u>	28.28	9.95	43.62	67.50
Summaries	60.50	78.11	46.50	62.35	33.50	53.12	17.00	31.79	12.04	43.93	75.00
Mixed	67.00	82.11	51.00	68.15	39.00	61.38	12.50	28.36	7.85	45.25	79.50

Table 1: Overall Performance (%) of various memory structures utilizing different retrieval methods across six datasets. The best performance is marked in boldface, while the second-best performance is underlined.



3 论文精读 - 1

- 论文: [StructuralMemory \[79\]](#)

- 发现

1. 记忆结构方面:

混合记忆始终提供平衡的性能, 各类结构适合于不同任务场景。

块、摘要 -> 阅读理解、对话理解 (涉及冗长上下文的任务)

知识三元组、原子事实 -> 多跳QA、单跳QA

2. 抗噪声方面:

混合记忆还表现出对噪音的非凡恢复能力, 其次为块。

3. 记忆检索方法方面:

迭代检索始终是最有效的记忆检索方法。



3 论文精读 - 2



● 论文: [KG-Retriever \[82\]](#)

■ 相关背景

■ 传统RAG --- 多跳场景劣势

□ 索引

长文本->块->索引向量

□ 检索

根据语义相似度, 将待查询的嵌入与已向量化的文本块进行匹配(相同模型进行编码)

□ 生成

检索到的块+原始查询->LLM输入

■ 改进RAG -> 尝试结合检索与生成过程 --- 难应用于零样本设置

□ ITRG - 每次输出用于下次检索, 迭代

□ IRCOT - 检索与思维链(CoT)相互优化

■ 基于图的RAG -> 侧重索引阶段

□ KGP - 多文档->文档级KG

□ KG-Retriever(本文) - 多文档->提取实体关系->实体级KG



3 论文精读 - 2



● 论文: [KG-Retriever \[82\]](#)

■ KG-Retriever 结构:

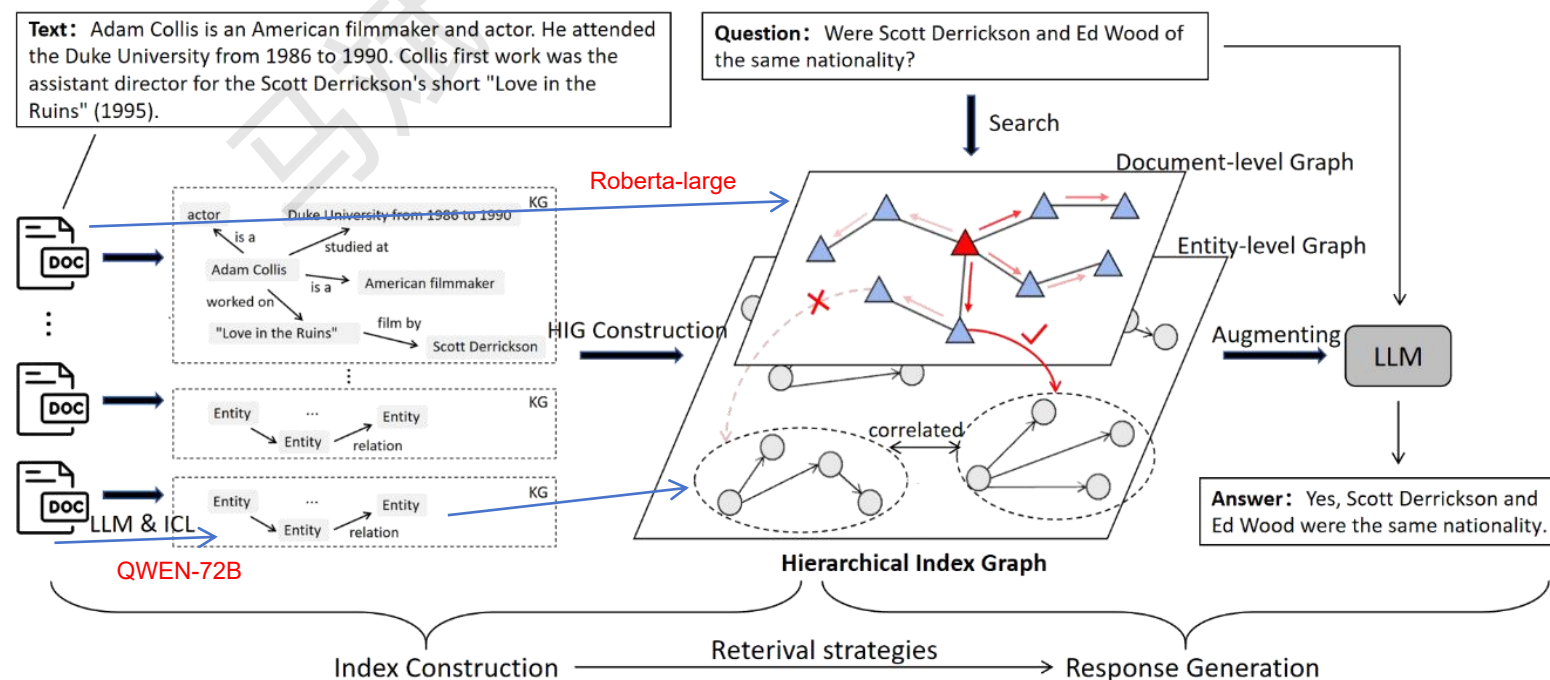
- 索引构建组件 --- for 文档图 \rightarrow 对应索引, 形成分层索引图 (HIG) - 包含文档级和实体级
- 知识检索组件 --- HIG 中检索 - Ret1 \rightarrow 相关文档 - Ret2 \rightarrow 相关文档中的三元组
- 响应生成组件 --- 三元组 + 原始查询 \rightarrow LLM \rightarrow 输出结果

