# RAGFlow—Leading the Open-Source Revolution in Enterprise-Grade RAG

Yingfeng@InfiniFlow.ai
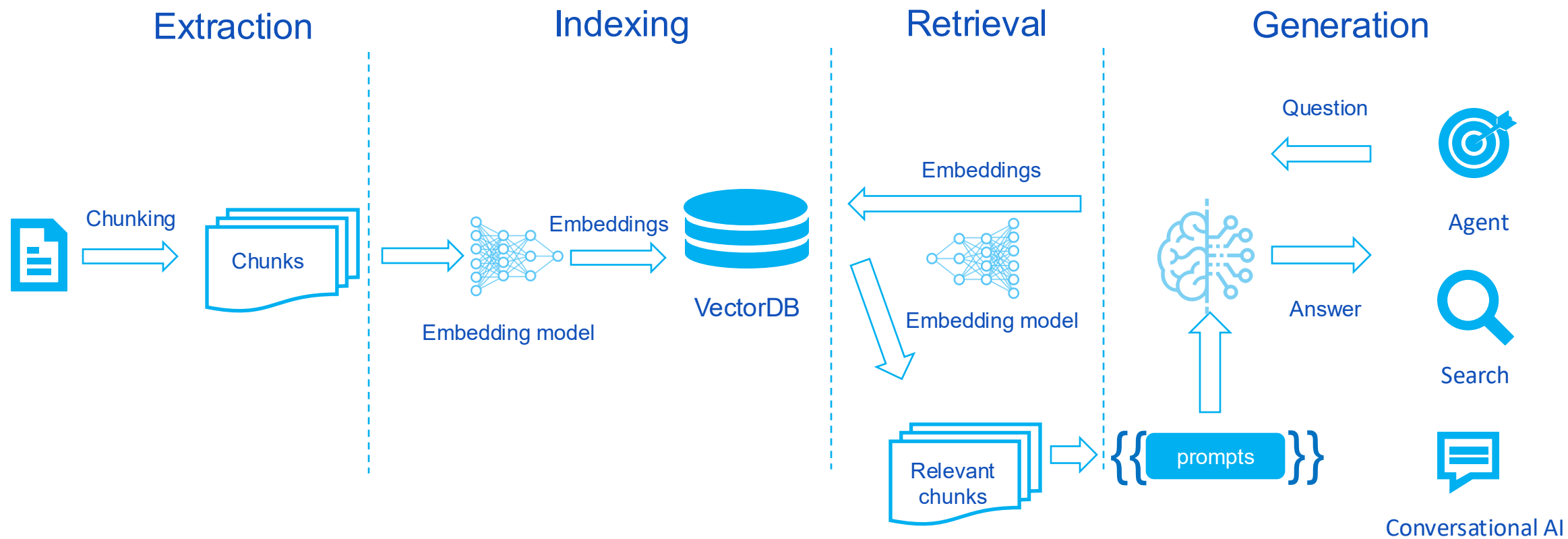
① **Introduction**

② **Chunking**

③ **Retrieval**

④ **Agentic and future**

⑤ **Summary**

# Introduction

# Paradigm of RAG



Extraction • Indexing • Retrieval • Generation

Chunking → Chunks → Embedding model → Embeddings → VectorDB → Embeddings → Embedding model → Relevant chunks → prompts → Answer → Question → Agent, Search, Conversational AI

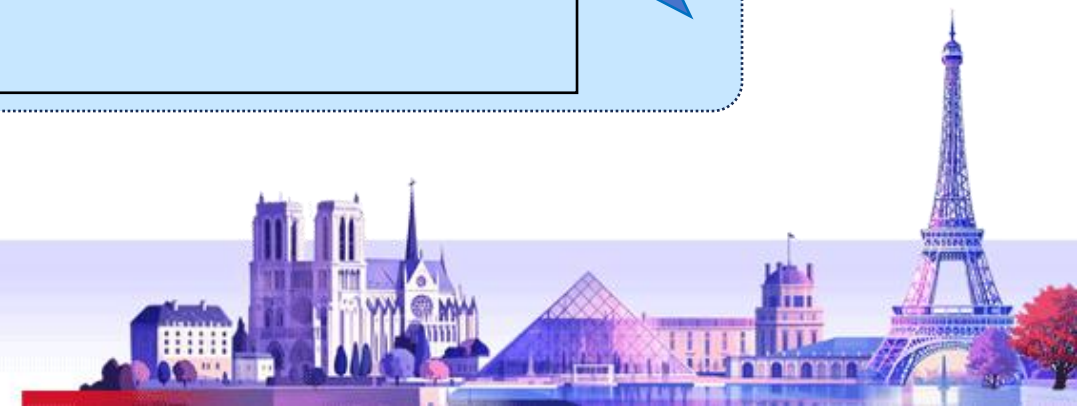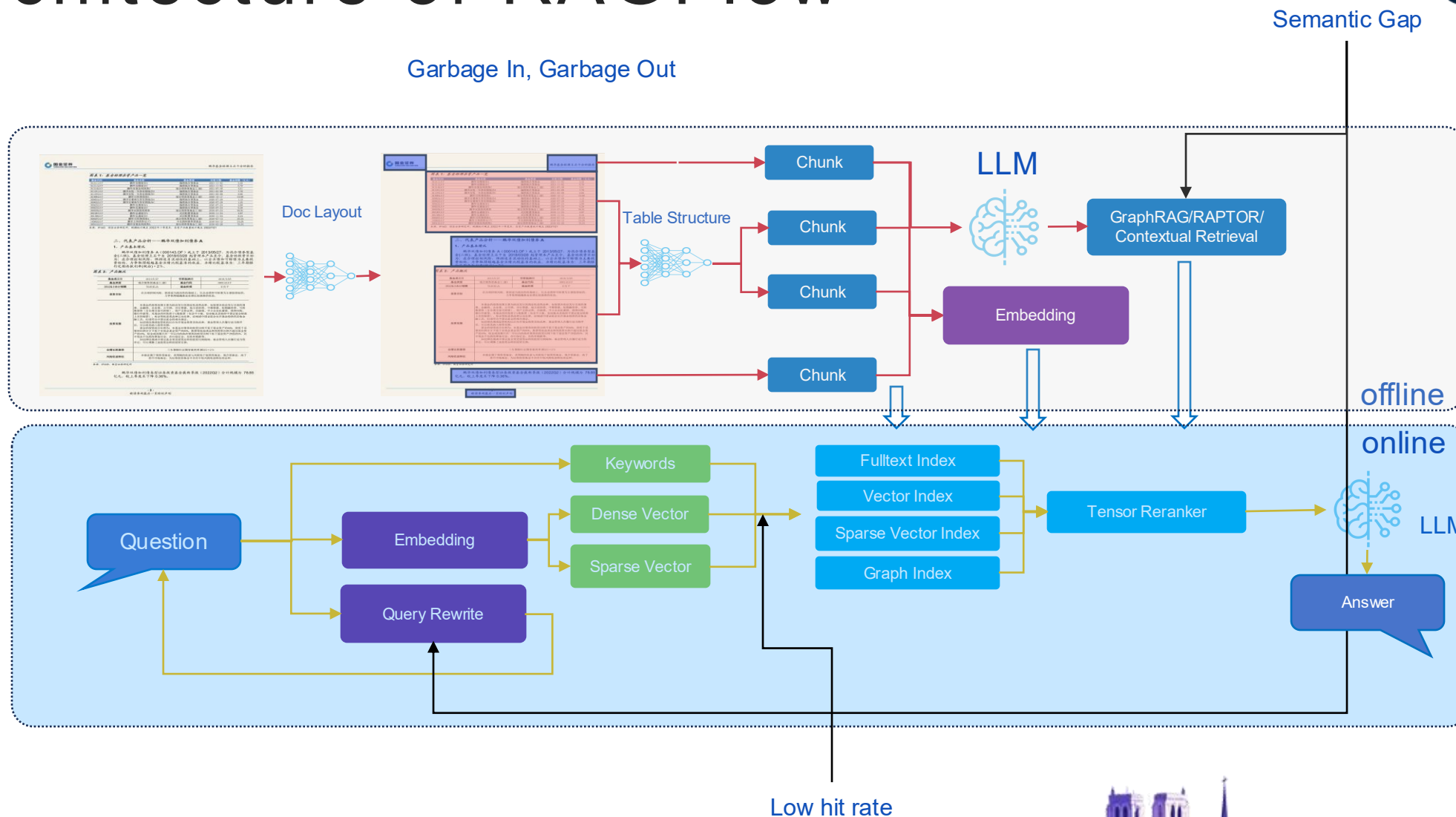# Challenges of RAG

