

# Pre-Training

Text corpus



(Self-supervised)  
Training

Pretrained LM



Adaptation

Tasks

Question  
Answering



Text  
Classification

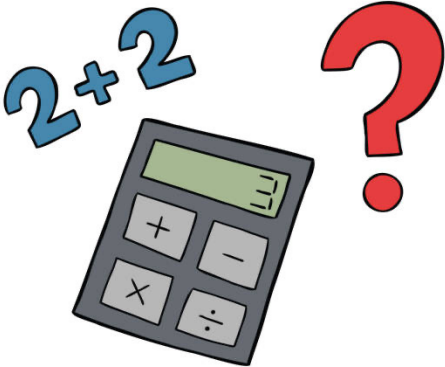


Information  
Retrieval



⋮

# Limitations of LLMs

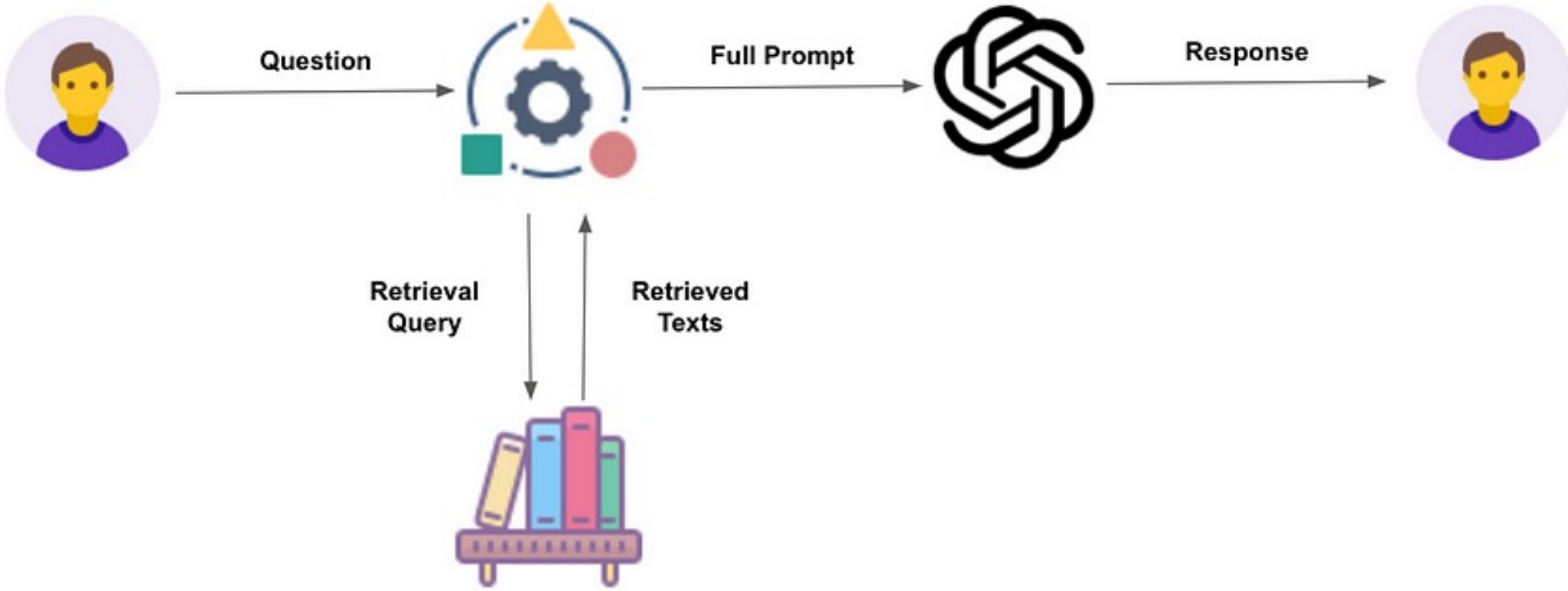


# Retrieval-Augmented Generation

Private Knowledge

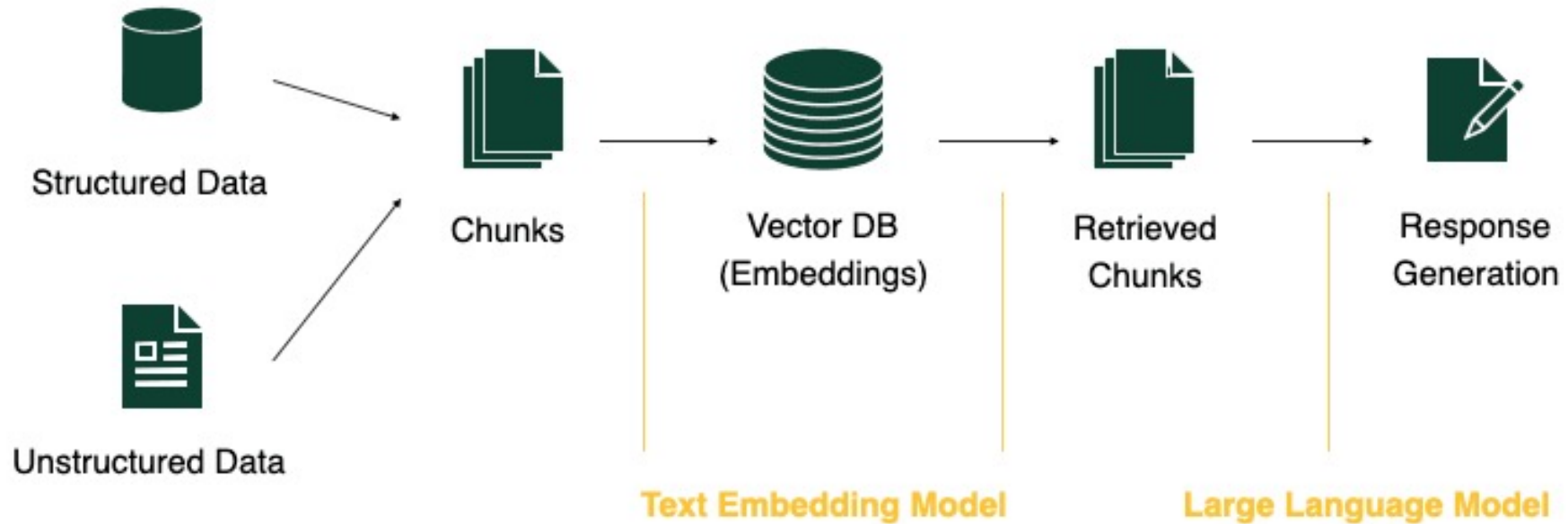


# Retrieval-Augmented Generation

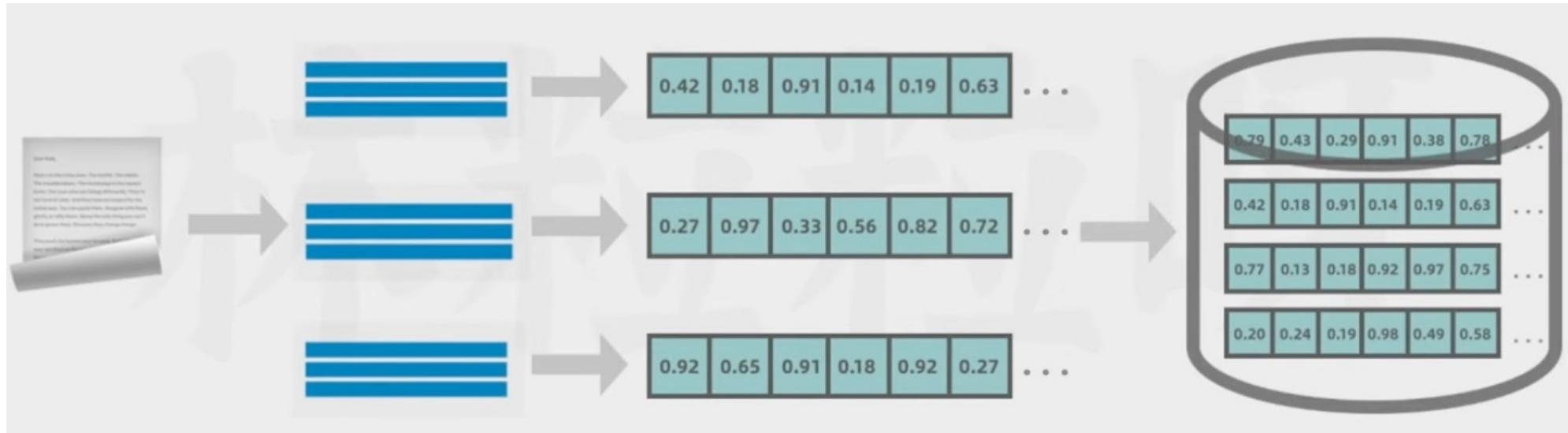


# Retrieval-Augmented Generation (RAG)

## Simple RAG



# Retrieval-Augmented Generation (RAG)

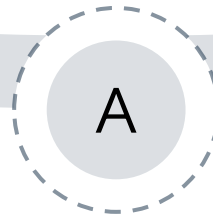


# RAG (Retrieval-Augmented Generation)



## (Retrieval)

Finds the most relevant information from external knowledge bases or document collections. By vectorizing and storing text data, the system efficiently retrieves the best matches using query vectors from user input.



## (Augmentation)

Uses retrieved information as additional context for the language model, enhancing its understanding and improving answer accuracy.



## (Generation)

Powered by LLMs, it generates coherent, context-aware, high-quality answers by combining user queries, retrieved information, and its language generation capabilities.

# Naive RAG Limitation

1

**Relies on flat data, limiting complex relationships.**

Naive RAG uses text chunks, ignoring entity relations and complex cross-document dependencies.

2

**Lacks context awareness, leading to incoherent answers.**

Retrieval feeds generation linearly, lacking global context integration.

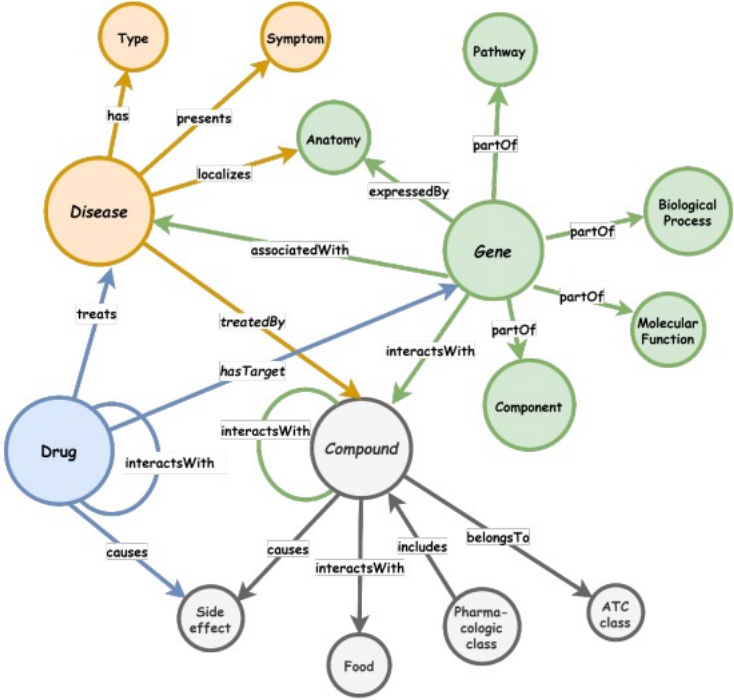
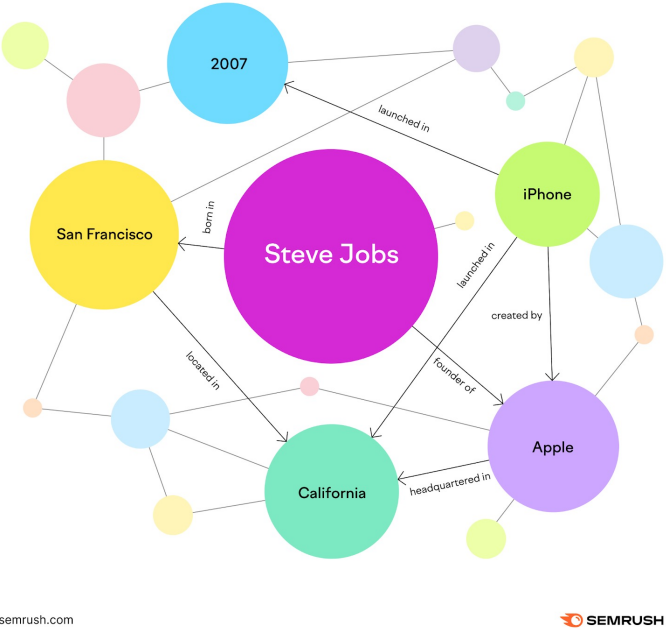
3