□ 文档级检索 Ret1

单跳协作、注意力协作、多跳协作

□ KG 级检索 Ret2

$$\text{Retrieved?} = \begin{cases} \text{yes} & \text{if } w * \text{CosSim}(\mathbf{v}_e, \mathbf{v}_q) > \lambda; \\ \text{no} & \text{else,} \end{cases}$$





3 论文精读 - 2



● 论文: [KG-Retriever \[82\]](#)

■ 实验设计

- 数据集: 5个公开QA数据集。
- 基线: 传统RAG、高级RAG(其他研究)、Graph-RAG(本文)。
- 指标: EM(二元评估匹配率)、BLEU(生成文本与参考译文的n-gram相似度)、Rouge-L(文本摘要召回率)、Time(生成响应平均用时)。

Datasets	HotpotQA		MuSiQue		2WikiMultiHopQA		CRUD-QA1			CRUD-QA2		
	EM	Time	EM	Time	EM	Time	BLEU	Rouge-L	Time	BLEU	Rouge-L	Time
Naive LLM	0.102	0.69s	0.040	0.31s	0.170	0.28s	0.073	0.24	1.75s	0.079	0.237	0.95s
COT + LLM	0.172	2.11s	0.060	1.32s	0.180	1.45s	0.035	0.146	2.11s	0.030	0.151	4.06s
Graph-guided reasoning	0.197	40.03s	0.070	42.30s	0.210	39.08s	0.311	0.393	34.59s	0.095	0.244	27.18s
BM25	0.236	1.25s	0.070	0.67s	0.250	0.73s	0.209	0.509	1.08s	0.069	0.233	1.37s
DenseRetriever	0.282	0.25s	0.050	0.35s	0.220	0.45s	0.23	0.395	1.06s	0.155	0.259	1.44s
ITRG (5-Iteration)	0.306	10.99s	0.120	9.60s	0.320	6.30s	0.333	0.457	14.7s	0.172	0.277	13.55s
ITER-RETGEN (3-Iteration)	0.323	6.65s	0.170	9.71s	0.290	10.03s	0.24	0.551	6.65s	0.069	0.236	8.95s
KGP (3-Iteration)	0.278	6.34s	0.140	6.98s	0.220	7.03s	0.275	0.376	7.07s	0.155	0.235	7.30s
KG-Retriever	0.328	0.93s	0.210	1.21s	0.350	0.74s	0.449	0.611	0.95s	0.233	0.353	1.46s
KG-Retriever (Attention)	0.322	1.14s	0.200	0.88s	0.340	0.83s	0.458	0.600	0.99s	0.239	0.357	1.63s
KG-Retriever (Multi-Hop)	0.328	1.19s	0.210	1.03s	0.350	0.83s	0.458	0.600	1.15s	0.238	0.354	2.20s