❌  Garbage in, garbage out

❌  Low hit rate

❌  Semantic gap between question and answer

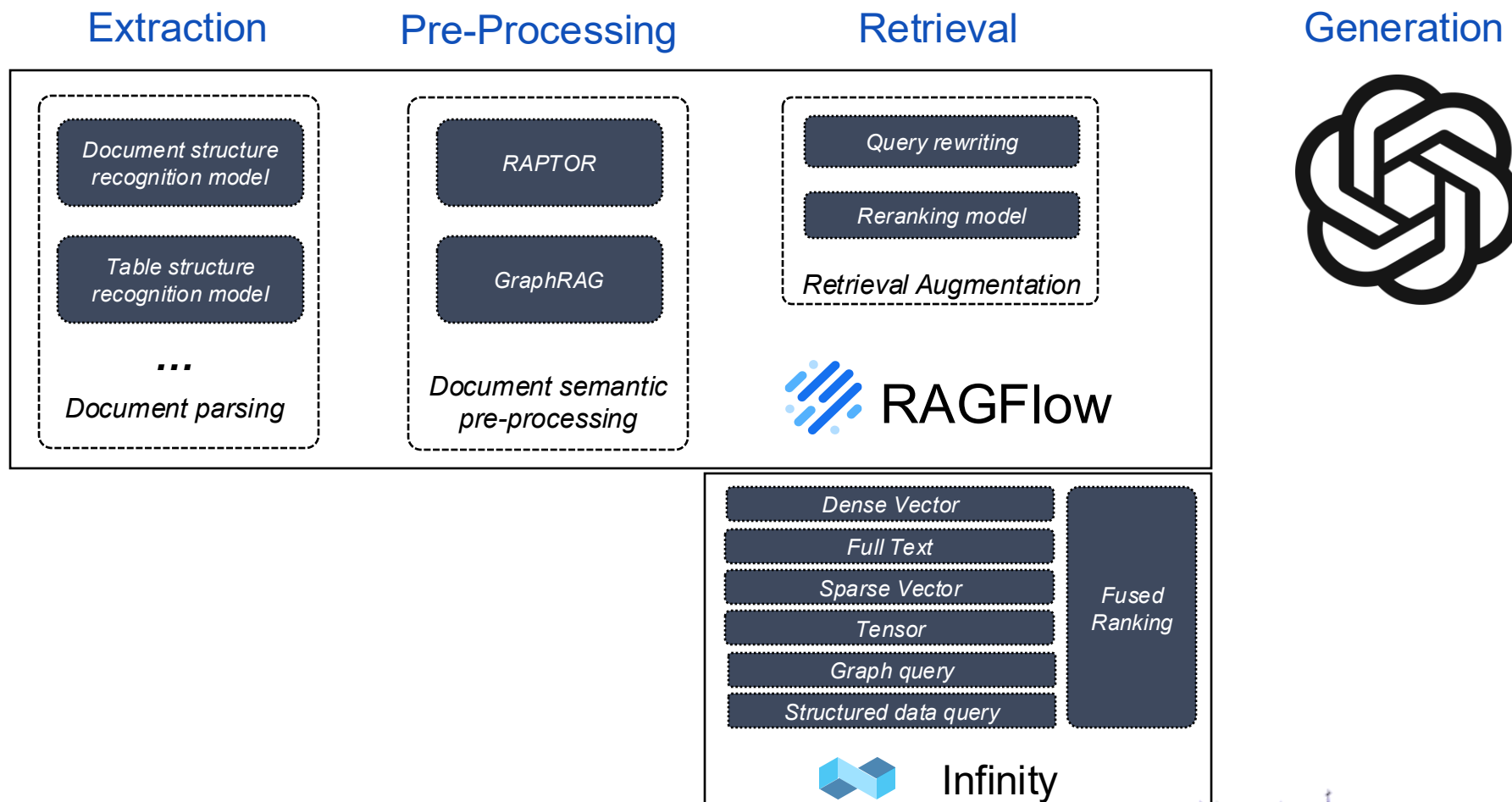# Architecture of RAGFlow

# RAG ≠ LLMOps
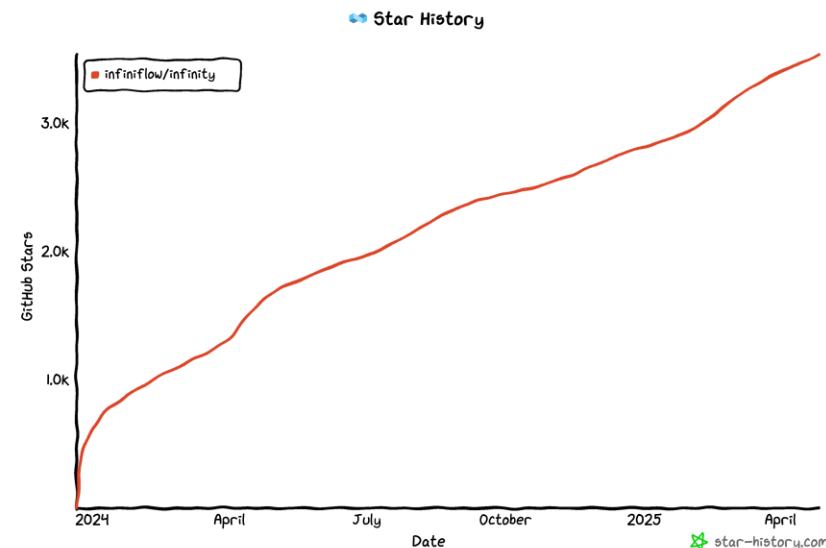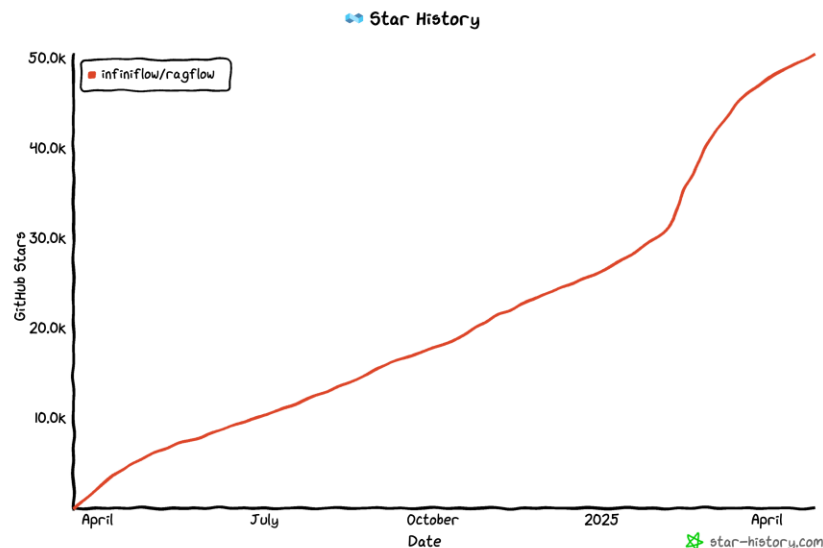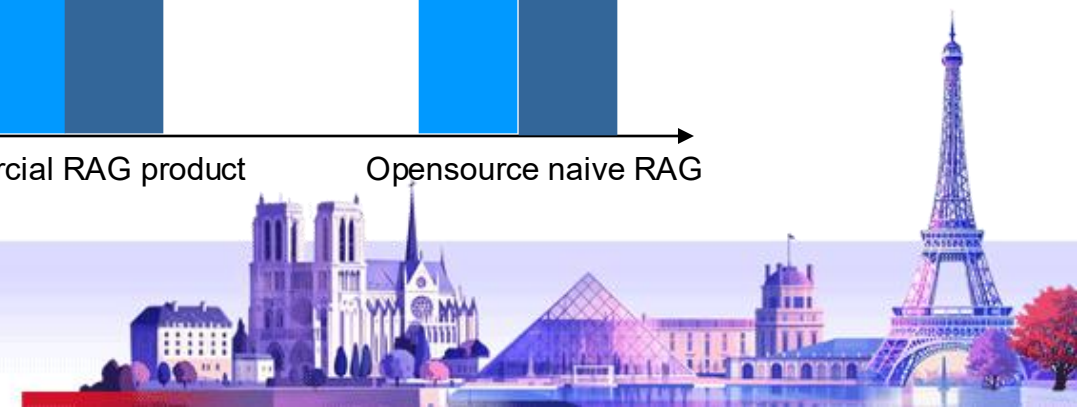
## Extraction

**Document parsing**

- Document structure recognition model
- Table structure recognition model
- ...

