

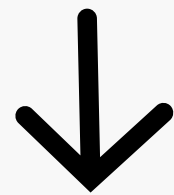
**ОТ ПОИСКА К ПОНИМАНИЮ:
ЭВОЛЮЦИЯ RAG — КАК
ПРЕВРАТИТЬ КОРПОРАТИВНЫЕ
ЗНАНИЯ В ИНТЕЛЛЕКТУАЛЬНОГО
СОБЕСЕДНИКА**

Жарков Олег

ПОЧЕМУ LLM НЕ МОЖЕТ БЫТЬ БЕЗ RAG?

01

Галлюцинации = ущерб
бизнесу



- Юридические риски
- Финансовые риски
- Репутационные риски

02

Устаревание знаний

04

Отсутствие
прозрачности и аудита

03

Модели недоступны
внутренние данные
компании

LONG CONTEXT WINDOW VS RAG

“

Зачем сейчас нужен RAG, если у моделей уже контекст на миллион токенов? Давайте просто положим туда все документы.

”

LONG CONTEXT WINDOW VS RAG

[Question]

"What was the best writing advice I got from my college classmate?"

[Needle]

[Haystack]

I've discovered a handy test for figuring out what you're addicted to.

Imagine you were going to spend the weekend at a friend's house on a little island off the coast of Maine. There are no shops on the island and you won't be able to leave while you're there. Also, you've never been to this house before, so you can't assume it will have more than any house might.

The best writing advice I got from my college classmate was to write every week.

What, besides clothes and toiletries, do you make a point of packing? That's what you're addicted to.

For example, if you find yourself packing a bottle of vodka (just in case), you may want to stop and think about that. For me the list is four things: books, earplugs, a notebook, and a pen. There are other things I might bring if I thought of it, like music, or tea, but I can live without them.

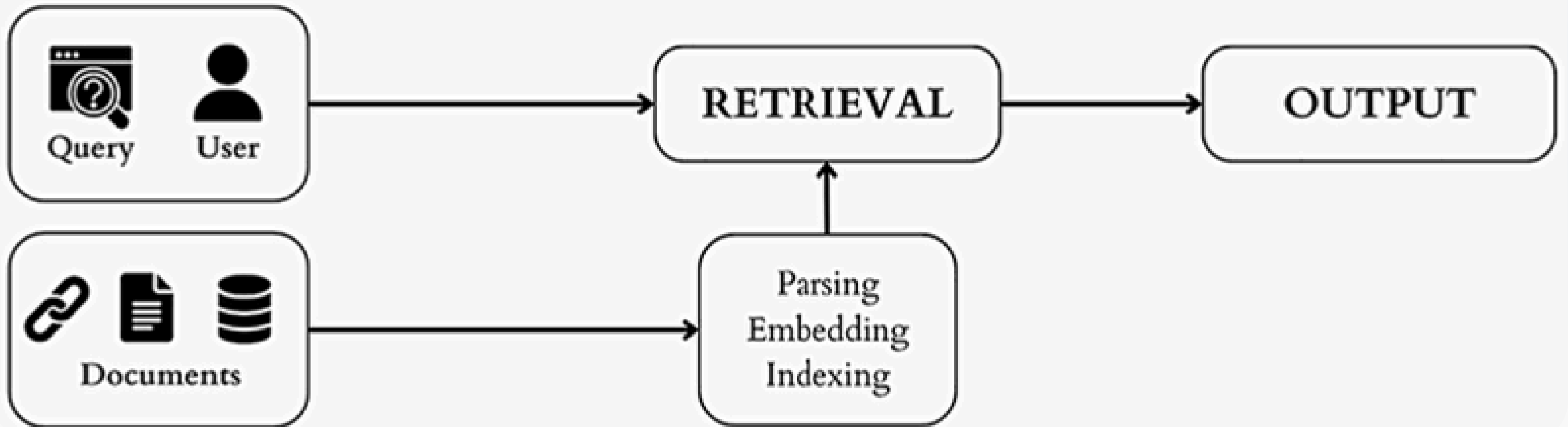
I'm not so addicted to caffeine that I wouldn't risk the house not having any tea, just for a weekend.

Миф о “бесконечном контексте”:

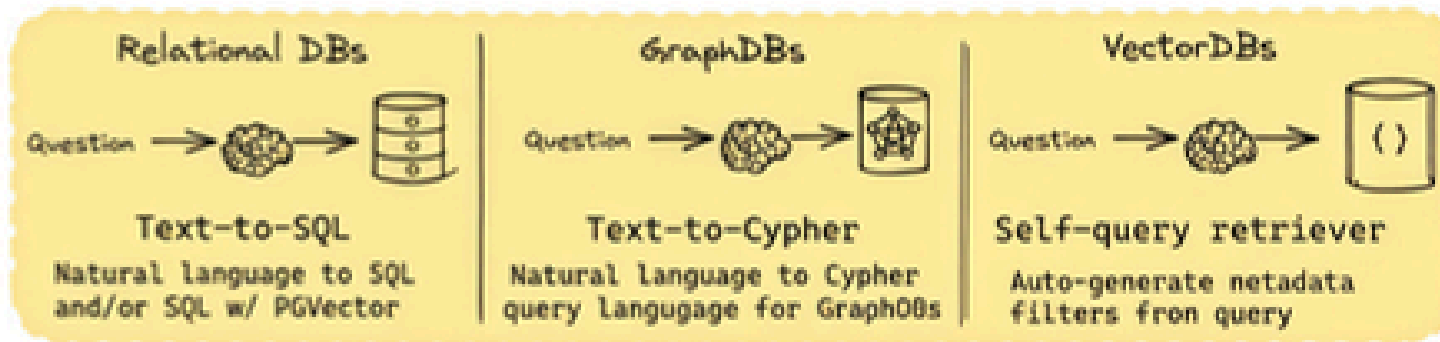
- Lost in the middle
- Цена и latency
- Масштабирование