■ 本文局限性

- 对骨干模型的推理能力有要求;
- 构建HIG过程(离线)的成本大;
- 静态索引结构可能并非适用于动态语料库。

■ 发现

- 本文研究以KGP(2024)为主线, 加以其他研究进行改进;
- 本文主要针对小米公司需求研发。



4 总结



本次调研论文均以利用图技术优化智能体记忆组织为中心，

主要可分为 优化记忆表示、优化记忆结构 两种策略。

且多以 零提示模型、传统RAG 为对照评测效果，用 消融实验 的方法得出各因素的影响。

研究探索了一系列记忆表示， 优化智能体的长期记忆组织

StructuralMemory [79]
2024-12-17

汇总对比四种记忆表示。

AriGraph [74]
2025-05-15

结合认知科学，
实现整合结构化和非结
构化知识。

应用场景

IKG [75]
2021-10-01

工业多模态数据

【该领域初尝试】

Graphusion [76]
2025-02-03

教育

MemGraph [77]
2025-04-21

专利分析

Mind Map [78]
2025-07-14

网络搜索

研究探索了一系列记忆结构， 实现不同的语义分辨率下的整体检索

DAMCS [80]
2025-02-08

分层图

GraphRAG [81]
2025-02-19

混合图

KG-Retriever [82]
2025-05-05

分层图

基于KG的文档协作管理，更高效地
调用相关文档内容。



5 引用目录



- [74] P. Anokhin, N. Semenov, A. Sorokin, D. Evseev, M. Burtsev, and E. Burnaev, “Arigraph: Learning knowledge graph world models with episodic memory for llm agents,” arXiv preprint arXiv:2407.04363, 2024.
- [75] P. Zheng, L. Xia, C. Li, X. Li, and B. Liu, “Towards self-x cognitive manufacturing network: An industrial knowledge graph-based multiagent reinforcement learning approach,” Journal of Manufacturing Systems, vol. 61, pp. 16–26, 2021.
- [76] R. Yang, B. Yang, A. Feng, S. Ouyang, M. Blum, T. She, Y. Jiang, F. Lecue, J. Lu, and I. Li, “Graphusion: a rag framework for knowledge graph construction with a global perspective,” arXiv preprint arXiv:2410.17600, 2024.
- [77] Q. Xiong, Z. Xu, Z. Liu, M. Wang, Z. Chen, Y. Sun, Y. Gu, X. Li, and G. Yu, “Enhancing the patent matching capability of large language models via memory graph,” in Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval, 2025.
- [78] J. Wu, J. Zhu, and Y. Liu, “Agentic reasoning: Reasoning llms with tools for the deep research,” arXiv preprint arXiv:2502.04644, 2025.
- [79] R. Zeng, J. Fang, S. Liu, and Z. Meng, “On the structural memory of llm agents,” arXiv preprint arXiv:2412.15266, 2024.
- [80] H. Yang, J. Chen, M. Siew, T. Llorido-Botran, and C. Joe-Wong, “Llm-powered decentralized generative agents with adaptive hierarchical knowledge graph for cooperative planning,” arXiv preprint arXiv:2502.05453, 2025.
- [81] D. Edge, H. Trinh, N. Cheng, J. Bradley, A. Chao, A. Mody, S. Truitt, D. Metropolitansky, R. O. Ness, and J. Larson, “From local to global: A graph rag approach to query-focused summarization,” arXiv preprint arXiv:2404.16130, 2024.
- [82] W. Chen, T. Bai, J. Su, J. Luan, W. Liu, and C. Shi, “Kg-retriever: Efficient knowledge indexing for retrieval-augmented large language models,” arXiv preprint arXiv:2412.05547, 2024.



谢谢！

2025年8月4日

马斌