## Pre-Processing

**Document semantic pre-processing**

- RAPTOR
- GraphRAG

## Retrieval

**Retrieval Augmentation**

- Query rewriting
- Reranking model

**RAGFlow**

- Dense Vector
- Full Text
- Sparse Vector
- Tensor
- Graph query
- Structured data query

- Fused Ranking

**Infinity**

## Generation

# InfiniFlow = Infinity + RAGFlow



## Top Star Geo-Locations



- China
- US
- India
- Germany
- Brazil
- France
- Other

# RAGFlow vs Others



GOSIM

Accuracy

■ Perfect answer num / total answer num

■ Perfect answer num + partial correct answer num / total answer num

1.0

0.5

0.0

| | RAGFlow Pro | RAGFlow | Commercial RAG product | Opensource naive RAG |
|---|---|---|---|---|
| Perfect | 0.85 | 0.65 | 0.35 | 0.15 |
| Perfect + partial | 0.97 | 0.8 | 0.65 | 0.5 |

# Chunking

# Deep Document Understanding

Documents

### Document Parsing & Chunking

Title

Paragraph

Table

Image

Chunk 1

Chunk 2

Chunk 3

Chunk 4

Document Layout Analysis

Header/Footer

Passages

Image

Table

Scan?

N

Y

OCR

Line feed detection

Chunking

Crop

Charts/ Diagram/ ...

VLM

Chunking

Table Structure Recognition

Chunking

# Visual Chunking Results

# Table Structure Recognition

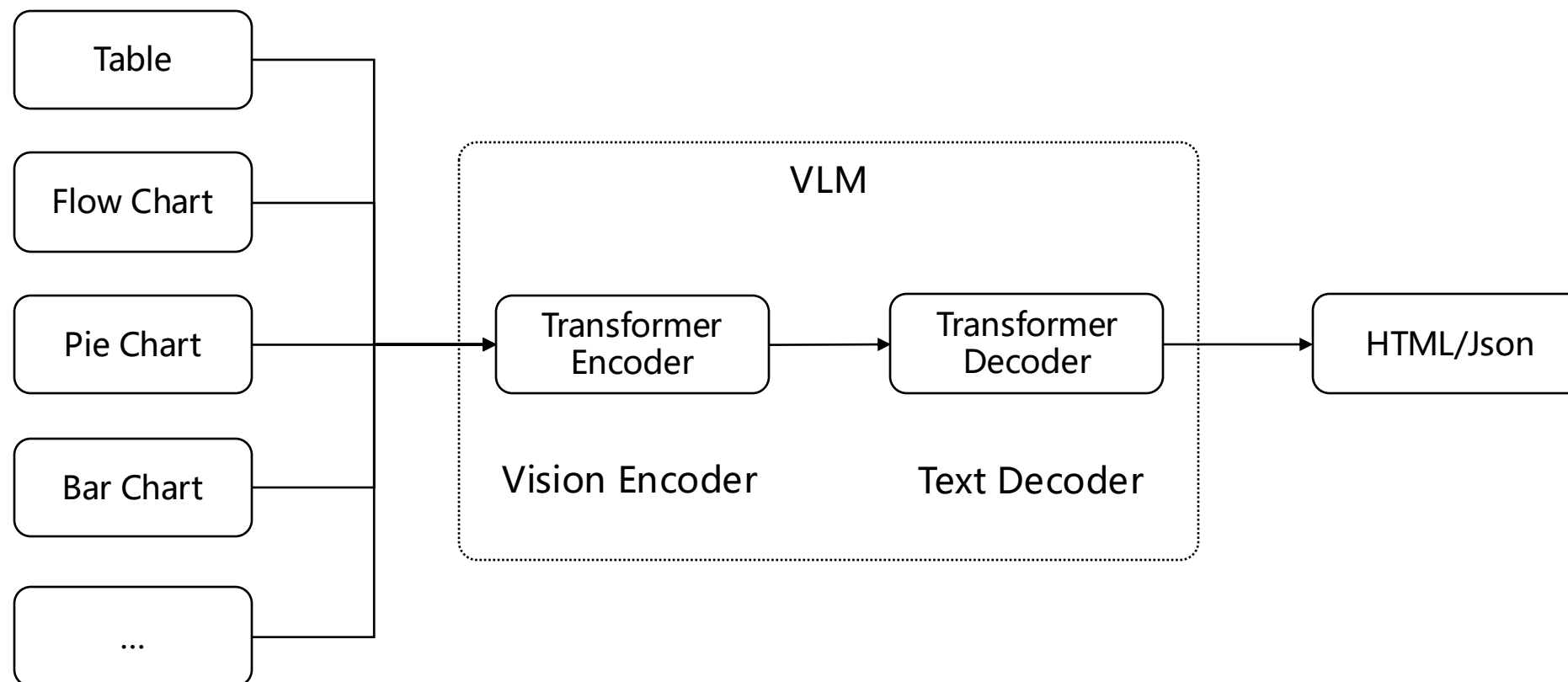| Model | CAD | | | | | | | MoCA-Mask-TE | | | | | | | ARG |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | $S_m$ ↑ | $F_\beta^\omega$ ↑ | MAE ↓ | $F_\beta$ ↑ | $E_m$ ↑ | mDice ↑ | mIoU ↑ | $S_m$ ↑ | $F_\beta^\omega$ ↑ | MAE ↓ | $F_\beta$ ↑ | $E_m$ ↑ | mDice ↑ | mIoU ↑ | |
| Ours[23] | **0.757** | **0.593** | **0.020** | **0.631** | **0.865** | **0.599** | **0.510** | **0.734** | **0.476** | **0.010** | **0.497** | **0.736** | **0.497** | **0.422** | 0.00% |
| w/ difference image | 0.724 | 0.528 | 0.023 | 0.569 | 0.761 | 0.526 | 0.441 | 0.693 | 0.401 | 0.017 | 0.429 | 0.704 | 0.425 | 0.360 | ↓15.46% |
| w/ optical flow | 0.729 | 0.534 | 0.022 | 0.573 | 0.765 | 0.532 | 0.450 | 0.697 | 0.409 | 0.016 | 0.436 | 0.708 | 0.433 | 0.370 | ↓13.41% |
| w/o intra-frame self-attention | 0.741 | 0.552 | 0.023 | 0.588 | 0.862 | 0.574 | 0.490 | 0.683 | 0.394 | 0.013 | 0.412 | 0.684 | 0.413 | 0.352 | ↓10.79% |
| w/o cross-frame cues diffusion | 0.736 | 0.546 | 0.024 | 0.584 | 0.860 | 0.569 | 0.481 | 0.679 | 0.384 | 0.013 | 0.404 | 0.681 | 0.405 | 0.340 | ↓12.17% |
| w/o temporal shifting | 0.742 | 0.571 | 0.021 | 0.610 | 0.828 | 0.576 | 0.488 | 0.690 | 0.408 | 0.013 | 0.426 | 0.692 | 0.427 | 0.366 | ↓8.88% |