**Text blocks are redundant with irrelevant content.**

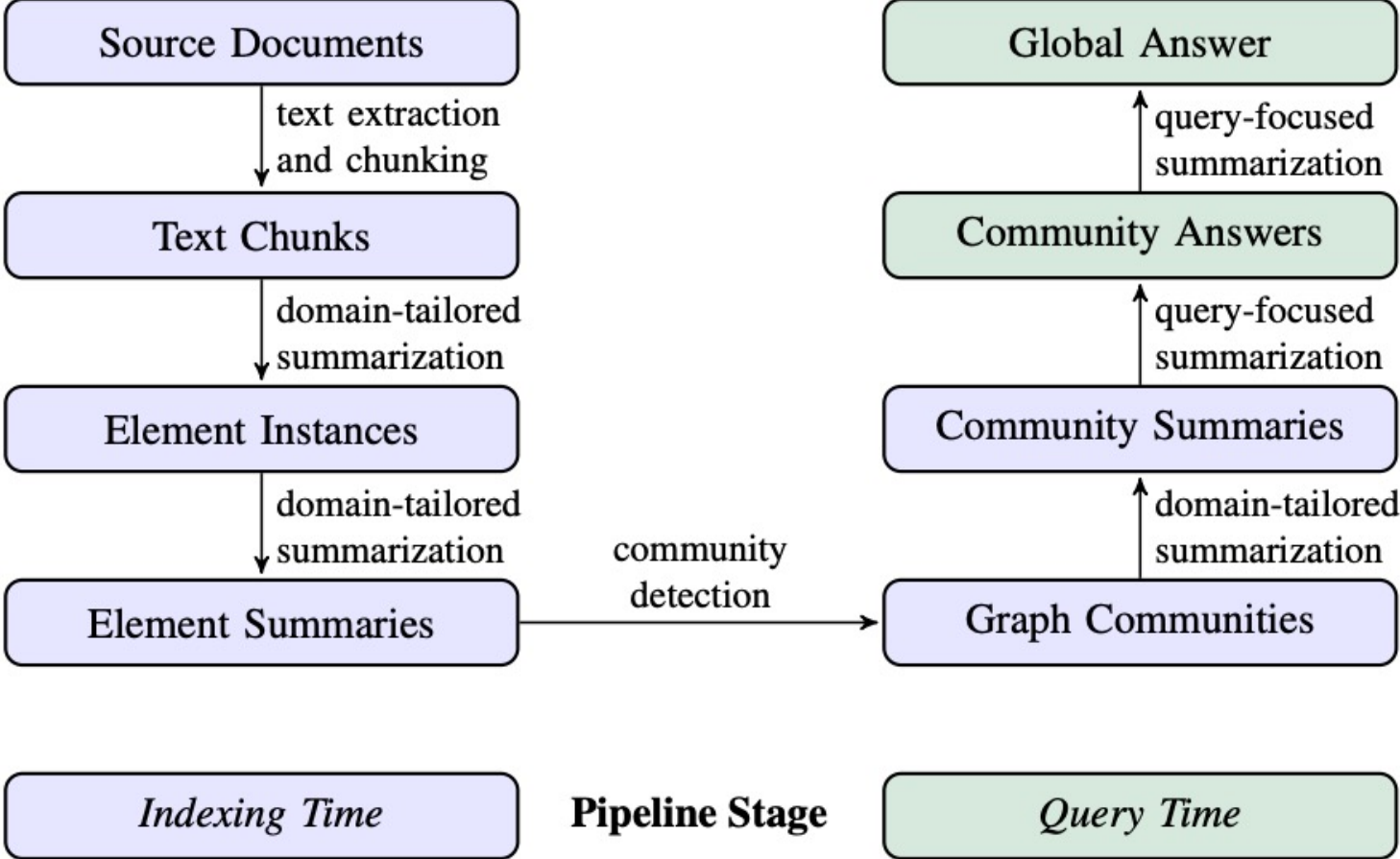
Naive RAG retrieval often includes irrelevant information, increasing computational load and reducing answer accuracy.

# Graph RAG

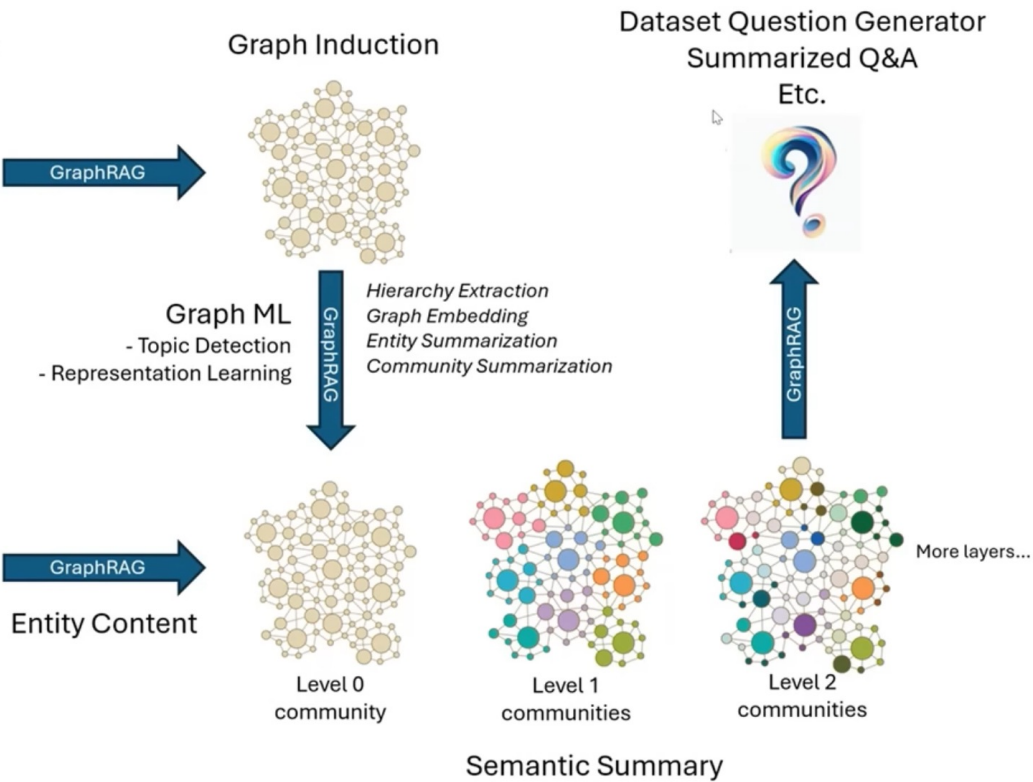
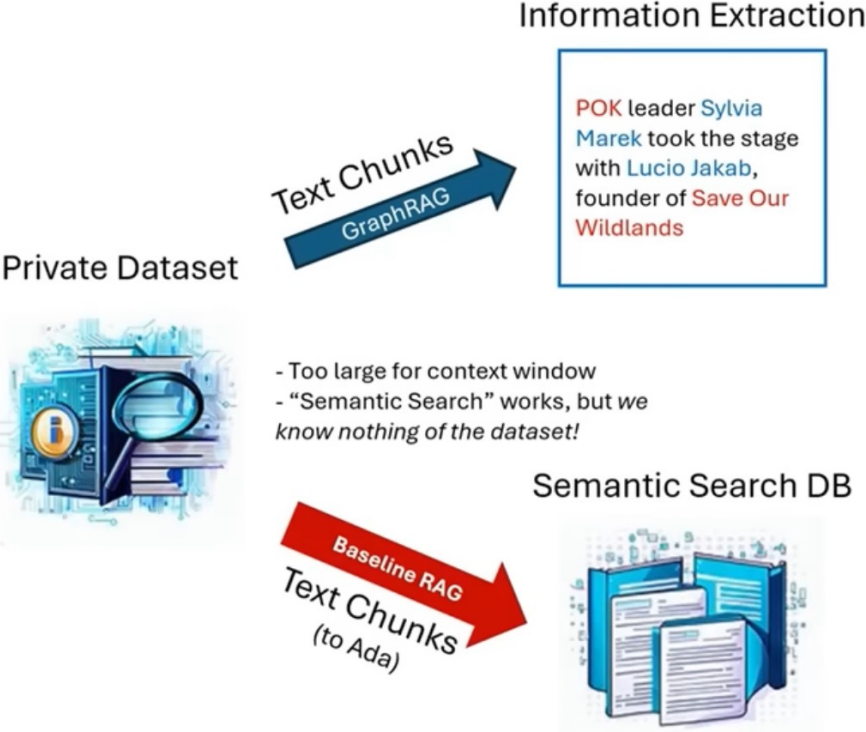
Unlike a baseline RAG that uses a vector database to retrieve semantically similar text, GraphRAG enhances RAG by incorporating knowledge graphs (KGs). Knowledge graphs are data structures that store and link related or unrelated data based on their relationships.



# Graph RAG



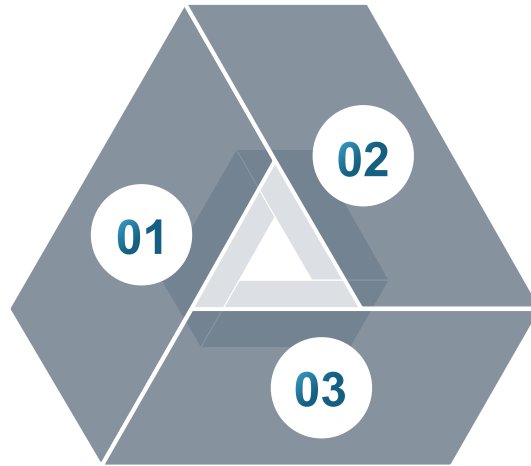
# Graph RAG



# GraphRAG—Limitation

## High Cost

**Generating and traversing community reports for retrieval and query stages leads to extremely high costs.**



## Inefficiency

**Both generating community reports and traversing the community require significant time, which greatly reduces system efficiency.**

## **Poor Scalability**

**Merging new data requires rebuilding communities, raising costs and limiting scalability.**

# LightRAG

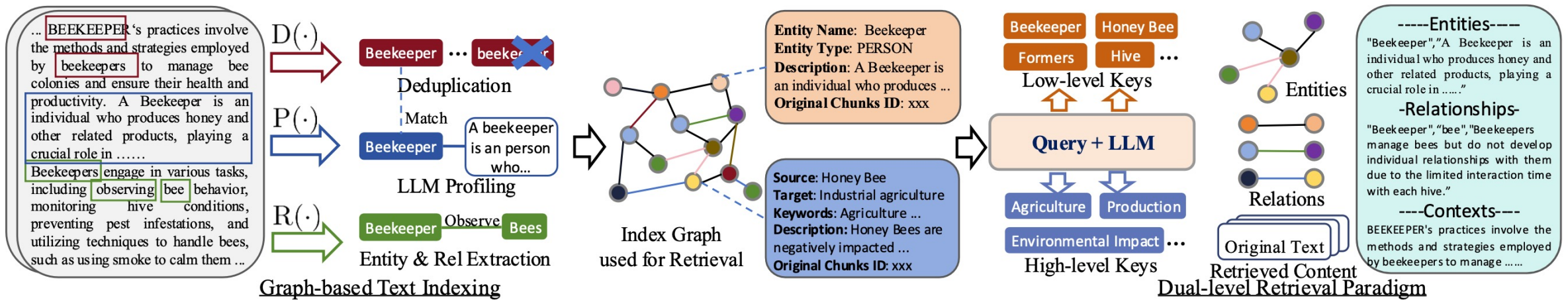


Figure 1: Overall architecture of the proposed LightRAG framework.

- 🔍 Comprehensive Information Retrieval with Complex Interdependencies
- ⚙️ Efficient Information Retrieval through a Dual-Level Retrieval Paradigm
- ⚡ Rapid Adaptability to Dynamic Data Changes

Table 4: Statistical information of the datasets.

<b>Statistics</b>	<b>Agriculture</b>	<b>CS</b>	<b>Legal</b>	<b>Mix</b>
Total Documents	12	10	94	61
Total Tokens	2,017,886	2,306,535	5,081,069	619,009

**Evaluation Datasets.** To conduct a comprehensive analysis of LightRAG, we selected four datasets from the UltraDomain benchmark (Qian et al., 2024). The UltraDomain data is sourced from 428 college textbooks and encompasses 18 distinct domains, including agriculture, social sciences, and humanities. From these, we chose the Agriculture, CS, Legal, and Mix datasets. Each dataset contains between 600,000 and 5,000,000 tokens, with detailed information provided in Table 4. Below is a specific introduction to the four domains utilized in our experiments:

- i) **Comprehensiveness:** How thoroughly does the answer address all aspects and details of the question?
- ii) **Diversity:** How varied and rich is the answer in offering different perspectives and insights related to the question?
- iii) **Empowerment:** How effectively does the answer enable the reader to understand the topic and make informed judgments?
- iv) **Overall:** This dimension assesses the cumulative performance across the three preceding criteria to identify the best overall answer.

Table 1: Win rates (%) of baselines v.s. LightRAG across four datasets and four evaluation dimensions.