ЭВОЛЮЦИЯ АРХИТЕКТУРЫ RAG

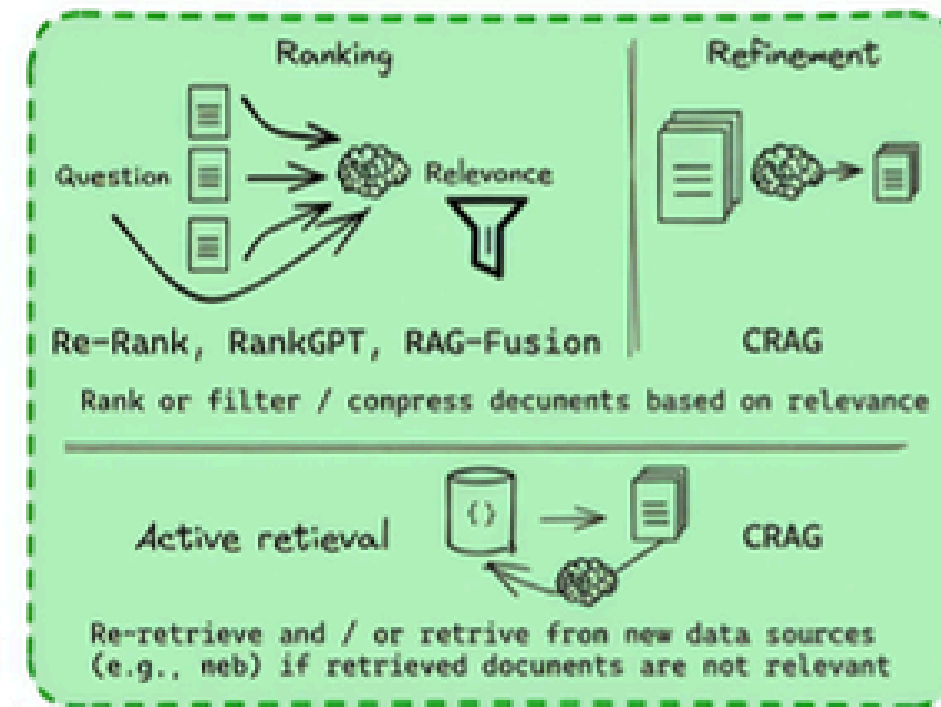
NAIVE RAG



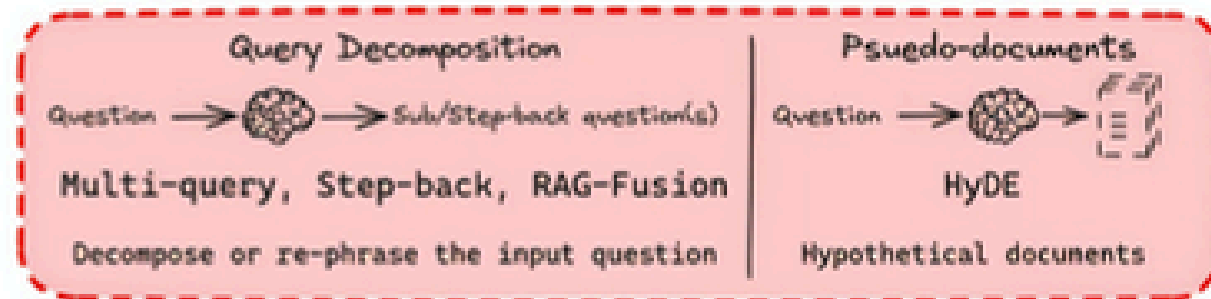
Query Construction



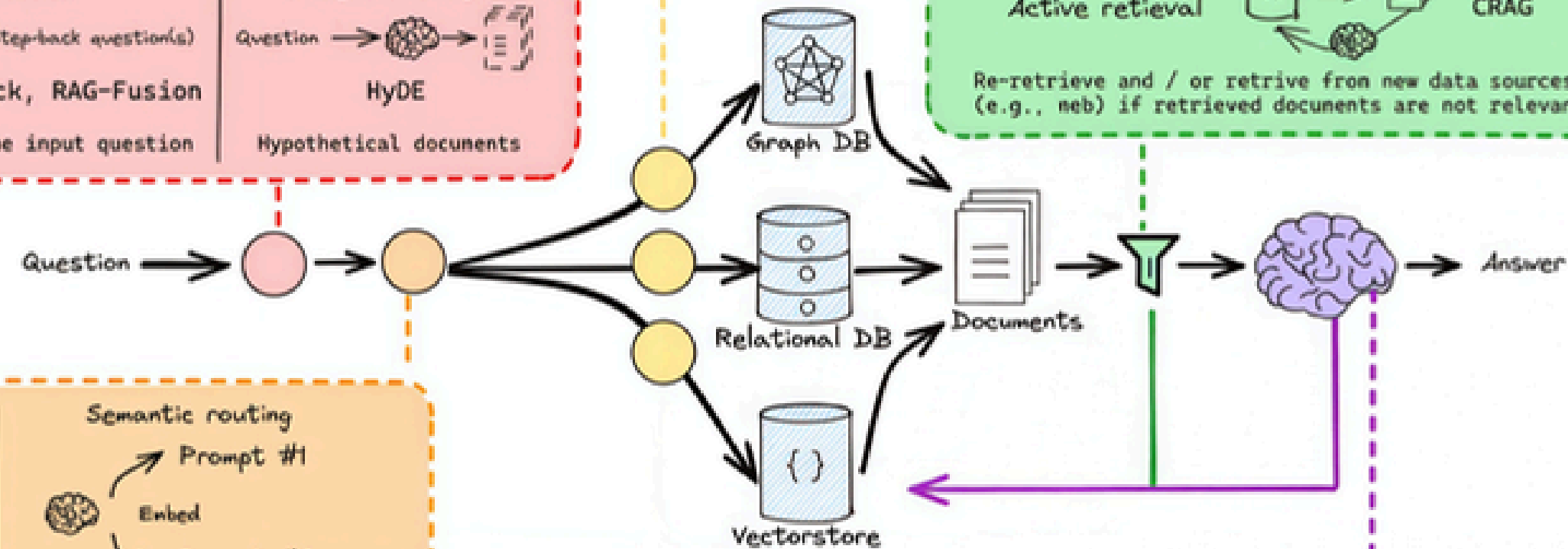
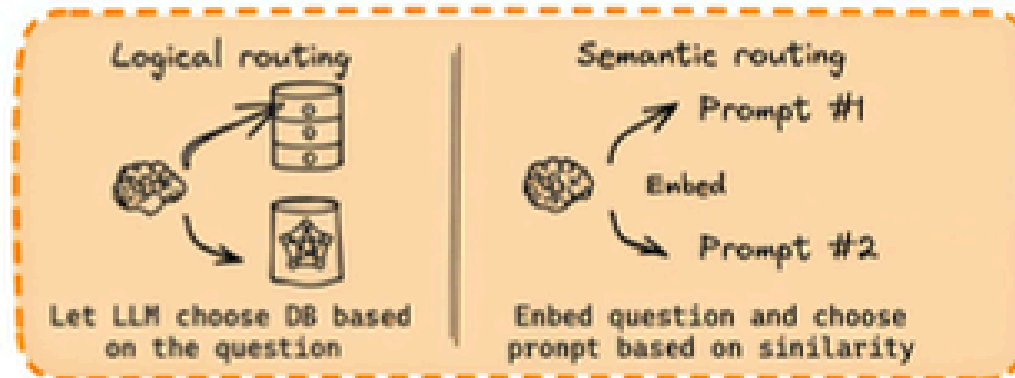
Retrieval