➢ Multi table header

➢ Border/Bordless

➢ Cell merge

➢ Cross page

# DeepDoc—CNN or GenAI?

# DeepDoc—CNN or GenAI?

GOSIM

| Table |
| Flow Chart |
| Pie Chart | → | VLM |
| Bar Chart |
| ... |

VLM

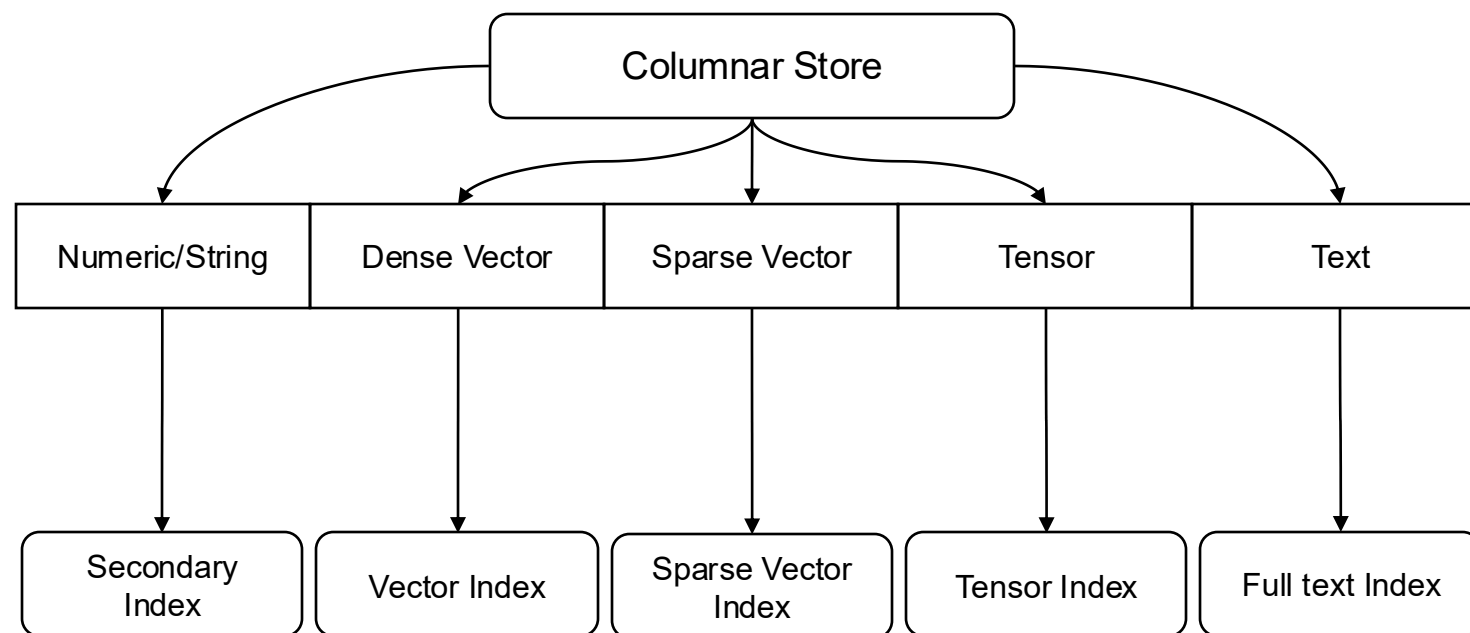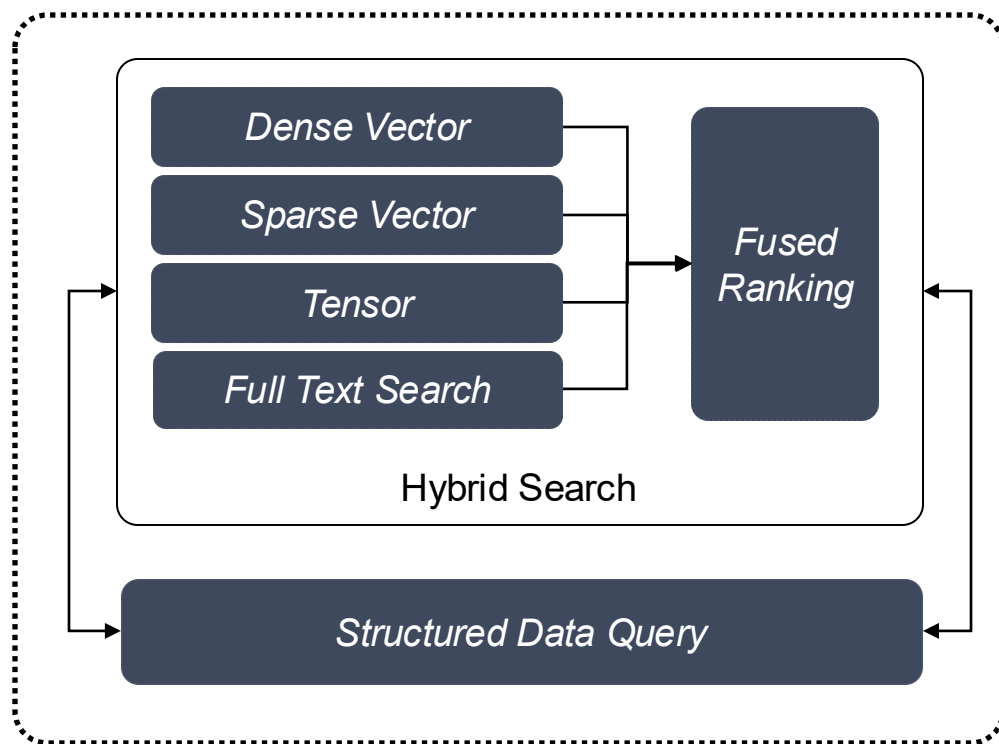| Transformer Encoder | → | Transformer Decoder | → | HTML/Json |

Vision Encoder          Text Decoder

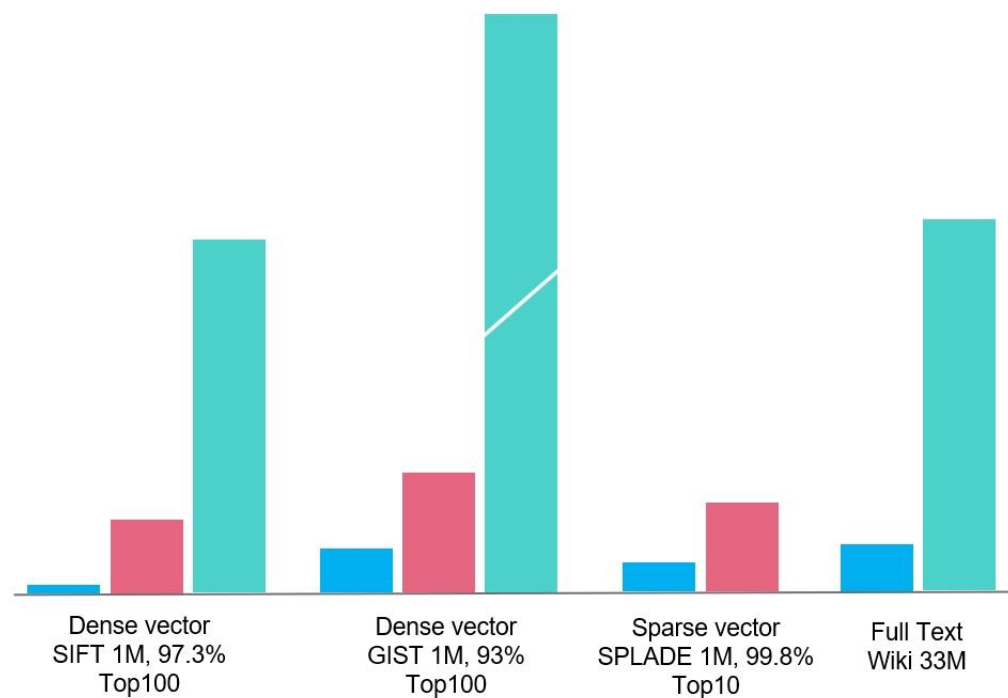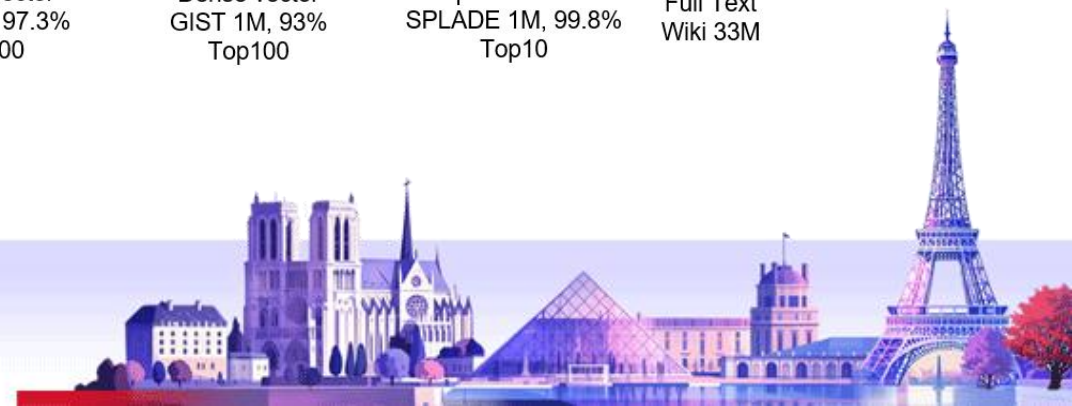# Multi-Modal RAG