	Agriculture		CS		Legal		Mix	
	NaiveRAG	<b>LightRAG</b>	NaiveRAG	<b>LightRAG</b>	NaiveRAG	<b>LightRAG</b>	NaiveRAG	<b>LightRAG</b>
Comprehensiveness	32.4%	<u>67.6%</u>	38.4%	<u>61.6%</u>	16.4%	<u>83.6%</u>	38.8%	<u>61.2%</u>
Diversity	23.6%	<u>76.4%</u>	38.0%	<u>62.0%</u>	13.6%	<u>86.4%</u>	32.4%	<u>67.6%</u>
Empowerment	32.4%	<u>67.6%</u>	38.8%	<u>61.2%</u>	16.4%	<u>83.6%</u>	42.8%	<u>57.2%</u>
Overall	32.4%	<u>67.6%</u>	38.8%	<u>61.2%</u>	15.2%	<u>84.8%</u>	40.0%	<u>60.0%</u>
	RQ-RAG	<b>LightRAG</b>	RQ-RAG	<b>LightRAG</b>	RQ-RAG	<b>LightRAG</b>	RQ-RAG	<b>LightRAG</b>
Comprehensiveness	31.6%	<u>68.4%</u>	38.8%	<u>61.2%</u>	15.2%	<u>84.8%</u>	39.2%	<u>60.8%</u>
Diversity	29.2%	<u>70.8%</u>	39.2%	<u>60.8%</u>	11.6%	<u>88.4%</u>	30.8%	<u>69.2%</u>
Empowerment	31.6%	<u>68.4%</u>	36.4%	<u>63.6%</u>	15.2%	<u>84.8%</u>	42.4%	<u>57.6%</u>
Overall	32.4%	<u>67.6%</u>	38.0%	<u>62.0%</u>	14.4%	<u>85.6%</u>	40.0%	<u>60.0%</u>
	HyDE	<b>LightRAG</b>	HyDE	<b>LightRAG</b>	HyDE	<b>LightRAG</b>	HyDE	<b>LightRAG</b>
Comprehensiveness	26.0%	<u>74.0%</u>	41.6%	<u>58.4%</u>	26.8%	<u>73.2%</u>	40.4%	<u>59.6%</u>
Diversity	24.0%	<u>76.0%</u>	38.8%	<u>61.2%</u>	20.0%	<u>80.0%</u>	32.4%	<u>67.6%</u>
Empowerment	25.2%	<u>74.8%</u>	40.8%	<u>59.2%</u>	26.0%	<u>74.0%</u>	46.0%	<u>54.0%</u>
Overall	24.8%	<u>75.2%</u>	41.6%	<u>58.4%</u>	26.4%	<u>73.6%</u>	42.4%	<u>57.6%</u>
	GraphRAG	<b>LightRAG</b>	GraphRAG	<b>LightRAG</b>	GraphRAG	<b>LightRAG</b>	GraphRAG	<b>LightRAG</b>
Comprehensiveness	45.6%	<u>54.4%</u>	48.4%	<u>51.6%</u>	48.4%	<u>51.6%</u>	<u>50.4%</u>	49.6%
Diversity	22.8%	<u>77.2%</u>	40.8%	<u>59.2%</u>	26.4%	<u>73.6%</u>	36.0%	<u>64.0%</u>
Empowerment	41.2%	<u>58.8%</u>	45.2%	<u>54.8%</u>	43.6%	<u>56.4%</u>	<u>50.8%</u>	49.2%
Overall	45.2%	<u>54.8%</u>	48.0%	<u>52.0%</u>	47.2%	<u>52.8%</u>	<u>50.4%</u>	49.6%

Table 2: Performance of ablated versions of LightRAG, using NaiveRAG as reference.

	Agriculture		CS		Legal		Mix	
	NaiveRAG	LightRAG	NaiveRAG	LightRAG	NaiveRAG	LightRAG	NaiveRAG	LightRAG
Comprehensiveness	32.4%	<u>67.6%</u>	38.4%	<u>61.6%</u>	16.4%	<u>83.6%</u>	38.8%	<u>61.2%</u>
Diversity	23.6%	<u>76.4%</u>	38.0%	<u>62.0%</u>	13.6%	<u>86.4%</u>	32.4%	<u>67.6%</u>
Empowerment	32.4%	<u>67.6%</u>	38.8%	<u>61.2%</u>	16.4%	<u>83.6%</u>	42.8%	<u>57.2%</u>
Overall	32.4%	<u>67.6%</u>	38.8%	<u>61.2%</u>	15.2%	<u>84.8%</u>	40.0%	<u>60.0%</u>
	NaiveRAG	<b>-High</b>	NaiveRAG	<b>-High</b>	NaiveRAG	<b>-High</b>	NaiveRAG	<b>-High</b>
Comprehensiveness	34.8%	<u>65.2%</u>	42.8%	<u>57.2%</u>	23.6%	<u>76.4%</u>	40.4%	<u>59.6%</u>
Diversity	27.2%	<u>72.8%</u>	36.8%	<u>63.2%</u>	16.8%	<u>83.2%</u>	36.0%	<u>64.0%</u>
Empowerment	36.0%	<u>64.0%</u>	42.4%	<u>57.6%</u>	22.8%	<u>77.2%</u>	47.6%	<u>52.4%</u>
Overall	35.2%	<u>64.8%</u>	44.0%	<u>56.0%</u>	22.0%	<u>78.0%</u>	42.4%	<u>57.6%</u>
	NaiveRAG	<b>-Low</b>	NaiveRAG	<b>-Low</b>	NaiveRAG	<b>-Low</b>	NaiveRAG	<b>-Low</b>
Comprehensiveness	36.0%	<u>64.0%</u>	43.2%	<u>56.8%</u>	19.2%	<u>80.8%</u>	36.0%	<u>64.0%</u>
Diversity	28.0%	<u>72.0%</u>	39.6%	<u>60.4%</u>	13.6%	<u>86.4%</u>	33.2%	<u>66.8%</u>
Empowerment	34.8%	<u>65.2%</u>	42.8%	<u>57.2%</u>	16.4%	<u>83.6%</u>	35.2%	<u>64.8%</u>
Overall	34.8%	<u>65.2%</u>	43.6%	<u>56.4%</u>	18.8%	<u>81.2%</u>	35.2%	<u>64.8%</u>
	NaiveRAG	<b>-Origin</b>	NaiveRAG	<b>-Origin</b>	NaiveRAG	<b>-Origin</b>	NaiveRAG	<b>-Origin</b>
Comprehensiveness	24.8%	<u>75.2%</u>	39.2%	<u>60.8%</u>	16.4%	<u>83.6%</u>	44.4%	<u>55.6%</u>
Diversity	26.4%	<u>73.6%</u>	44.8%	<u>55.2%</u>	14.4%	<u>85.6%</u>	25.6%	<u>74.4%</u>
Empowerment	32.0%	<u>68.0%</u>	43.2%	<u>56.8%</u>	17.2%	<u>82.8%</u>	45.2%	<u>54.8%</u>
Overall	25.6%	<u>74.4%</u>	39.2%	<u>60.8%</u>	15.6%	<u>84.4%</u>	44.4%	<u>55.6%</u>

# On-Device RAG



## Data Privacy



## Resource Limited



# Multi-Modal RAG



2025/8/20

**IBM Technology**  
 @ibmtechnology 1.12M subscribers · 1.2K videos  
 Whether it's AI, automation, cybersecurity, data science, DevOps, quantum computing or...more  
 ibm.com/watsonx and 2 more links

Subscribed

Home Videos Shorts Podcasts Playlists Posts

Latest Popular Oldest

- DeepSeek-V3-0324, Gemini 2.5 and Canvas, Extropic's thermodynamic chip, and OpenAI image generation 41:40
- Building Text-to-SQL Agent for Smarter Database Queries 26:18
- Are LLMs losing their minds? Exploring Generative AI 6:43
- LLMs Simplifying Application Modernization 2:16
- DeepSeek-V3-0324, Gemini Canvas and GPT-4o image generation 3:2K views · 1 day ago
- Build a Text-to-SQL Agent for Smarter Database Queries 1.9K views · 2 days ago
- Can LLMs Learn Without Losing Their Minds? Exploring Generative AI 5.8K views · 3 days ago
- Large Language Models Simplifying App Modernization 7.4K views · 4 days ago
- What is a Vector Database 9:49
- Mixture of Experts | Ep. 47 39:13
- Disaster Recovery vs Operational Resilience 7:24
- Data Visualization Made Simple: Do's and Don'ts 5:38
- What is a Vector Database? Powering Semantic Search & AI Applications 93K views · 5 days ago
- NVIDIA GTC, Baidu reasoning models and Gemini AI image generation 4.7K views · 8 days ago
- Disaster Recovery vs Operational Resilience: Protecting Your Data 3.9K views · 10 days ago
- Data Visualization Made Simple: Do's & Don'ts for Clear Insights 8.9K views · 11 days ago
- RAG vs. CAG: Solving Knowledge Gaps in AI Models 16:00
- Manus, vibe coding, scaling laws and Perplexity's AI phone 19K views · 2 weeks ago
- The Truth About Ground Truth 10:05
- Building Trustworthy AI: Avoid Model Drift & Unsafe Outputs 4K views · 2 months ago

**Sequoia Capital**  
 @sequoiacapital 55.9K subscribers · 75 videos  
 We help the daring build legendary companies – from idea to IPO and beyond...more  
 sequoiacap.com

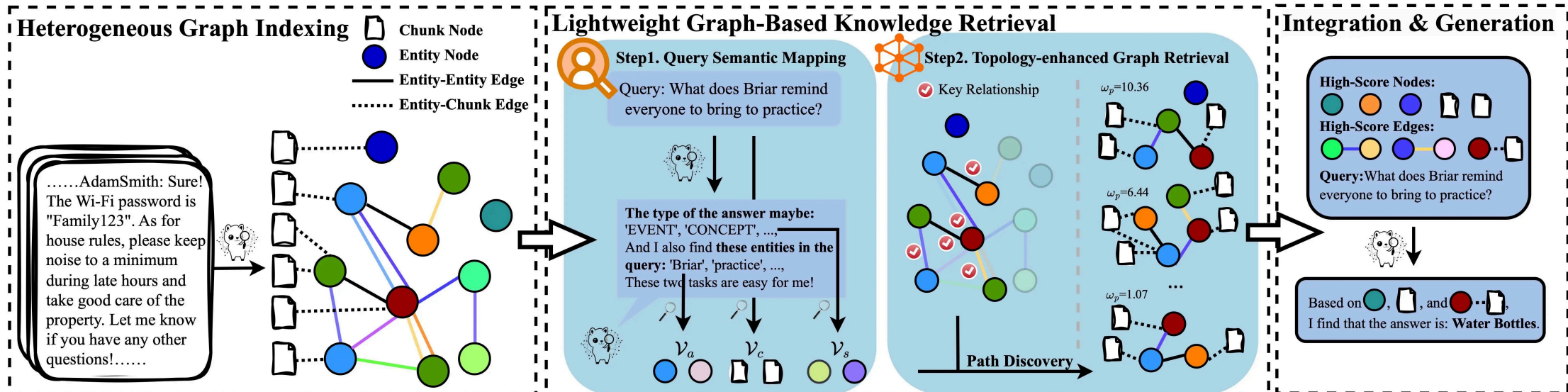
Subscribe

Home Videos Shorts Podcasts Playlists

Latest Popular Oldest

- TRAINING DATA 43:05
- TRAINING DATA 61:17
- TRAINING DATA 54:10
- TRAINING DATA 1:04:32
- From Software Engineers to AI Word Artists: Filip Kozera of Wordware 2.5K views · 4 days ago
- Josh Woodward: Google Labs is Rapidly Building AI Products from 0-to-1 10K views · 11 days ago
- How AI Breakout Harvey is Transforming Legal Services, with CEO Winston Weinberg 7.6K views · 2 weeks ago
- The AI Product Going Viral With Doctors: OpenEvidence, with CEO Daniel Nadler 4.4K views · 3 weeks ago
- TRAINING DATA 32:46
- TRAINING DATA 1:00:09
- TRAINING DATA 44:27
- TRAINING DATA 54:57
- OpenAI's Deep Research Team on Why Reinforcement Learning is the Future for All... 12K views · 1 month ago
- AI Security and the New World Order ft. Palo Alto Networks's Hitesh Arora 10K views · 1 month ago
- Vector Databases and the Data Structure of AI ft. MongoDB's Sahir Azam 9K views · 1 month ago
- Using AI to Empower Creators ft Roblox Studio Head Stef Corazza 831 views · 1 month ago
- TRAINING DATA 52:30
- CRUCIBLE MOMENTS 49:33
- TRAINING DATA 52:57
- TRAINING DATA 1:00:06
- From AlphaGo to AGI ft ReflectionAI Founder Ioannis Antonoglou 2.4K views · 2 months ago
- Nubank ft. David Vélez: An Outsider Upends the Brazilian Banking System 8K views · 2 months ago
- Turning Graph AI into ROI ft Kum'o's Hema Raghavan 1.6K views · 2 months ago
- Turning Academic Open Source into Startup Success ft Databricks Founder Ion Stoica 4.3K views · 2 months ago

# MiniRAG



# What is MiniRAG?



MiniRAG is a **free, lightweight, yet powerful** RAG framework.

**Free:** No need for API-based LLMs; local SLMs can also handle RAG tasks.

**Lightweight:** Requires only 25% of the storage space compared to leading baselines.

**Strong:** MiniRAG+3B model > GraphRAG+gpt-4o-mini

# Observation of SLMs on RAG tasks

## Limitation of Indexing

### Original Text:

LiHua: Hey Adam! Hope you're doing well. Just wanted to ask for the Wi-Fi password when you get a chance. Also, are there any house rules I should know about? Thanks!

AdamSmith: Sure! The Wi-Fi password is "Family123". As for house rules, please keep noise to a minimum during late hours and take good care of the property. Let me know if you have any other questions!

LiHua: Got it, thanks! I'll make sure to keep it down at night. Just one more thing, is it okay to have friends over occasionally?

AdamSmith: You're welcome! Yes, having friends over occasionally is fine, just try to keep the gatherings reasonable. Enjoy!

LiHua: Awesome, appreciate it! Looking forward to settling in.

### gpt-4o-mini response:

Entity: <HOUSE RULES>

Description: <House Rules refer to the guidelines provided by Adam Smith regarding noise and care of the property, aimed at maintaining a peaceful living environment.>

### Phi-3.5-mini response:

Entity: <HOUSE RULES>

Description: <House rules are guidelines provided by AdamSmith for LiHua to follow while living in the house.>

## Limitation of Answering

### Input Prompt:

You are an assistant designed to answer questions. Please respond to the user's query based on the information provided by the user. If there is insufficient information, please reply with "I don't have enough information."

Here is the information:

**Entities:** <HOUSE RULES> ... and 49 other entities.

**Relations:** <LiHua, AdamSmith> ... and 49 other relationships.

**Chunk:** 10 chunks

**Query:** What are the House Rules mentioned by Adam Smith?

### gpt-4o-mini response:

Based on <HOUSE RULES> and <LiHua, AdamSmith>, I know that the House Rules mentioned by Adam Smith is about...

### Phi-3.5-mini response:

I'm sorry, I don't have enough information.

**Left:** SLMs generate notably lower-quality descriptions than LLMs. **Right:** When processing identical inputs, SLMs struggle to locate relevant information in large contexts, while LLMs perform this task effectively.