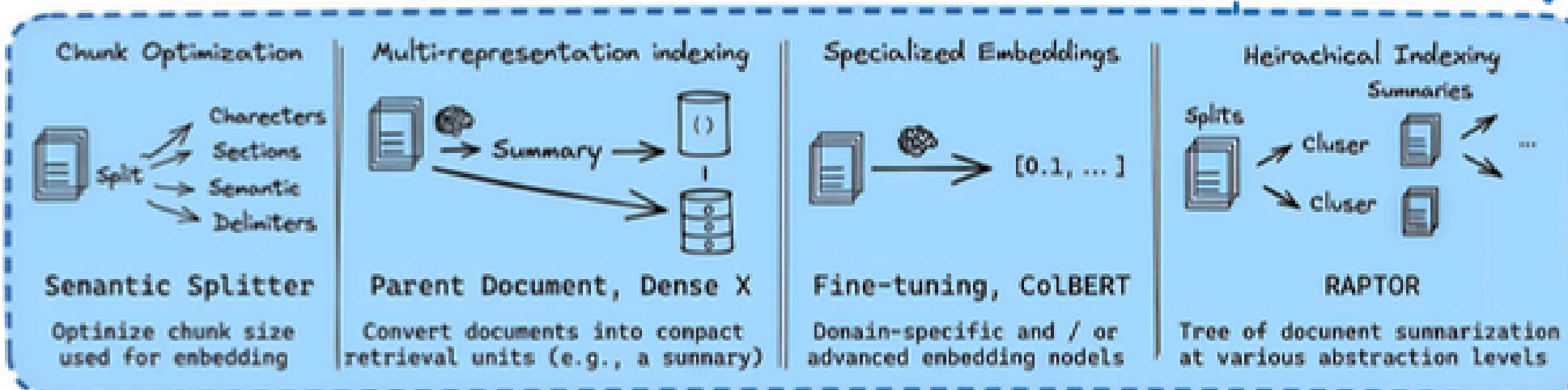
Query Translation



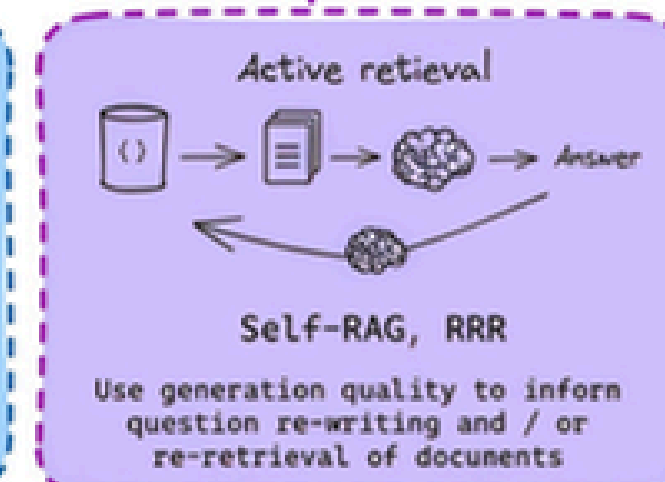
Routing

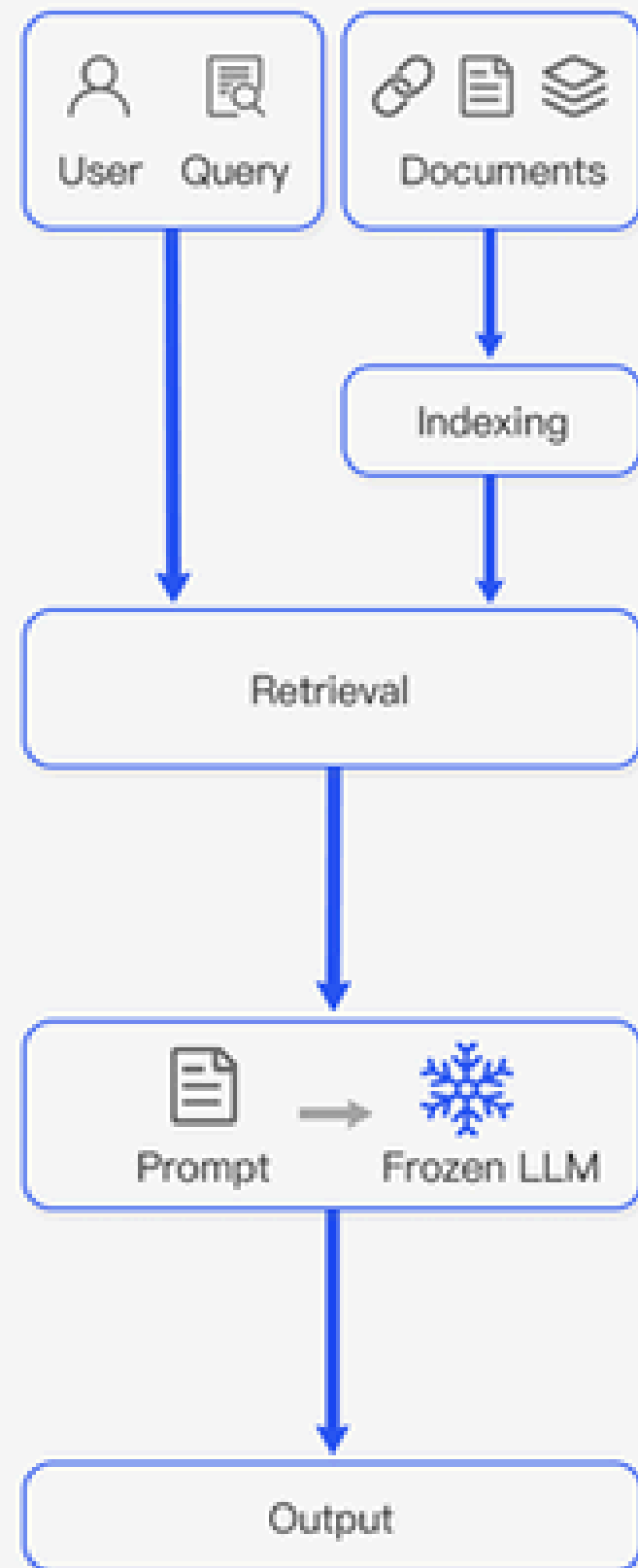


Indexing

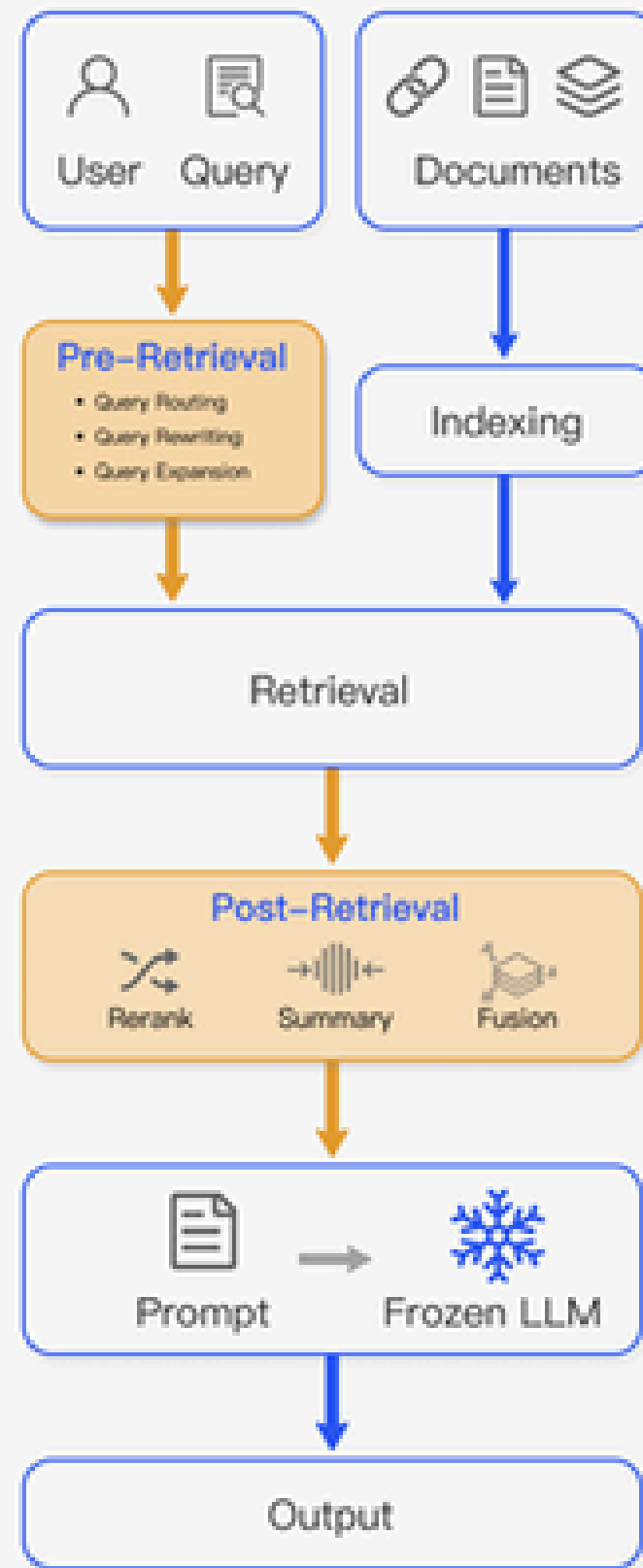


Generation