Retrieval

# Infinity—Indexing Database

# Benchmark



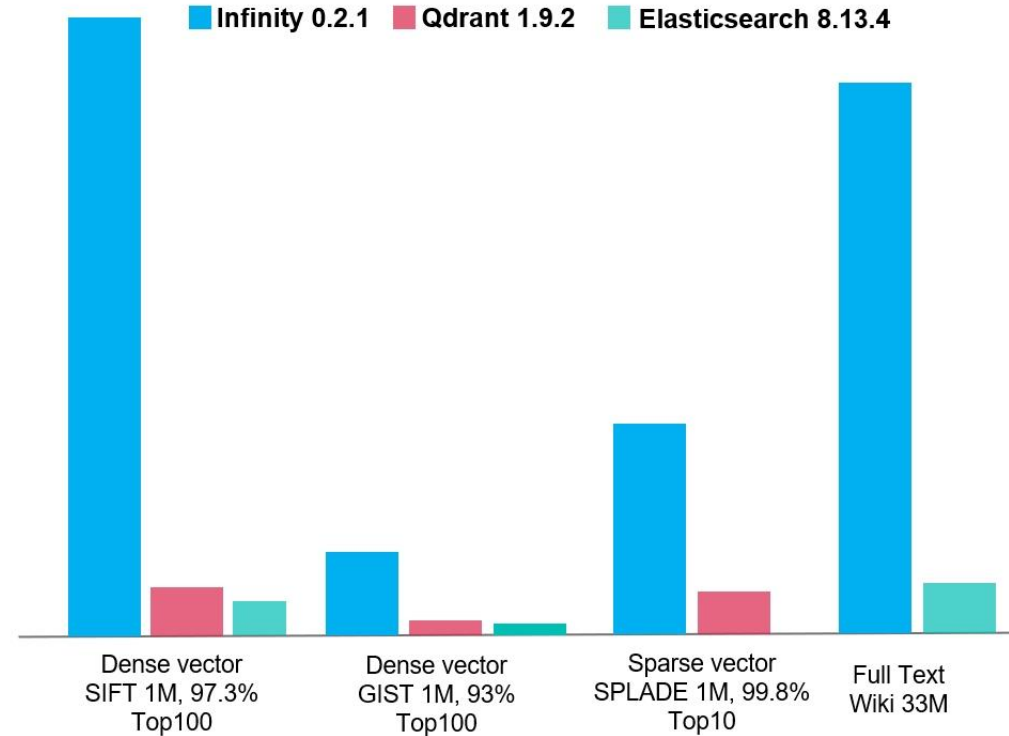**Query latency benchmark: Lower is better**
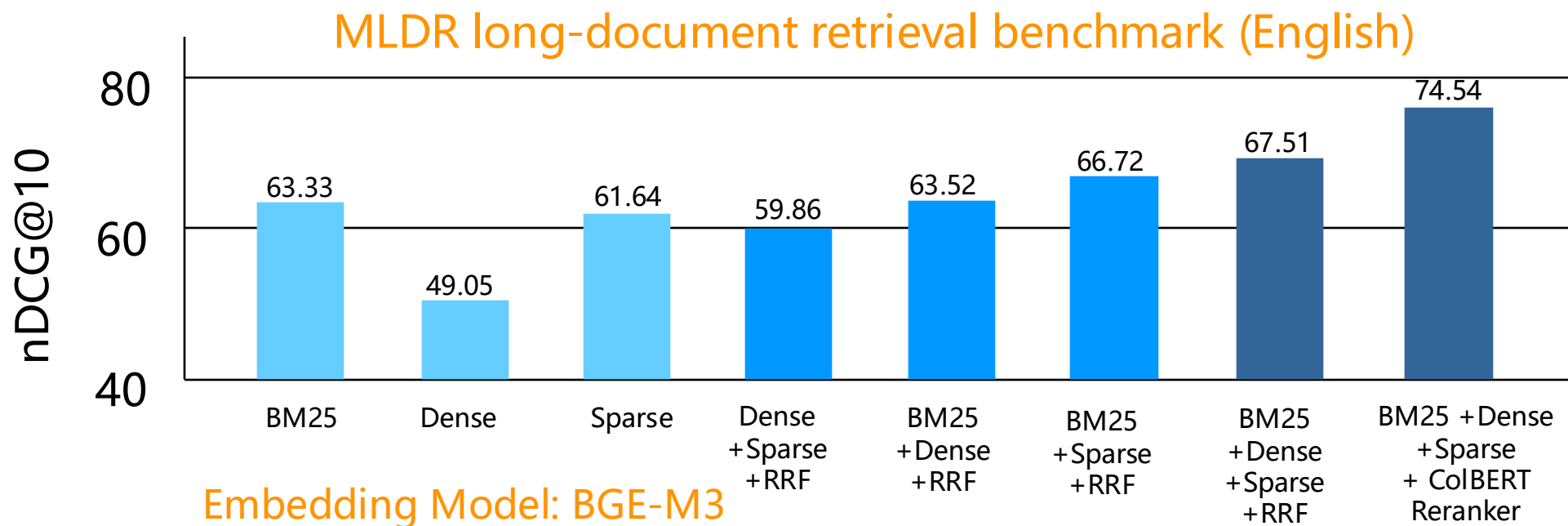
■ Infinity 0.2.1  ■ Qdrant 1.9.2  ■ Elasticsearch 8.13.4

Dense vector
SIFT 1M, 97.3%
Top100

Dense vector
GIST 1M, 93%
Top100

Sparse vector
SPLADE 1M, 99.8%
Top10

Full Text
Wiki 33M

**Query QPS benchmark: Higher is better**

■ Infinity 0.2.1  ■ Qdrant 1.9.2  ■ Elasticsearch 8.13.4

Dense vector
SIFT 1M, 97.3%
Top100

Dense vector
GIST 1M, 93%
Top100

Sparse vector
SPLADE 1M, 99.8%
Top10

Full Text
Wiki 33M

# Hybrid Search

**GOSIM**

## MLDR long-document retrieval benchmark (English)



Embedding Model: BGE-M3

| Method | nDCG@10 |
|---|---|
| BM25 | 63.33 |
| Dense | 49.05 |
| Sparse | 61.64 |
| Dense +Sparse +RRF | 59.86 |
| BM25 +Dense +RRF | 63.52 |
| BM25 +Sparse +RRF | 66.72 |
| BM25 +Dense +Sparse +RRF | 67.51 |
| BM25 +Dense +Sparse + ColBERT Reranker | 74.54 |

IBM Blended RAG
https://arxiv.org/abs/2404.07220

# Benefits of ColBERT



MLDR long-document retrieval benchmark (English)