**Observation 1:** while SLMs struggle with sophisticated semantic understanding, they excel at pattern matching and localized text processing

**Observation 2:** explicit structural information can effectively compensate for limited semantic capabilities

**Observation 3:** decomposing complex RAG operations into simpler, well-defined steps can maintain robustness without requiring advanced reasoning capabilities.

# Motivation of MiniRAG

**Observation 1:** while SLMs struggle with sophisticated semantic understanding, they excel at pattern matching and localized text processing

**Observation 2:** explicit structural information can effectively compensate for limited semantic capabilities

**Observation 3:** decomposing complex RAG operations into simpler, well-defined steps can maintain robustness without requiring advanced reasoning capabilities.

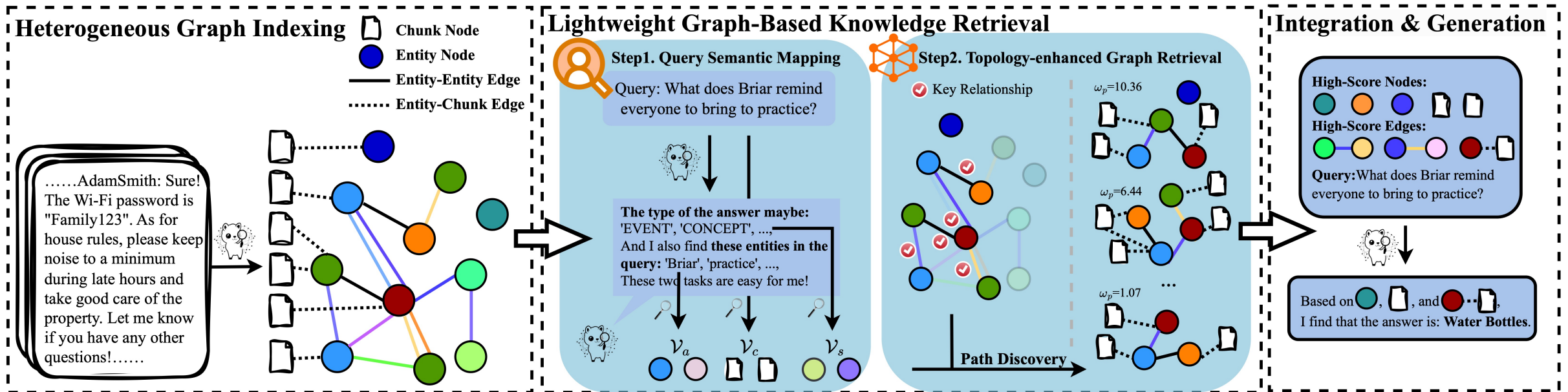


To alleviate the **high demands on model capabilities** imposed by the RAG framework, We should use

- (1) a new indexing mechanism** to reduce reliance on complex semantic understanding
- (2) a new retrieval approach** for efficient knowledge discovery.

# The MiniRAG framework

- (1) a new indexing mechanism to reduce reliance on complex semantic understanding
- (2) a new retrieval approach for efficient knowledge discovery.



The MiniRAG framework employs a streamlined workflow built on the key components: (1) **heterogeneous graph indexing** and (2) **lightweight graph-based knowledge retrieval**.

# Details of MiniRAG: Heterogeneous Graph Indexing



## Target:

- (1) The indexing mechanism should **extract the key relationships and contextual connections within the data.**
- (2) The indexing approach should **condense the retrieved content to its most query-relevant elements.**



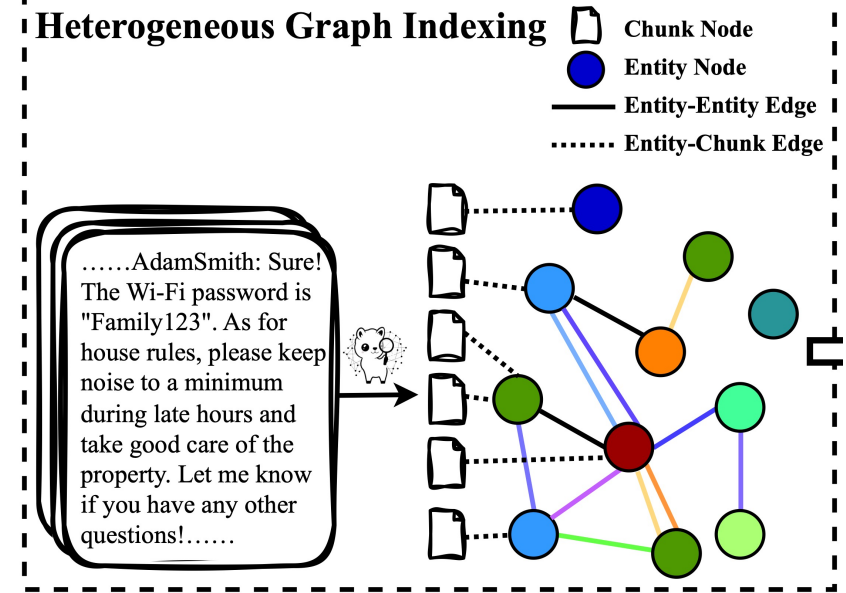
## Key components:

- Text Chunk Node  $\mathcal{V}_c$
- Entity Node  $\mathcal{V}_e$
- Entity-Entity Connection  $\mathcal{E}_\alpha$
- Entity-Chunk Connection  $\mathcal{E}_\beta$



## Advantages:

1. It systematically incorporates both text chunks and named entities, **creating a rich semantic network**
2. It enables data chunks to directly participate in retrieval and effectively **mitigates information distortion**
3. It **avoids challenging tasks** (e.g., sophisticated semantic understanding) that could cause SLMs to fail.



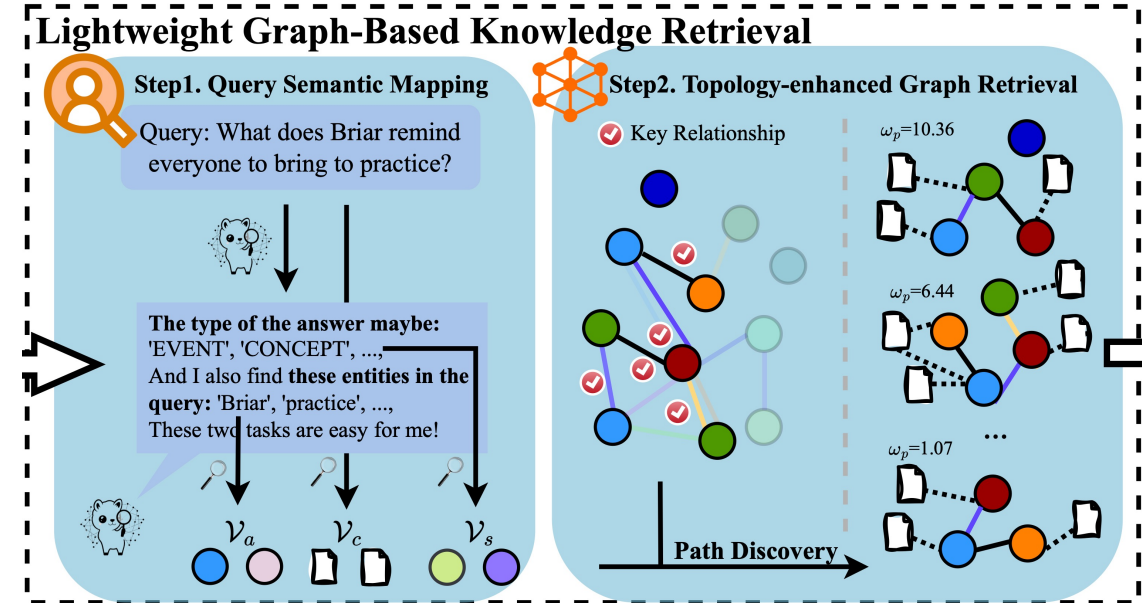
# Details of MiniRAG: lightweight graph-based knowledge retrieval.

## Target:

- (1) **Reduce the complexity of input content.** For generation, ensuring that semantic information is clear and concise;
- (2) **Shorten the length of input content.** For SLMs, facilitating improved comprehension and retrieval accuracy.

## Step 1: Query Semantic Mapping

- (1) Extract relevant entities  $\mathcal{V}_q$  and predict potential types of query
- (2) Query-driven reasoning path discovery:
  - Initial Entity Identification  $\hat{\mathcal{V}}_s$
  - Answer-Aware Entity Selection  $\hat{\mathcal{V}}_a$
  - Context-Rich Path Formation  $\hat{\mathcal{V}}_c$



# Details of MiniRAG: lightweight graph-based knowledge retrieval.

## Target:

- (1) **Reduce the complexity of input content.** For generation, ensuring that semantic information is clear and concise; ✓
- (2) **Shorten the length of input content.** For SLMs, facilitating improved comprehension and retrieval accuracy.

## Step 2: Topology-Enhanced Graph Retrieval

### (1) Key Relationship Identification

$$\omega_e(e) = \sum_{\hat{v}_s \in \hat{V}_s} \text{count}(\hat{v}_s, \hat{G}_{e,k}) + \sum_{\hat{v}_a \in \hat{V}_a} \text{count}(\hat{v}_a, \hat{G}_{e,k}),$$

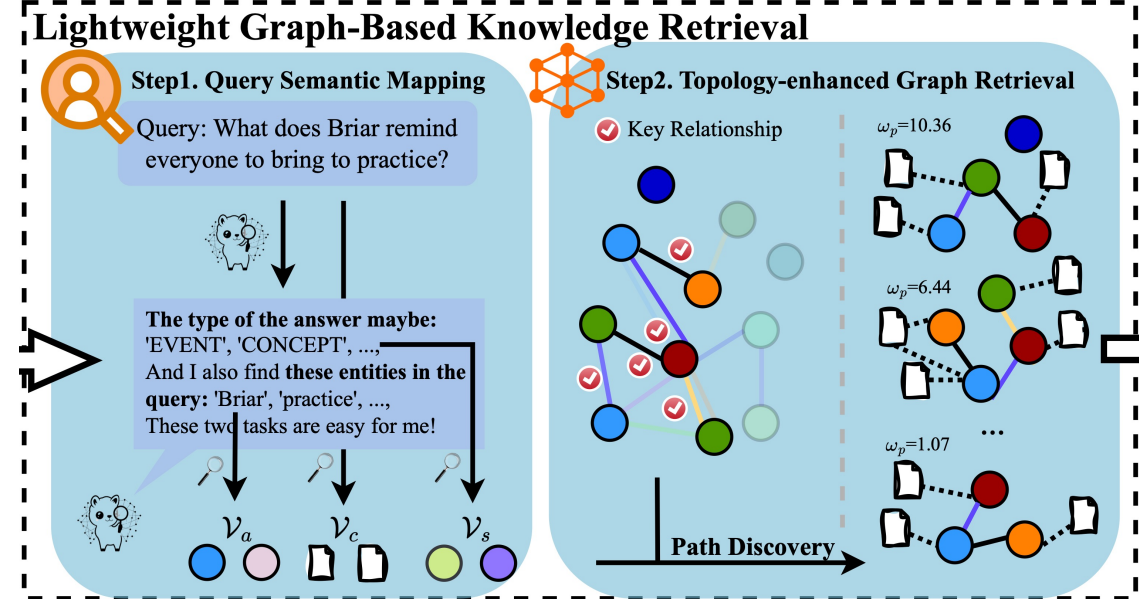
### (2) Query-guided path discovery

$$\omega_p(p | v_q) = \omega_v(\hat{v}_s | v_q) \cdot (1 + \sum_{v \in (p \wedge \hat{V}_a)} \text{count}(v, p) + \sum_{e \in (p \wedge \hat{E}_\alpha)} \omega_e(e)).$$

### (3) Retrieval of Query-Relevant Text Chunks

① Candidate Filtering  $\hat{v}_c \wedge v_c^q$    ② Similarity Computation   ③ Ranking & Selection

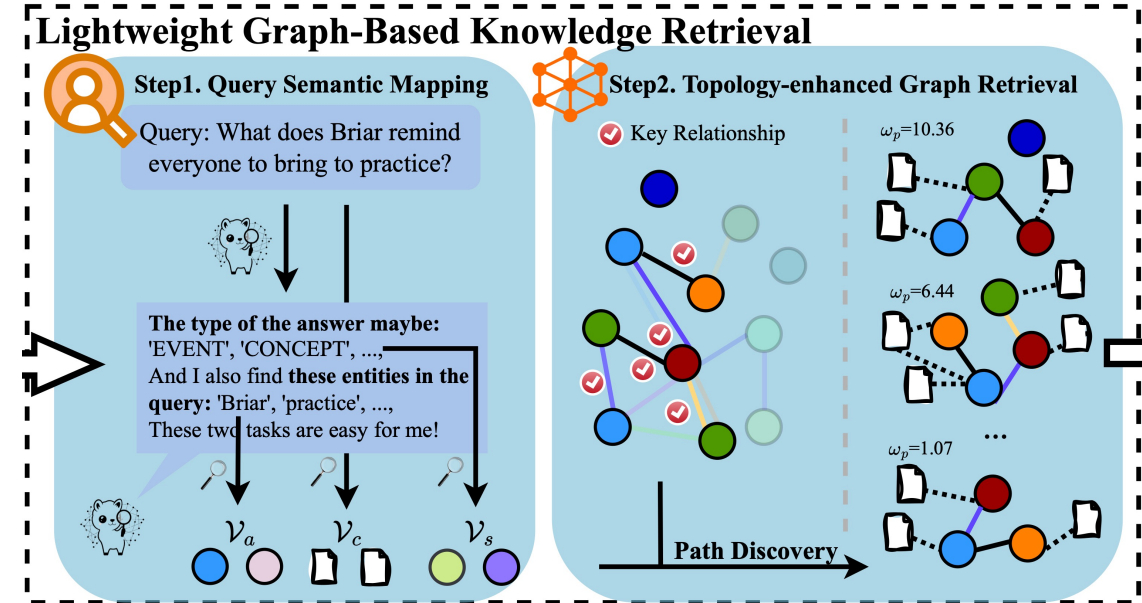
### (4) Integration for Augmented Generation



# Details of MiniRAG: lightweight graph-based knowledge retrieval.

## Target:

- (1) **Reduce the complexity of input content.** For generation, ensuring that semantic information is clear and concise; ✓
- (2) **Shorten the length of input content.** For SLMs, facilitating improved comprehension and retrieval accuracy. ✓



## Advantages:

1. It effectively **mitigates the issue of significant noise** introduced into the retrieval process due to the limited semantic understanding capabilities of SLMs.
2. It **avoids challenging tasks** (e.g., predicting the high-level information of the query) that could cause SLMs to fail.