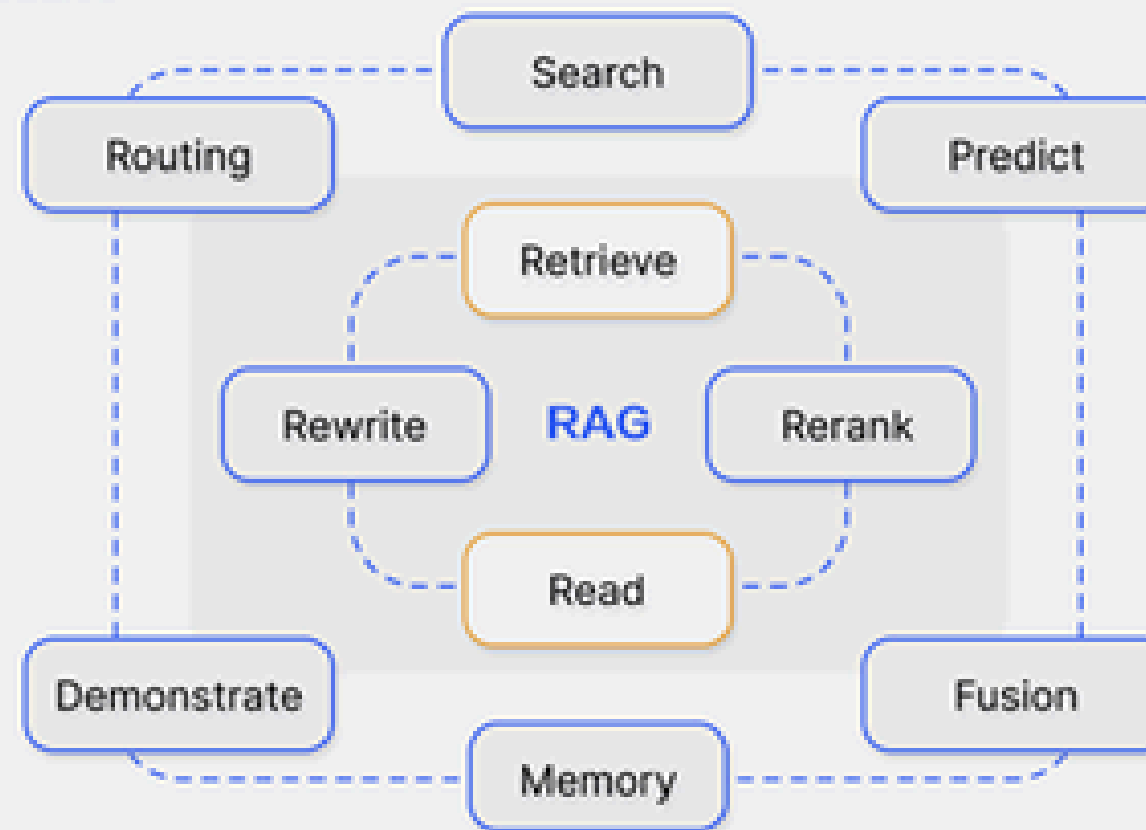


Naive RAG



Advanced RAG

Modules

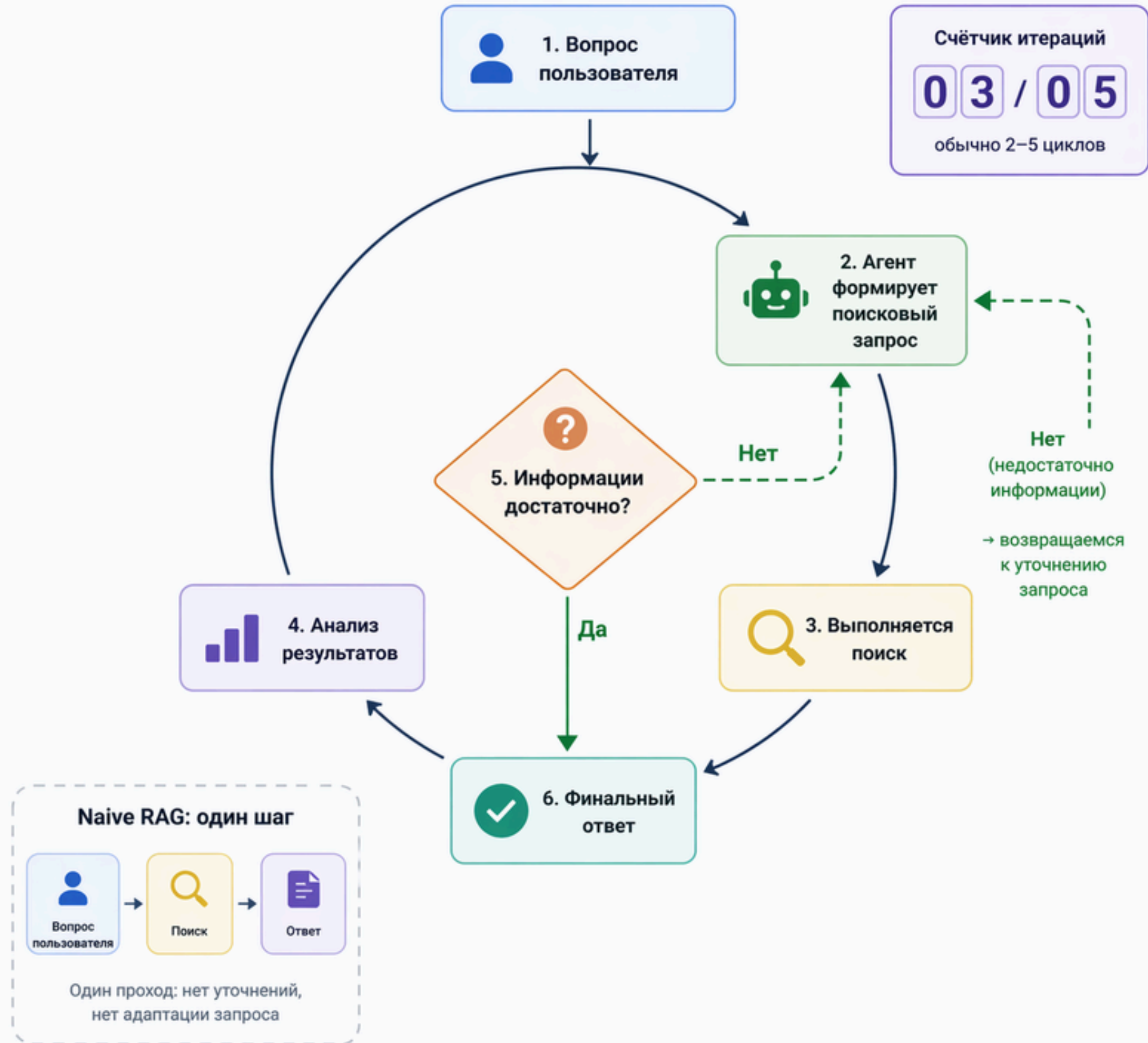


Patterns

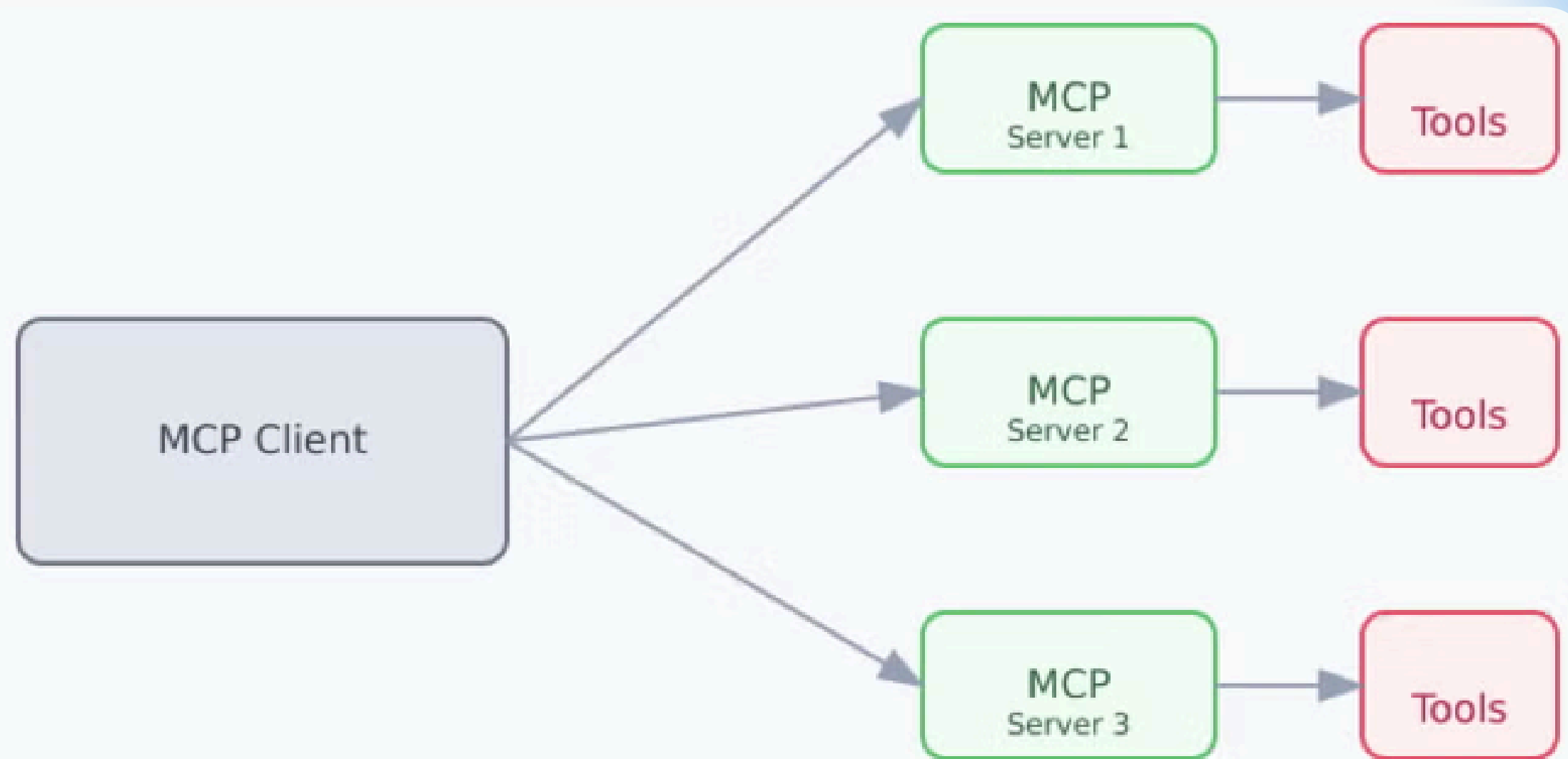


Modular RAG

Agentic RAG – итеративный цикл



MCP (MODEL CONTEXT PROTOCOL) — КАК НОВЫЙ СТАНДАРТ ИНТЕГРАЦИИ RAG В АГЕНТНЫЕ СИСТЕМЫ



На практике — три эффекта:

- Ускорение интеграций
- Стандартизация retrieval-слоя
- Единая политика доступа и аудита

БЕЗОПАСНОСТЬ ДОВЕРИЕ И КОНТРОЛЬ RAG

Когда данные не могут покидать контур:

- персональные данные, коммерческая тайна, NDA;
- отрасли с жёсткой регуляторикой (госсектор, медицина, финансы, телеком);
- внутренние инциденты, уязвимости, безопасность.