Embedding Model: BGE-M3

# Ranking Model



**Dual Encoder**

**Cross Encoder**

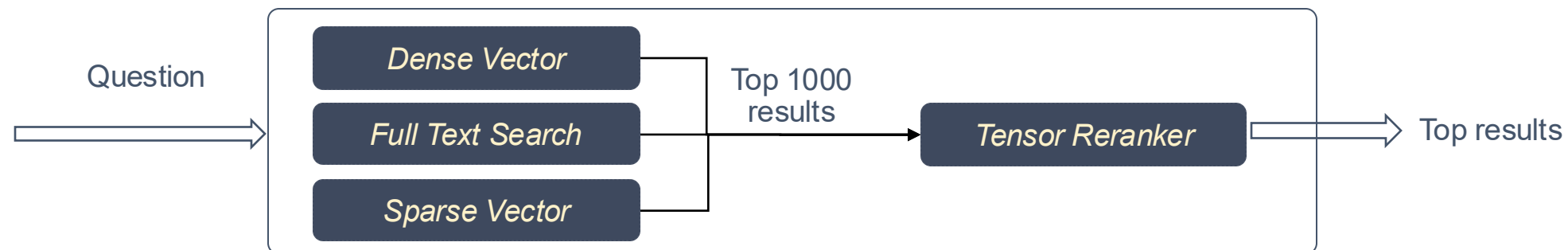**Late Interaction Encoder ColBERT**

# ColBERT Ranker vs Reranker

GOSIM



MLDR long-document retrieval benchmark (English)

nDCG@10

| | 72.23 | 73.35 | 74.11 |
|---|---|---|---|

ColBERT EMVB Index | BM25 +ColBERT Reranker | ColBERT Brute force

Embedding Model: BGE-M3

# Benefits of Col-xxx Rerankers



GOSIM

Question → VectorDB → Question / Top 10 results → Cross Encoder Reranker → Top results

**VS**

Question → [ Dense Vector / Full Text Search / Sparse Vector ] → Top 1000 results → Tensor Reranker → Top results

# Col-xxx in Multi-Modal-RAG

Divides a PDF page into 32 × 32 = 1024 patches

Each patch is represented as a 128-dim vector

# ViDoRe Benchmark

Qwen2 VL 7B => ColQwen2

PaliGemma => ColPali

| Rank | Model | Average | TAT-DQA | Shift Project | Artificial Intelligence | Government Reports | ArxivQA | DocVQA | Healthcare Industry |
|---|---|---|---|---|---|---|---|---|---|
| 1 | Metric-AI_ColQwen2.5-7b-multi | 90.4 | 80.9 | 87.7 | 99.3 | 95.2 | 91.7 | 65.1 | 97.8 |
| 2 | Metric-AI_ColQwen2.5-3b-multi | 90.3 | 80.7 | 88.4 | 98.9 | 96.4 | 92.2 | 64.4 | 98.4 |
| 3 | yydxlv_colqwen2.5-7b-v0.1 | 90.2 | 81.9 | 89.1 | 98.9 | 96.2 | 91.1 | 63.1 | 98.5 |
| 4 | tsystems_colqwen2-7b-v1.0 | 90.1 | 78.6 | 89.3 | 99.3 | 96.3 | 90.7 | 64.5 | 99.3 |
| 5 | Metric-AI_colqwen2.5-3b-multi | 90 | 80.7 | 86.5 | 99.4 | 97.7 | 93 | 65.1 | 98.8 |
| 6 | yydxlv_colqwen2-7b-v1.0 | 90 | 81.7 | 87.8 | 100 | 96.7 | 89.3 | 63.1 | 98.7 |
| 7 | Alibaba-NLP_gme-Qwen2-VL-7B-I | 89.9 | 76.1 | 97.3 | 99.6 | 99.5 | 87.6 | 56.6 | 99.6 |
| 8 | Metric-AI_ColQwenStella-2b-mu | 89.8 | 78.7 | 88.1 | 99.6 | 97.1 | 89.5 | 62.6 | 99.6 |
| 9 | tsystems_colqwen2-2b-v1.0 | 89.6 | 79.5 | 89.9 | 99.6 | 95.2 | 90 | 60.9 | 98.9 |
| 10 | vidore_colqwen2.5-v0.2 | 89.4 | 81.1 | 87.3 | 99.6 | 96.4 | 89.2 | 63.2 | 97.9 |
| 11 | vidore_colqwen2-v1.0 | 89.3 | 81.4 | 90.7 | 99.4 | 96.3 | 88.1 | 60.6 | 98.1 |
| 12 | vidore_colqwen2.5-v0.1 | 88.8 | 80.8 | 85.5 | 99.3 | 95.3 | 88.5 | 61.9 | 98.8 |
| 13 | Alibaba-NLP_gme-Qwen2-VL-2B-I | 87.8 | 71.1 | 94.3 | 99 | 97.9 | 83.9 | 54.6 | 98.9 |
| 14 | vidore_colqwen2-v0.1 | 87.3 | 75.9 | 86 | 98.7 | 92.8 | 86.1 | 61.4 | 98 |
| 15 | vidore_colsmolvlm-v0.1 | 86.1 | 79.5 | 79.5 | 98.1 | 96.9 | 79.4 | 60 | 99.6 |
| 16 | MrLight_dse-qwen2-2b-mrl-v1 | 85.8 | 69.4 | 82 | 97.5 | 96 | 85.6 | 57.1 | 96.4 |
| 17 | vidore_colpali-v1.3 | 84.8 | 70.4 | 77.4 | 97.4 | 96.2 | 83 | 58.5 | 96.9 |
| 18 | vidore_colpali2-3b-pt-448 | 84.5 | 68.6 | 75.9 | 98 | 94.1 | 82.5 | 59 | 97.2 |
| 19 | vidore_colpali-v1.2 | 83.9 | 68 | 79.1 | 98.1 | 94.8 | 78 | 57.2 | 96.7 |
| 20 | yydxlv_colphi3.5 | 83.7 | 73.1 | 68.2 | 98.5 | 95.2 | 87.2 | 56.9 | 95.9 |

# Col-xxx based Multi-Modal RAG



GOSIM

page1

page2

PDF1

PDF2

page1

page2

ColPali

Tensor

VLM

What's China's IDC business market size in 2018?
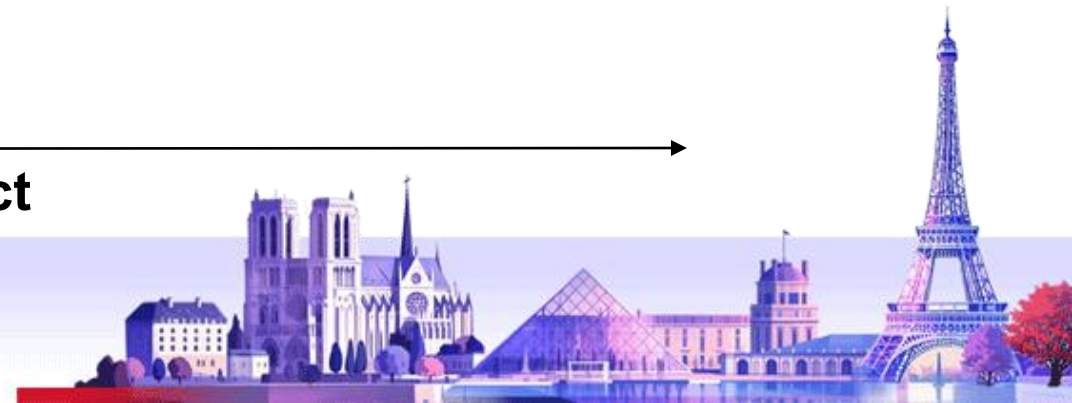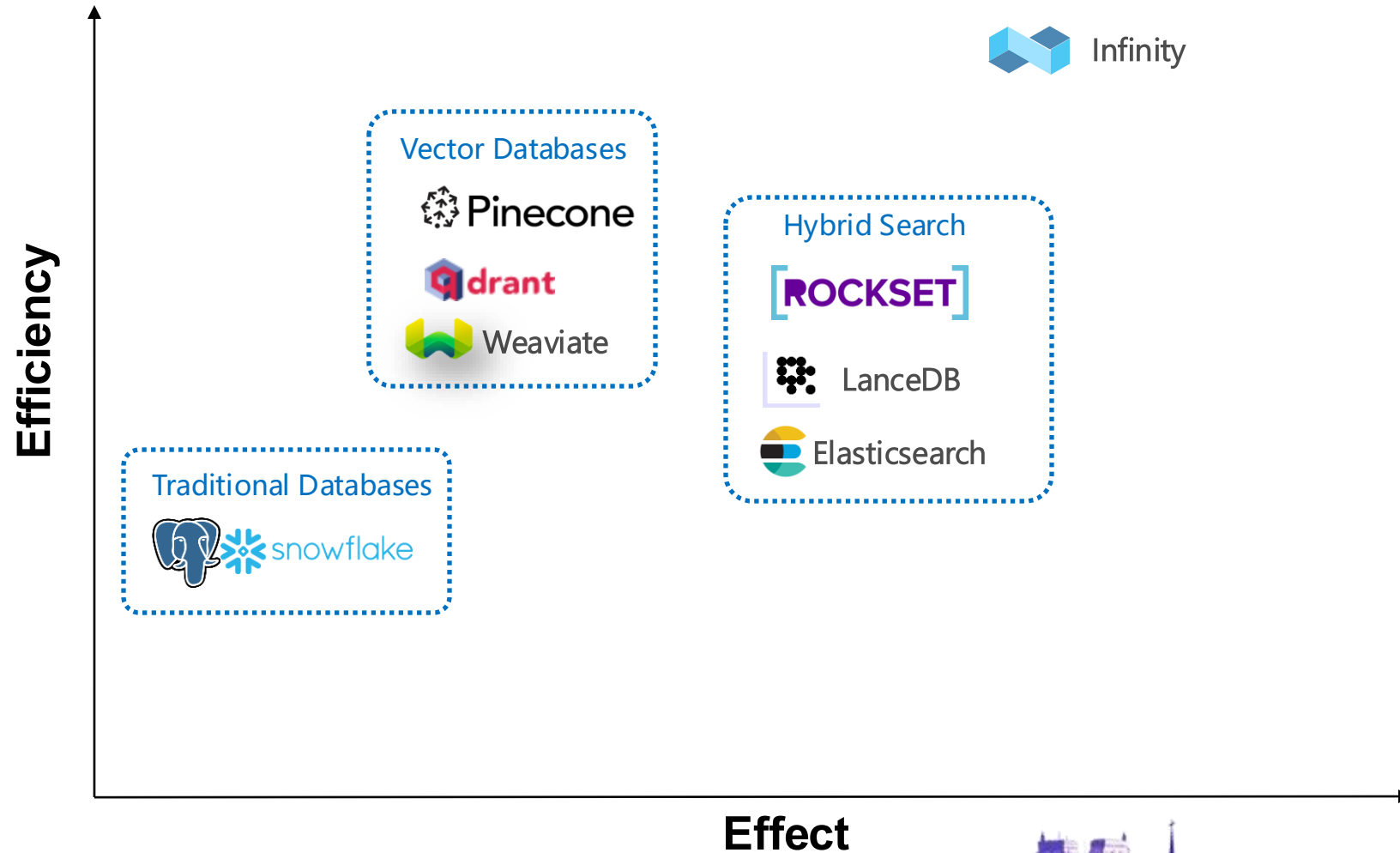
Answer: 1228