# New Datasets: LiHua-World



**[Design for on-device RAG]** The LiHua-World dataset is specifically designed for local RAG scenarios, containing **one year of chat records**.

**[Diverse types of questions]** The dataset includes single-hop, multi-hop, and summarization questions, each accompanied by manually annotated answers and supporting documents.

**[Wide variety of events]** The content covers various aspects of daily life, such as social interactions, fitness training, recreational activities, and personal affairs.

# Experiment

**RQ1: Comparative Performance.** How does MiniRAG perform against state-of-the-art alternatives in terms of retrieval accuracy and efficiency?

**RQ2: Component Analysis.** What is the contribution of key components to MiniRAG’s overall effectiveness?

**RQ3: Case Studies.** How effectively does MiniRAG handle complex, multi-step reasoning tasks with small language models, as demonstrated through practical case studies?

Table 1: Performance evaluation using accuracy (acc) and error (err) rates, measured as percentages (%). Higher accuracy and lower error rates indicate better RAG performance. Results compare various baseline methods against our MiniRAG across multiple datasets. Bold values indicate best performance, while “/” denotes cases where methods failed to generate effective responses.

LiHuaWorld	NaiveRAG		GraphRAG		LightRAG		MiniRAG	
	acc↑	err↓	acc↑	err↓	acc↑	err↓	acc↑	err↓
Phi-3.5-mini-instruct	41.22%	23.20%	/	/	39.81%	25.39%	<b>53.29%</b>	23.35%
GLM-Edge-1.5B-Chat	42.79%	24.76%	/	/	35.74%	25.86%	<b>52.51%</b>	25.71%
Qwen2.5-3B-Instruct	43.73%	24.14%	/	/	39.18%	28.68%	<b>48.75%</b>	26.02%
MiniCPM3-4B	43.42%	17.08%	/	/	35.42%	21.94%	<b>51.25%</b>	21.79%
gpt-4o-mini	46.55%	19.12%	35.27%	37.77%	<b>56.90%</b>	20.85%	54.08%	19.44%

MultiHop-RAG	NaiveRAG		GraphRAG		LightRAG		MiniRAG	
	acc↑	err↓	acc↑	err↓	acc↑	err↓	acc↑	err↓
Phi-3.5-mini-instruct	42.72%	31.34%	/	/	27.03%	11.78%	<b>49.96%</b>	28.44%
GLM-Edge-1.5B-Chat	44.44%	24.26%	/	/	/	/	<b>51.41%</b>	23.44%
Qwen2.5-3B-Instruct	39.48%	31.69%	/	/	21.91%	13.73%	<b>48.55%</b>	33.10%
MiniCPM3-4B	39.24%	31.42%	/	/	19.48%	10.41%	<b>47.77%</b>	26.88%
gpt-4o-mini	53.60%	27.19%	60.92%	16.86%	64.91%	19.37%	<b>68.43%</b>	19.41%

# Experiment

**RQ1: Comparative Performance.** How does MiniRAG perform against state-of-the-art alternatives in terms of retrieval accuracy and efficiency?

Table 1: Performance evaluation using accuracy (acc) and error (err) rates, measured as percentages (%). Higher accuracy and lower error rates indicate better RAG performance. Results compare various baseline methods against our MiniRAG across multiple datasets. Bold values indicate best performance, while “/” denotes cases where methods failed to generate effective responses.

LiHuaWorld	NaiveRAG		GraphRAG		LightRAG		MiniRAG	
	acc↑	err↓	acc↑	err↓	acc↑	err↓	acc↑	err↓
Phi-3.5-mini-instruct	41.22%	23.20%	/	/	39.81%	25.39%	<b>53.29%</b>	23.35%
GLM-Edge-1.5B-Chat	42.79%	24.76%	/	/	35.74%	25.86%	<b>52.51%</b>	25.71%
Qwen2.5-3B-Instruct	43.73%	24.14%	/	/	39.18%	28.68%	<b>48.75%</b>	26.02%
MiniCPM3-4B	43.42%	17.08%	/	/	35.42%	21.94%	<b>51.25%</b>	21.79%
gpt-4o-mini	46.55%	19.12%	35.27%	37.77%	<b>56.90%</b>	20.85%	54.08%	19.44%

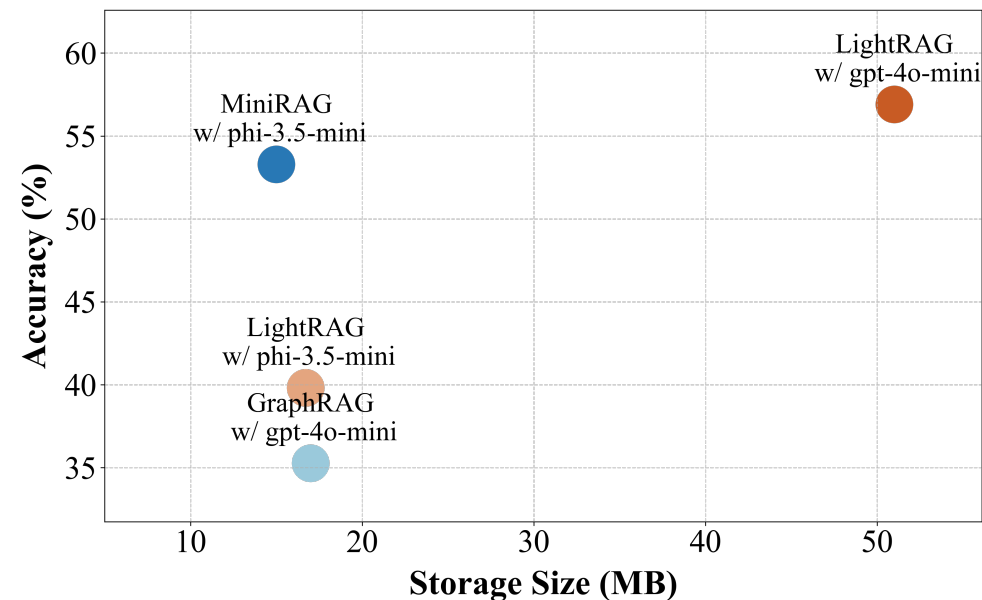
  

MultiHop-RAG	NaiveRAG		GraphRAG		LightRAG		MiniRAG	
	acc↑	err↓	acc↑	err↓	acc↑	err↓	acc↑	err↓
Phi-3.5-mini-instruct	42.72%	31.34%	/	/	27.03%	11.78%	<b>49.96%</b>	28.44%
GLM-Edge-1.5B-Chat	44.44%	24.26%	/	/	/	/	<b>51.41%</b>	23.44%
Qwen2.5-3B-Instruct	39.48%	31.69%	/	/	21.91%	13.73%	<b>48.55%</b>	33.10%
MiniCPM3-4B	39.24%	31.42%	/	/	19.48%	10.41%	<b>47.77%</b>	26.88%
gpt-4o-mini	53.60%	27.19%	60.92%	16.86%	64.91%	19.37%	<b>68.43%</b>	19.41%

## Observation:

1. Current RAG systems **face critical challenges** when operating with SLMs.
2. MiniRAG **maintain strong performance** even with SLMs.
3. MiniRAG demonstrates exceptional **storage efficiency** while preserving high accuracy levels.

Accuracy & Storage Comparison of different RAG methods



# Experiment

**RQ2: Component Analysis.** What is the contribution of key components to MiniRAG’s overall effectiveness?

Table 2: Ablation study results comparing accuracy ( $acc$ ,  $\uparrow$ ) and error rate ( $err$ ,  $\downarrow$ ) (%) across architectural variants: baseline MiniRAG versus variants with (i) semantic-driven indexing replacement ( $-\mathcal{I}$ ), (ii) edge information removal ( $-\mathcal{R}_{edge}$ ), and (iii) chunk nodes removal ( $-\mathcal{R}_{chunk}$ ). Results validate SLM limitations and the effectiveness of query-guided reasoning path components.

LiHuaWorld	MiniRAG		$-\mathcal{I}$		$-\mathcal{R}_{chunk}$		$-\mathcal{R}_{edge}$	
	$acc\uparrow$	$err\downarrow$	$acc\uparrow$	$err\downarrow$	$acc\uparrow$	$err\downarrow$	$acc\uparrow$	$err\downarrow$
Phi-3.5-mini-instruct	53.29%	23.35%	26.02%	19.12%	48.90%	17.40%	50.47%	15.36%
GLM-Edge-1.5B-Chat	52.51%	25.71%	25.08%	31.50%	46.24%	16.77%	47.81%	20.53%
Qwen2.5-3B-Instruct	48.75%	26.02%	24.14%	15.67%	40.91%	16.14%	48.43%	18.65%
MiniCPM3-4B	51.25%	21.79%	26.18%	15.52%	46.39%	15.83%	48.59%	19.44%

## Observation:

**1. SLM Limitations.** Replacing indexing with text semantic-driven techniques ( $-\mathcal{I}$ ) causes significant performance drop, **validating SLMs’ constraints in semantic understanding**, impacting knowledge graph generation and text descriptions.

**2. Effectiveness of Query-guided Reasoning Path Discovery.** Removing edge info ( $-\mathcal{R}_{edge}$ ) or chunk nodes ( $-\mathcal{R}_{chunk}$ ) harms performance. These components enable query matching and **compensate for SLMs’ limitations during indexing**.

# 🚀 RAG-Anything: All-in-One RAG System

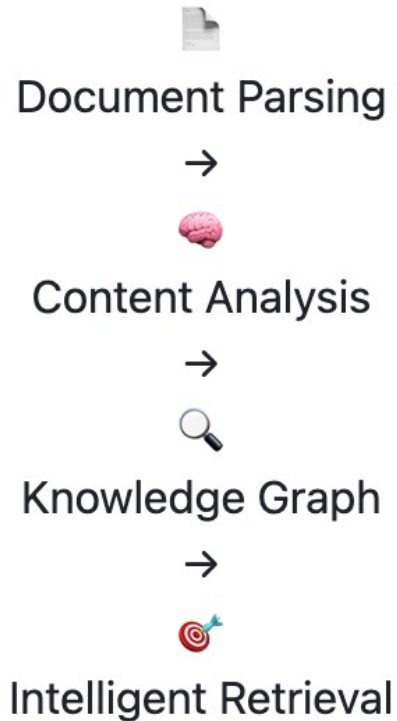
## 🎯 Key Features

- 🔄 **End-to-End Multimodal Pipeline** - Complete workflow from document ingestion and parsing to intelligent multimodal query answering
- 📄 **Universal Document Support** - Seamless processing of PDFs, Office documents, images, and diverse file formats
- 🧠 **Specialized Content Analysis** - Dedicated processors for images, tables, mathematical equations, and heterogeneous content types
- 🔗 **Multimodal Knowledge Graph** - Automatic entity extraction and cross-modal relationship discovery for enhanced understanding
- ⚡ **Adaptive Processing Modes** - Flexible MinerU-based parsing or direct multimodal content injection workflows
- 📋 **Direct Content List Insertion** - Bypass document parsing by directly inserting pre-parsed content lists from external sources
- 🎯 **Hybrid Intelligent Retrieval** - Advanced search capabilities spanning textual and multimodal content with contextual understanding



## Core Algorithm




RAG-Anything implements an effective **multi-stage multimodal pipeline** that fundamentally extends traditional RAG architectures to seamlessly handle diverse content modalities through intelligent orchestration and cross-modal understanding.



# 1. Document Parsing Stage

The system provides high-fidelity document extraction through adaptive content decomposition. It intelligently segments heterogeneous elements while preserving contextual relationships. Universal format compatibility is achieved via specialized optimized parsers.

## Key Components:

-  **MinerU Integration:** Leverages [MinerU](#) for high-fidelity document structure extraction and semantic preservation across complex layouts.
-  **Adaptive Content Decomposition:** Automatically segments documents into coherent text blocks, visual elements, structured tables, mathematical equations, and specialized content types while preserving contextual relationships.
-  **Universal Format Support:** Provides comprehensive handling of PDFs, Office documents (DOC/DOCX/PPT/PPTX/XLS/XLSX), images, and emerging formats through specialized parsers with format-specific optimization.

## 2. Multi-Modal Content Understanding & Processing

The system automatically categorizes and routes content through optimized channels. It uses concurrent pipelines for parallel text and multimodal processing. Document hierarchy and relationships are preserved during transformation.




### Key Components:

- 🎯 **Autonomous Content Categorization and Routing:** Automatically identify, categorize, and route different content types through optimized execution channels.
- ⚡ **Concurrent Multi-Pipeline Architecture:** Implements concurrent execution of textual and multimodal content through dedicated processing pipelines. This approach maximizes throughput efficiency while preserving content integrity.
- 📄 **Document Hierarchy Extraction:** Extracts and preserves original document hierarchy and inter-element relationships during content transformation.