GOSIM AI Paris 2025

# Scaling of Col-xxx Support

Data base Side
- ➤ Tensor Reranker
- ➤ Binary Quantization
- ➤ Multi-Vector Index
- ➤ Full text index

Model Side
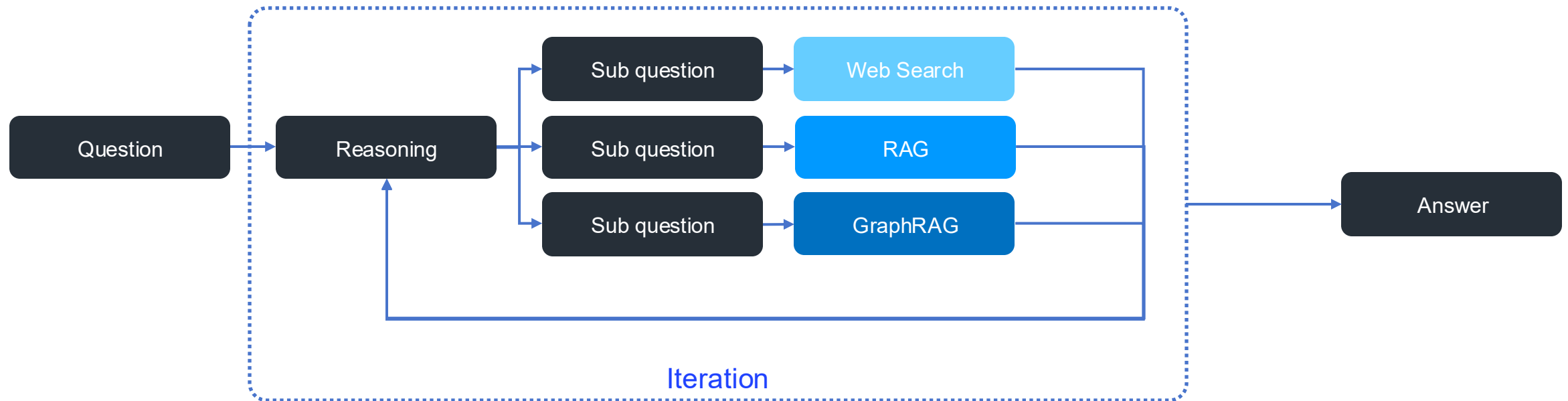- ➤ Dimensional Reduction（MRL）
- ➤ Tokens Reduction（ConstBERT/SVD）

# Comparison

# Agentic and Future

# Agentic Reasoning

# Agentic Reasoning

# RAG ⇔ Agent

➤ RAG As a Distinct Layer Grows More Vital

➤ RAG & Agent:
Frequent Interaction, Collaboration and Competition Coexist

# The Essence of Memory



LLM — Reasoning — Infinite Context

Filtering:
- Hybrid Search
- Reranker
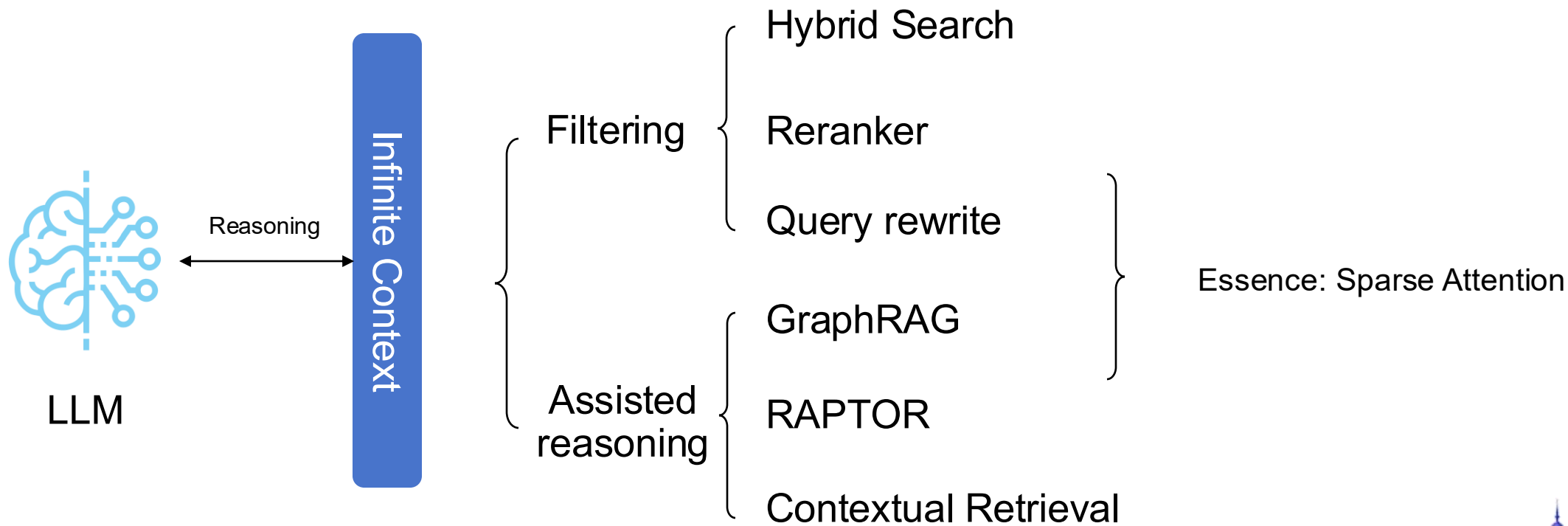- Query rewrite