### 3. Multimodal Analysis Engine

The system deploys modality-aware processing units for heterogeneous data modalities:

#### Specialized Analyzers:

-  **Visual Content Analyzer:**
  - Integrate vision model for image analysis.
  - Generates context-aware descriptive captions based on visual semantics.
  - Extracts spatial relationships and hierarchical structures between visual elements.
-  **Structured Data Interpreter:**
  - Performs systematic interpretation of tabular and structured data formats.
  - Implements statistical pattern recognition algorithms for data trend analysis.
  - Identifies semantic relationships and dependencies across multiple tabular datasets.
-  **Mathematical Expression Parser:**
  - Parses complex mathematical expressions and formulas with high accuracy.
  - Provides native LaTeX format support for seamless integration with academic workflows.
  - Establishes conceptual mappings between mathematical equations and domain-specific knowledge bases.

## 4. Multimodal Knowledge Graph Index

The multi-modal knowledge graph construction module transforms document content into structured semantic representations. It extracts multimodal entities, establishes cross-modal relationships, and preserves hierarchical organization. The system applies weighted relevance scoring for optimized knowledge retrieval.




### Core Functions:

- 🔍 **Multi-Modal Entity Extraction:** Transforms significant multimodal elements into structured knowledge graph entities. The process includes semantic annotations and metadata preservation.
- 🔗 **Cross-Modal Relationship Mapping:** Establishes semantic connections and dependencies between textual entities and multimodal components. This is achieved through automated relationship inference algorithms.
- 🏗️ **Hierarchical Structure Preservation:** Maintains original document organization through "belongs\_to" relationship chains. These chains preserve logical content hierarchy and sectional dependencies.
- ⚖️ **Weighted Relationship Scoring:** Assigns quantitative relevance scores to relationship types. Scoring is based on semantic proximity and contextual significance within the document structure.

## 5. Modality-Aware Retrieval

The hybrid retrieval system combines vector similarity search with graph traversal algorithms for comprehensive content retrieval. It implements modality-aware ranking mechanisms and maintains relational coherence between retrieved elements to ensure contextually integrated information delivery.

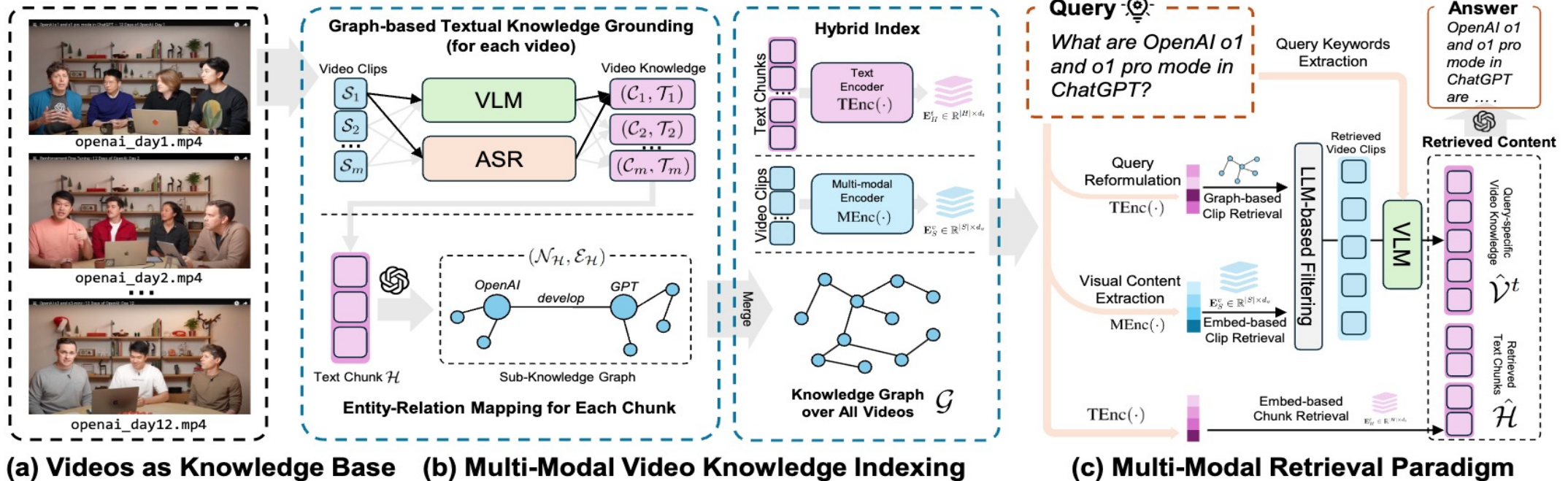
### Retrieval Mechanisms:

-  **Vector-Graph Fusion:** Integrates vector similarity search with graph traversal algorithms. This approach leverages both semantic embeddings and structural relationships for comprehensive content retrieval.
-  **Modality-Aware Ranking:** Implements adaptive scoring mechanisms that weight retrieval results based on content type relevance. The system adjusts rankings according to query-specific modality preferences.
-  **Relational Coherence Maintenance:** Maintains semantic and structural relationships between retrieved elements. This ensures coherent information delivery and contextual integrity.

# VideoRAG



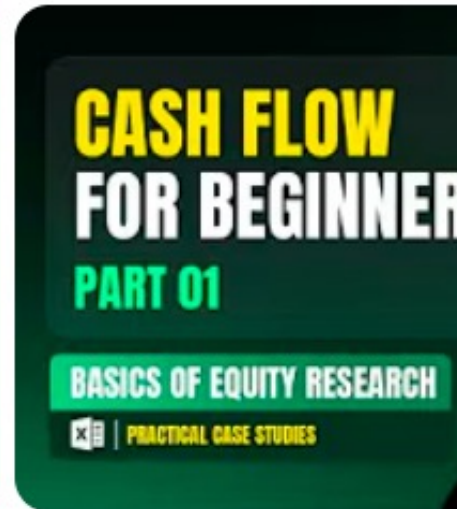
Video Type	#video list	#video	#query	#avg. queries per list	#overall duration
Lecture	12	135	376	31.3	~ 64.3 hours
Documentary	5	12	114	22.8	~ 28.5 hours
Entertainment	5	17	112	22.4	~ 41.9 hours
All	22	164	602	27.4	~ 134.6 hours





Jerry from Uganda: An Open Learner's Story

Lotfullah from Afghani An Open Learner's Stor



OUR OWN PDF Langchain

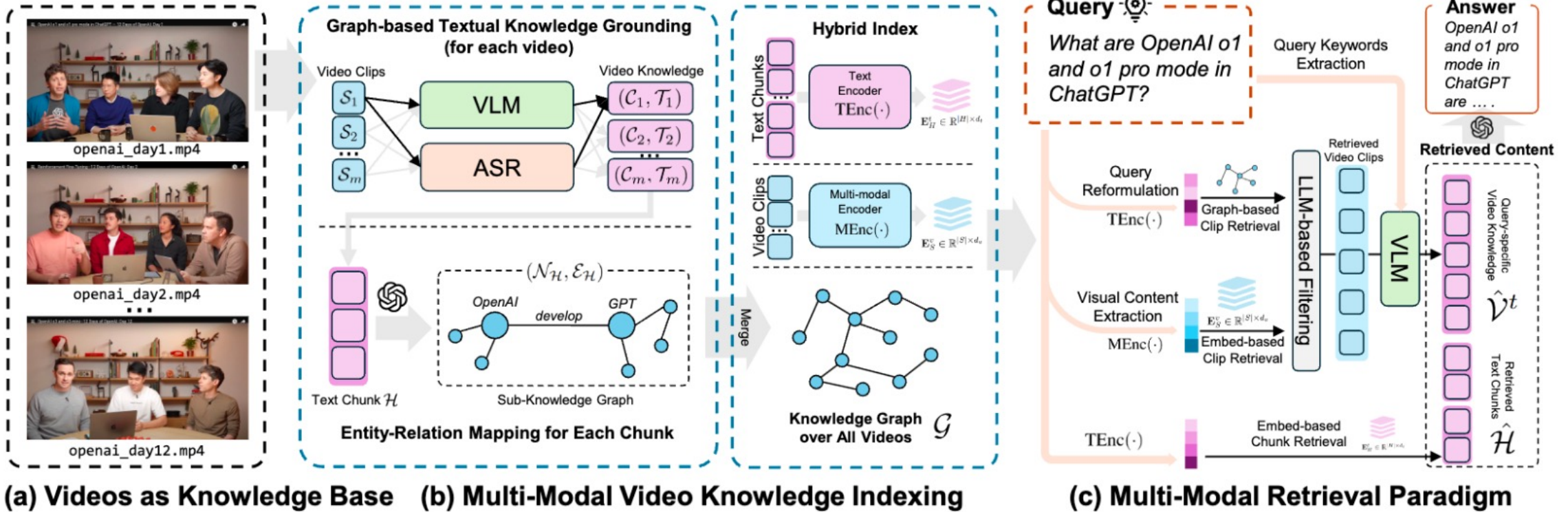
Langchain - 2 | (GUI) | Chourse...

Cash flow for Beg Equity Research I

Videos on the Internet contain a wealth of knowledge.

Can AI watch all the videos and learn from them?

# We present VideoRAG



An unlimited AI framework for watching long-context videos on a single RTX 3090.

# What is Retrieval-Augmented Generation (RAG)?

Knowledge Data Base (e.g., Documents, Videos)  $\mathcal{D}$



Indexing Process  $\hat{\mathcal{D}} = \varphi(\mathcal{D})$



Retrieval Knowledge for Input Query  $\psi(q, \hat{\mathcal{D}})$



Answer Generation  $\text{LLM}(q, \psi(q, \hat{\mathcal{D}}))$

# VideoRAG Knowledge Base

- We are considering leveraging **video list** as the knowledge base for the VideoRAG framework, such as OpenAI's 12 Days Show.
- There are no limitations on the duration of each video or the number of videos.

The screenshot shows a video player interface for a playlist titled "12 Days of OpenAI". The current video is "Work with Apps—12 Days of OpenAI: Day 11" with a duration of 19:04. Below it, a list of other videos in the playlist is shown, including "1-800-CHAT-GPT—12 Days of OpenAI: Day 10", "Dev Day Holiday Edition—12 Days of OpenAI: Day 9", "Search—12 Days of OpenAI: Day 8", "Projects—12 Days of OpenAI: Day 7", and "Santa Mode & Video in Advanced Voice—12 Days of OpenAI: Day 6".

Day	Video Title	Duration
11	Work with Apps—12 Days of OpenAI: Day 11	19:04
10	1-800-CHAT-GPT—12 Days of OpenAI: Day 10	11:17
9	Dev Day Holiday Edition—12 Days of OpenAI: Day 9	22:15
8	Search—12 Days of OpenAI: Day 8	13:35
7	Projects—12 Days of OpenAI: Day 7	18:18
6	Santa Mode & Video in Advanced Voice—12 Days of OpenAI: Day 6	

# VideoRAG

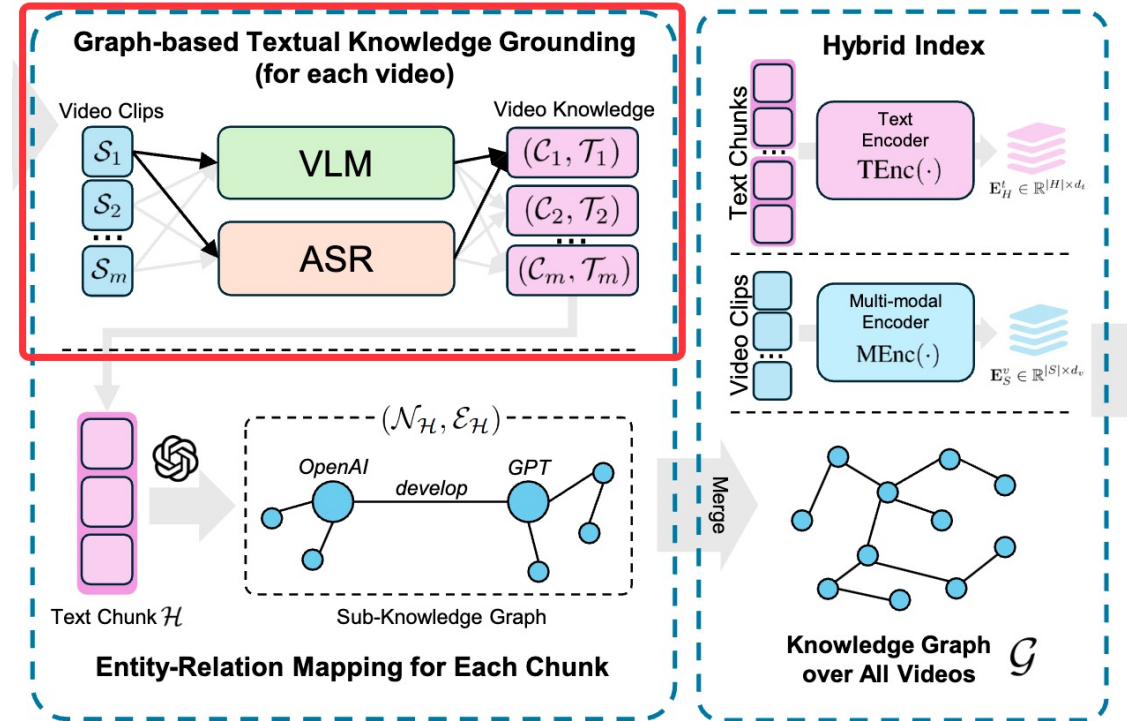
## Multi-Modal Video Knowledge Indexing

- For each video, we split it into sub-clips and use VLMs and ASR models to extract both **visual and audio information**.

$$\mathcal{C}_j = \text{VLM}(\mathcal{T}_j, \{\mathbf{F}_1, \dots, \mathbf{F}_k\} \mid \mathbf{F} \in \mathcal{S}_j)$$

$$\mathcal{T}_j = \text{ASR}(\mathcal{S}_j)$$

$$\mathcal{V}^t = \{(\mathcal{C}_l, \mathcal{T}_l) \mid l \in [1, m]\}$$

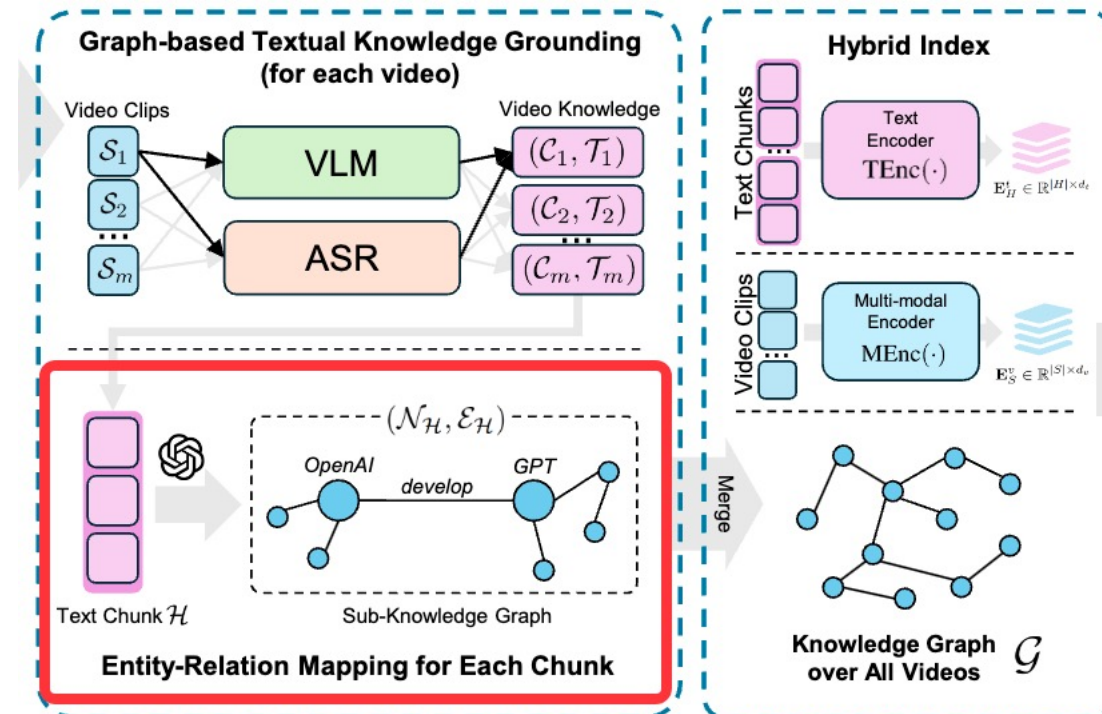


# VideoRAG

## Multi-Modal Video Knowledge Indexing

- Next, we construct a **knowledge graph** from the extracted textual knowledge of the sub clips to link information across multiple videos.

$$\mathcal{G} = (\mathcal{N}, \mathcal{E}) = \bigcup_{\mathcal{H} \in \{\mathcal{V}_1^t, \dots, \mathcal{V}_n^t\}} (\mathcal{N}_{\mathcal{H}}, \mathcal{E}_{\mathcal{H}}),$$



# VideoRAG

## Multi-Modal Video Knowledge Indexing

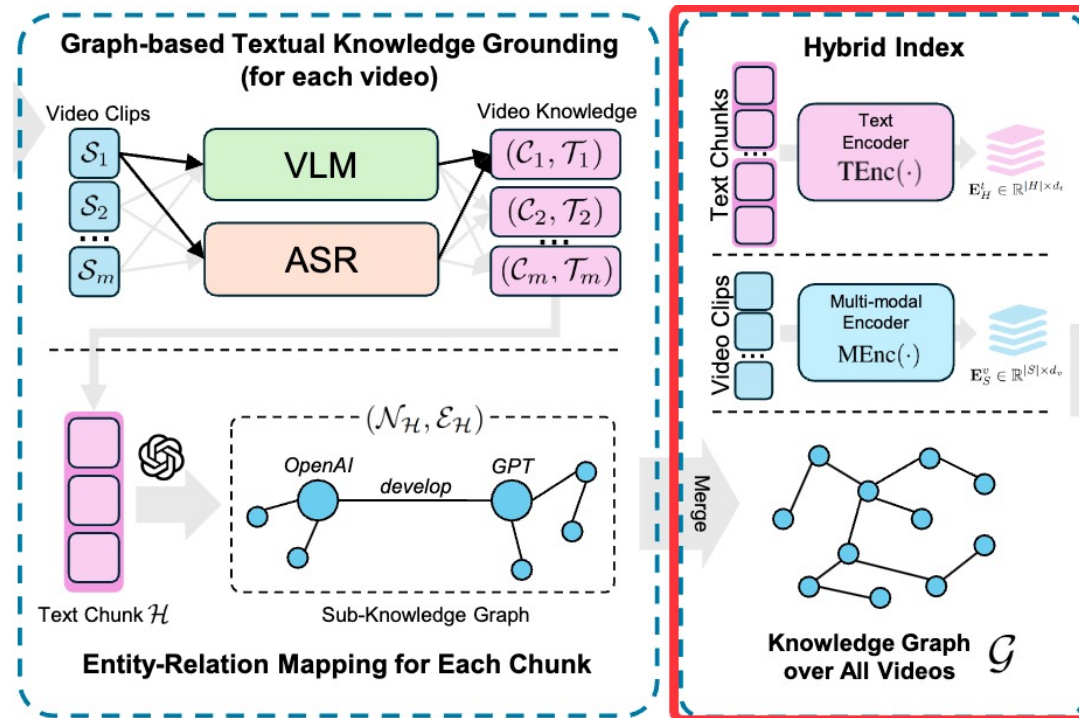
- Meanwhile, We encoded multi-modal embeddings for textual chunks and video clips separately with multi-modal embedders.

$$\mathbf{E}_H^t \in \mathbb{R}^{|H| \times d_t}$$

$$\mathbf{E}_S^v \in \mathbb{R}^{|S| \times d_v} \quad \text{w.r.t.} \quad \mathbf{e}_S^v = \text{MEnc}(\mathcal{S}).$$

- The final indexing results in a **hybrid system** that incorporates both multi-modal embeddings and knowledge graphs (KGs).

$$\hat{\mathcal{D}} = \varphi(\mathcal{D}) = (\mathcal{G}, \mathbf{E}_H^t, \mathbf{E}_S^v).$$

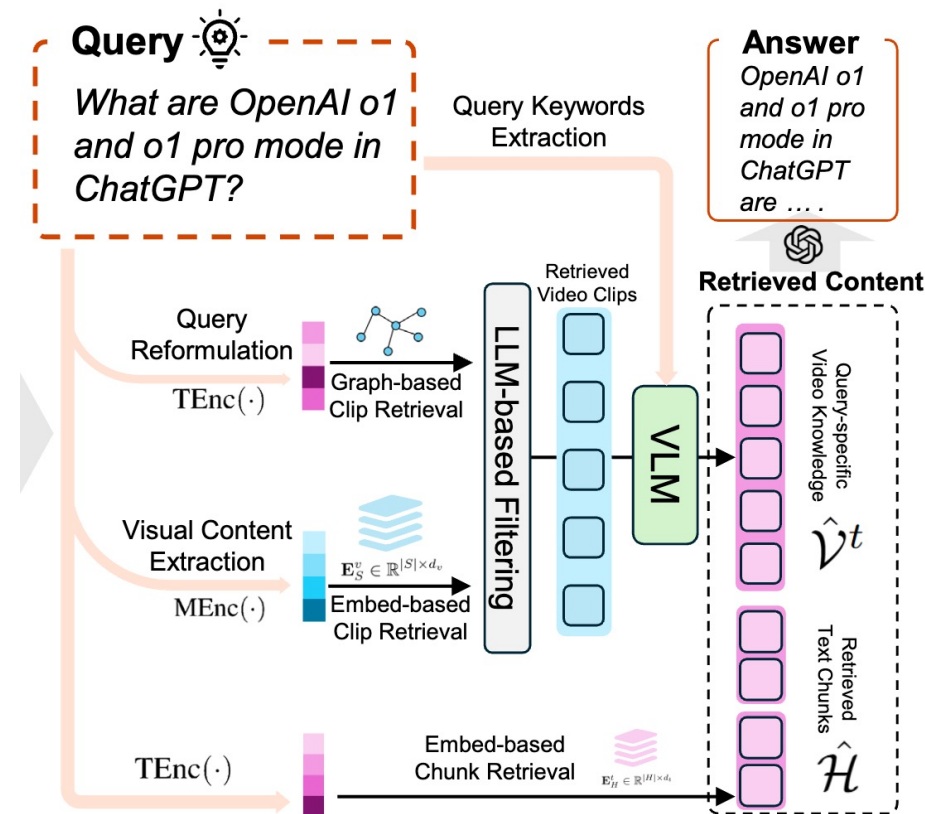


# VideoRAG

## Multi-Modal Retrieval Paradigm

- We perform video clip-level retrieval using both **textual semantic matching** and **visual retrieval** techniques.
- For the retrieved video clips, we then **filter** them by assessing their relevance to the input query.
- Finally, the VLM is employed once more to extract **fine-grained, query-specific knowledge** from the clips.

$$\hat{\mathcal{C}} = \text{VLM}(\mathcal{K}_q, \hat{\mathcal{T}}, \{\mathbf{F}_1, \dots, \mathbf{F}_{\hat{k}}\} \mid \mathbf{F} \in \hat{\mathcal{S}})$$



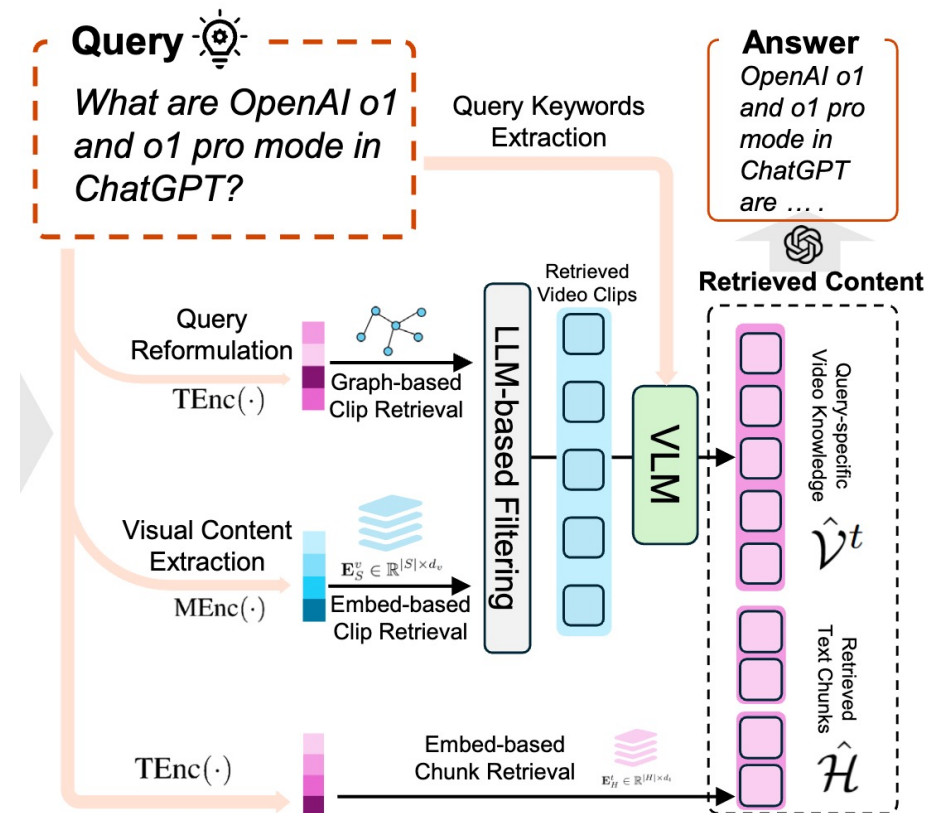
# VideoRAG

## Response Generation

- After retrieving the query-specific knowledge from the videos, we concatenate it with the extracted query-relevant text chunks for the LLMs to generate responses.

$$\psi(q, \hat{\mathcal{D}}) = (\hat{\mathcal{V}}^t, \hat{\mathcal{H}})$$

$$\text{LLM}(q, \psi(q, \hat{\mathcal{D}}))$$



# Benchmark (LongerVideos)

We constructed the *LongerVideos* benchmark to evaluate the model's performance in comprehending multiple long-context videos and answering open-ended queries.

Table 1: Statistics of the experimental dataset *LongerVideos*.

Video Type	#video list	#video	#query	#avg. queries per list	#overall duration
<b>Lecture</b>	12	135	376	31.3	~ 64.3 hours
<b>Documentary</b>	5	12	114	22.8	~ 28.5 hours
<b>Entertainment</b>	5	17	112	22.4	~ 41.9 hours
<b>All</b>	22	164	602	27.4	~ 134.6 hours

# Comparison with RAG baselines (Win- Rate)

Table 2: We analyze the performance of VideoRAG against RAG baselines on the LongerVideos dataset, presenting results both by individual video categories and across the complete dataset.