Assisted reasoning:
- GraphRAG
- RAPTOR
- Contextual Retrieval

Essence: Sparse Attention
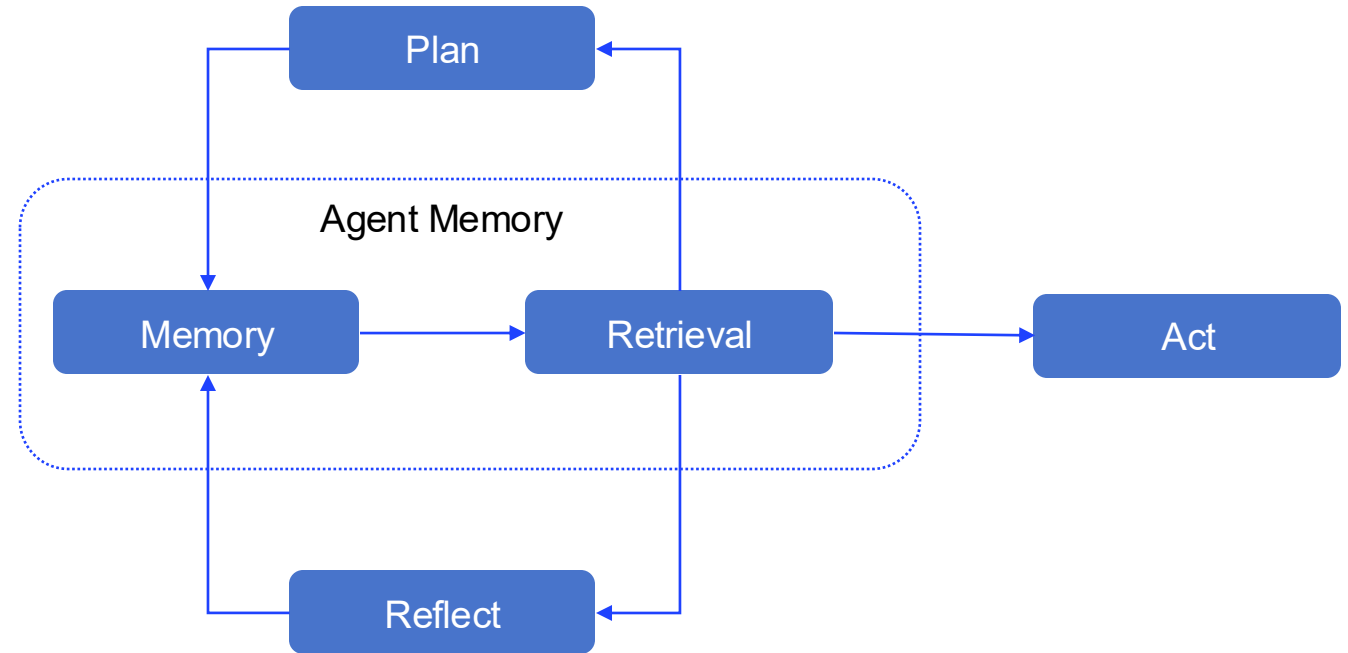
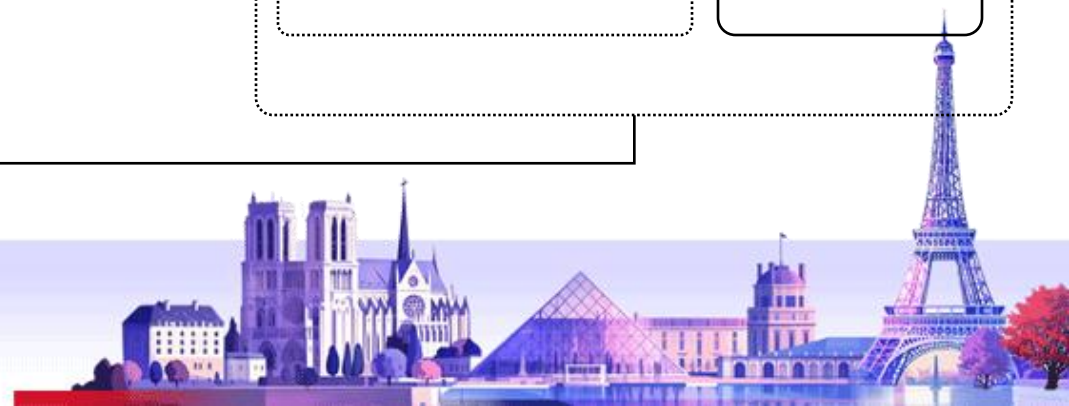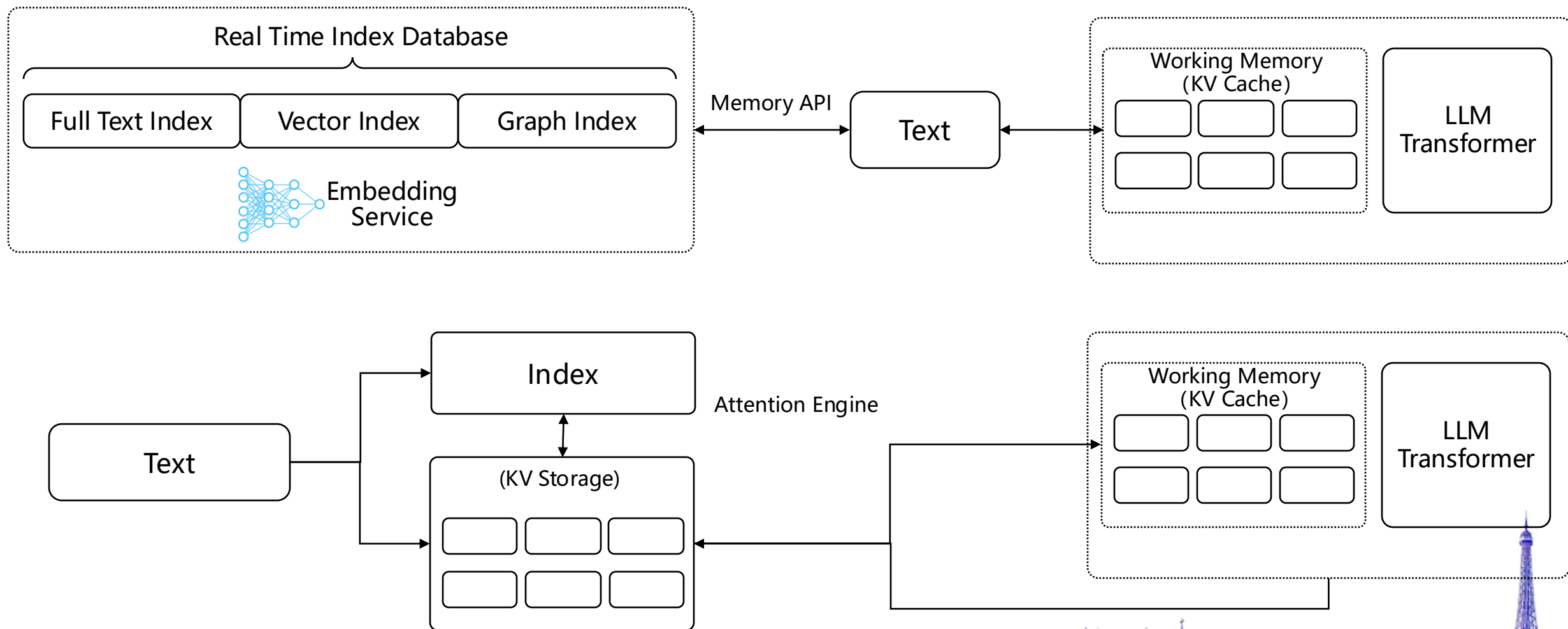# Agent Memory = RAG ?

➢ Retrieval Differences—Temporal/Weights

➢ Memory Decay

➢ Agent Facilities

# RAG/Memory Evolution



Real Time Index Database

Full Text Index | Vector Index | Graph Index

Embedding Service

Memory API

Text

Working Memory (KV Cache)

LLM Transformer

Index

Attention Engine

Text

(KV Storage)

Working Memory (KV Cache)

LLM Transformer

# Summary

# RAGFlow Ecosystem

GOSIM

| AI Search | Customer Service | Guaranteed Compliance | Decision-making Assistant | Talent Development |
|---|---|---|---|---|

**LLMs**

**RAGFlow**

*Agent Orchestration*

*Small AI Models*
*Document Understanding / Col-xxx Reranker / Knowledge Extraction / …*

*Infinity*

**Data Flows**

## RAG is not LLMOps, RAG is the database of LLM era

# THANK YOU

https://github.com/infiniflow/ragflow
https://github.com/infiniflow/infinity