	Lecture		Documentary		Entertainment		All	
	NaiveRAG	VideoRAG	NaiveRAG	VideoRAG	NaiveRAG	VideoRAG	NaiveRAG	VideoRAG
Comprehensiveness	47.63%	<u>52.37%</u>	44.04%	<u>55.96%</u>	46.43%	<u>53.57%</u>	46.73%	<u>53.27%</u>
Empowerment	45.85%	<u>54.15%</u>	40.00%	<u>60.00%</u>	45.36%	<u>54.64%</u>	44.65%	<u>55.35%</u>
Trustworthiness	46.73%	<u>53.27%</u>	42.54%	<u>57.46%</u>	44.46%	<u>55.54%</u>	45.51%	<u>54.49%</u>
Depth	46.70%	<u>53.30%</u>	43.25%	<u>56.75%</u>	46.07%	<u>53.93%</u>	45.93%	<u>54.07%</u>
Density	46.73%	<u>53.27%</u>	44.21%	<u>55.79%</u>	44.29%	<u>55.71%</u>	45.80%	<u>54.20%</u>
Overall Winner	47.66%	<u>52.34%</u>	44.04%	<u>55.96%</u>	46.43%	<u>53.57%</u>	46.74%	<u>53.26%</u>
	GraphRAG-l	VideoRAG	GraphRAG-l	VideoRAG	GraphRAG-l	VideoRAG	GraphRAG-l	VideoRAG
Comprehensiveness	44.60%	<u>55.40%</u>	48.68%	<u>51.32%</u>	49.29%	<u>50.71%</u>	46.25%	<u>53.75%</u>
Empowerment	42.34%	<u>57.66%</u>	47.54%	<u>52.46%</u>	49.02%	<u>50.98%</u>	44.57%	<u>55.43%</u>
Trustworthiness	42.79%	<u>57.21%</u>	47.11%	<u>52.89%</u>	46.07%	<u>53.93%</u>	44.22%	<u>55.78%</u>
Depth	42.34%	<u>57.66%</u>	48.33%	<u>51.67%</u>	49.55%	<u>50.45%</u>	44.82%	<u>55.18%</u>
Density	39.26%	<u>60.74%</u>	45.26%	<u>54.74%</u>	46.52%	<u>53.48%</u>	41.74%	<u>58.26%</u>
Overall Winner	44.44%	<u>55.56%</u>	48.68%	<u>51.32%</u>	49.20%	<u>50.80%</u>	46.13%	<u>53.87%</u>
	GraphRAG-g	VideoRAG	GraphRAG-g	VideoRAG	GraphRAG-g	VideoRAG	GraphRAG-g	VideoRAG
Comprehensiveness	42.66%	<u>57.34%</u>	46.23%	<u>53.77%</u>	48.48%	<u>51.52%</u>	44.42%	<u>55.58%</u>
Empowerment	39.55%	<u>60.45%</u>	44.04%	<u>55.96%</u>	48.30%	<u>51.70%</u>	42.03%	<u>57.97%</u>
Trustworthiness	38.54%	<u>61.46%</u>	41.49%	<u>58.51%</u>	43.48%	<u>56.52%</u>	40.02%	<u>59.98%</u>
Depth	40.61%	<u>59.39%</u>	45.26%	<u>54.74%</u>	47.23%	<u>52.77%</u>	42.72%	<u>57.28%</u>
Density	37.55%	<u>62.45%</u>	46.93%	<u>53.07%</u>	48.04%	<u>51.96%</u>	41.28%	<u>58.72%</u>
Overall Winner	42.23%	<u>57.77%</u>	46.32%	<u>53.68%</u>	48.75%	<u>51.25%</u>	44.22%	<u>55.78%</u>
	LightRAG	VideoRAG	LightRAG	VideoRAG	LightRAG	VideoRAG	LightRAG	VideoRAG
Comprehensiveness	42.42%	<u>57.58%</u>	45.09%	<u>54.91%</u>	43.84%	<u>56.16%</u>	43.19%	<u>56.81%</u>
Empowerment	39.55%	<u>60.45%</u>	38.95%	<u>61.05%</u>	42.05%	<u>57.95%</u>	39.90%	<u>60.10%</u>
Trustworthiness	39.52%	<u>60.48%</u>	42.11%	<u>57.89%</u>	40.00%	<u>60.00%</u>	40.10%	<u>59.90%</u>
Depth	40.13%	<u>59.87%</u>	41.93%	<u>58.07%</u>	41.96%	<u>58.04%</u>	40.81%	<u>59.19%</u>
Density	39.57%	<u>60.43%</u>	42.37%	<u>57.63%</u>	41.61%	<u>58.39%</u>	40.48%	<u>59.52%</u>
Overall Winner	42.15%	<u>57.85%</u>	44.30%	<u>55.70%</u>	43.75%	<u>56.25%</u>	42.86%	<u>57.14%</u>

# Comparison with Video Understanding baselines (Quantitative Comparison)

## Quantitative Long-context Video Understanding Performance

Metric	LLaMA-VID				VideoAgent				NotebookLM				VideoRAG Agentic thinking			
	LEC	DOC	ENT	ALL	LEC	DOC	ENT	ALL	LEC	DOC	ENT	ALL	LEC	DOC	ENT	ALL
Comprehensiveness	2.36	2.62	2.54	<b>2.44</b>	2.02	1.99	1.80	<b>1.98</b>	3.53	3.20	2.96	<b>3.36</b>	4.48	4.51	4.44	<b>4.48</b>
Empowerment	2.79	3.03	2.86	<b>2.85</b>	2.42	2.37	2.10	<b>2.35</b>	3.88	3.62	3.29	<b>3.72</b>	4.51	4.55	4.45	<b>4.51</b>
Trustworthiness	3.15	3.30	3.35	<b>3.22</b>	2.83	2.73	2.65	<b>2.78</b>	3.95	3.80	3.71	<b>3.88</b>	4.50	4.54	4.48	<b>4.50</b>
Depth	2.01	2.06	2.00	<b>2.02</b>	1.79	1.75	1.62	<b>1.75</b>	3.14	2.89	2.55	<b>2.98</b>	4.34	4.42	4.31	<b>4.35</b>
Density	3.15	3.28	3.21	<b>3.18</b>	2.82	2.73	2.52	<b>2.75</b>	4.07	3.82	3.61	<b>3.94</b>	4.59	4.63	4.56	<b>4.59</b>
Overall Score	2.36	2.61	2.54	<b>2.44</b>	2.03	2.01	1.80	<b>1.98</b>	3.54	3.21	2.97	<b>3.37</b>	4.45	4.49	4.41	<b>4.45</b>

# Ablation Study

Both variants that remove graph indexing (-Graph) or visual encoding (-Vision) result in decreased performance.

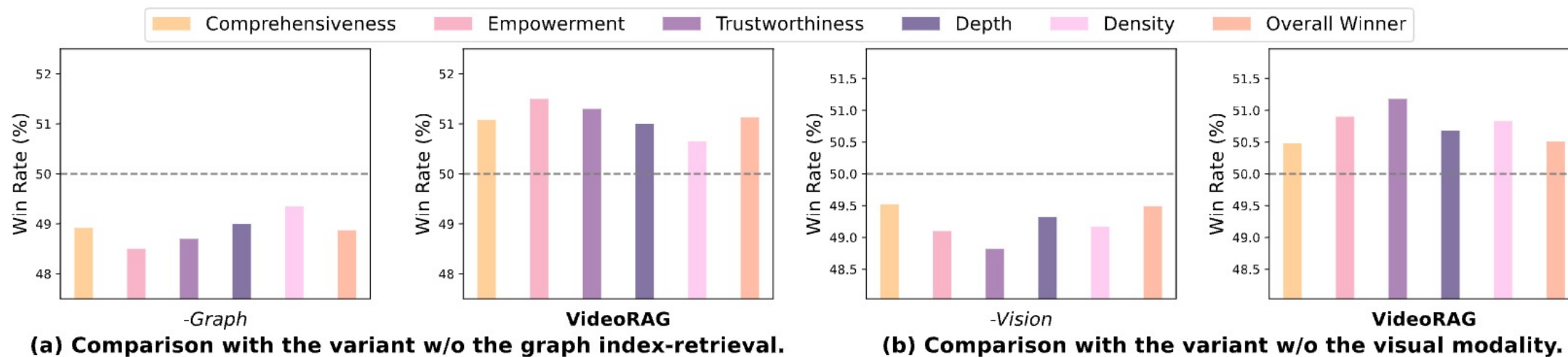


Figure 2: Ablation on graph-based knowledge grounding and cross-modal retrieval components.

# Case Study

VideoRAG can accurately retrieve four continuous video segments from OpenAI's 12-day show, specifically from 10:00 to 12:00 on Day 2 of the video titled "Reinforcement Fine-Tuning."



Reinforcement Fine-Tuning-12 Days of OpenAI: Day 2



Reinforcement Fine-Tuning-12 Days of OpenAI: Day 2

Table 4: A case study examining VideoRAG's response to a query about "graders" in OpenAI's 12-day educational series demonstrates its video retrieval capabilities and knowledge synthesis.

---

**Query:** Explain the purpose and functionality of 'graders' in the context of reinforcement fine-tuning.

---

**Input Videos:** 12 Days of OpenAI (A total of 12 videos)

---

**VideoRAG Response:**

### Understanding Gradation in Reinforcement Fine-Tuning

**\*\*Purpose of Graders\*\***

Graders are integral components in the process of reinforcement fine-tuning, mainly utilized to evaluate and improve the performance of machine learning models... (Omitted for brevity)

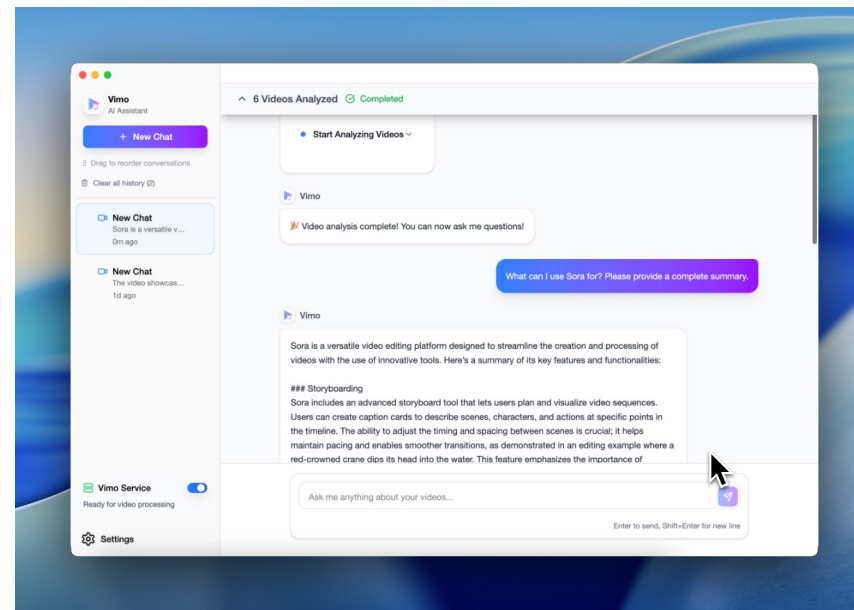
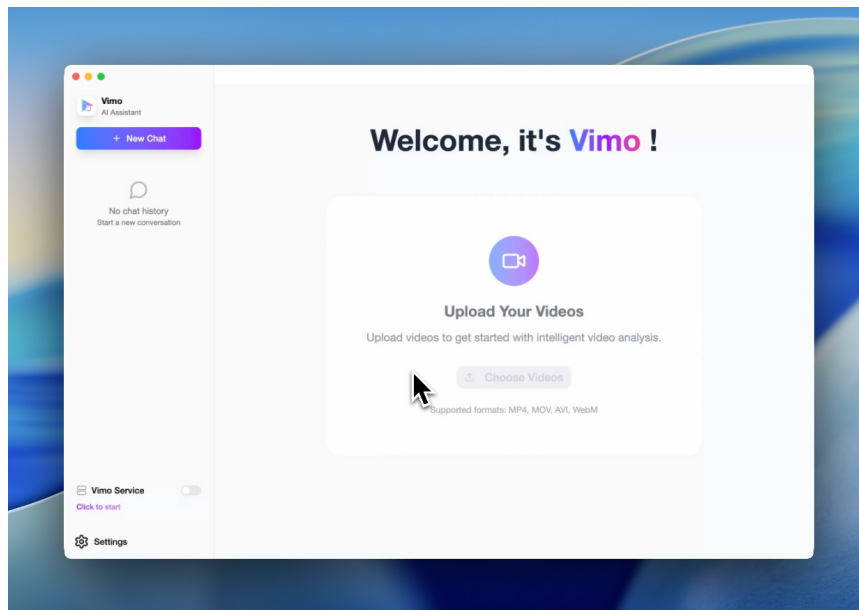
**\*\*Functionality of Graders\*\***

The grading system operates by comparing a model's outputs against those of verified correct responses. The mechanics of this assessment involve assigning scores to the model's outputs based on their correctness. These scores range from zero to one, where:

- A score of **\*\*0\*\*** indicates that the model's answer is completely incorrect.
- A score of **\*\*1\*\*** signals that the answer is fully correct.
- Scores can also be fractional, allowing for partial credit. For instance, a score of **\*\*0.7\*\*** may denote that the model's suggestion was correct but not the highest in a ranked list of possible answers.



We deliver an open-source chat platform, **Vimo**, that can understand and analyze videos of unlimited length and answer your questions with remarkable accuracy. Now everyone can **chat with videos**.

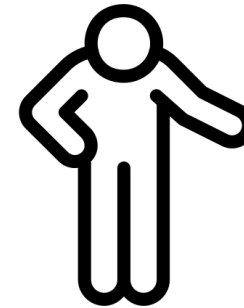




**vimo**

**Vimo: Chat with Your Videos**

✨ Data Intelligence Lab@HKU ✨



Chao Huang

Data Intelligence Lab

chaohuang75@gmail.com

<https://github.com/HKUDS>