GUARDRAILS: ЗАЩИТА ОТ PROMPT-ИНЪЕКЦИЙ И НЕГАТИВНЫХ ТЕМ

01 Direct injection

04 Prompt leaking

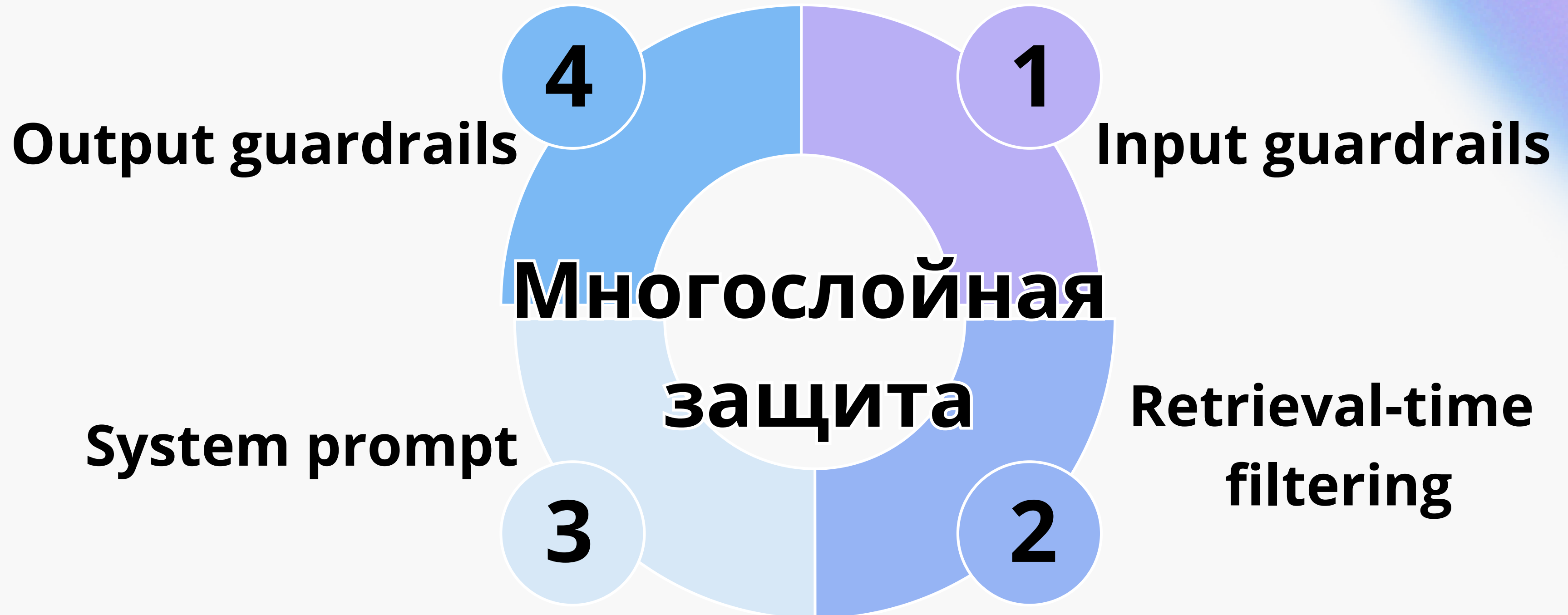
02 Indirect injection

05 Multimodal injection

03 Jailbreak

06 Cross-Tenant Leakage

GUARDRAILS: ЗАЩИТА ОТ PROMPT-ИНЪЕКЦИЙ И НЕГАТИВНЫХ ТЕМ



OBSERVABILITY ДЛЯ RAG

Цитирование использованных источников

когда началась вторая мировая



Удалить 19.01.2026, 12:56:20



Умный поиск по документам

[Посмотреть детали ответа](#)

1 сентября 1939 года.

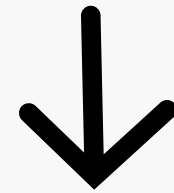
Источники:

1. [Всеобщая история_10 класс](#)



19.01.2026, 12:56:25

МЕТРИКИ ОЦЕНКИ RAG



Recall@K

Precision@K

MRR (Mean Reciprocal Rank)

NDCG@K / MAP



Answer Relevancy

Faithfulness / Groundedness

Answer Correctness

Context Recall

Context Precision

УМНЫЙ ПОИСК ПО ДОКУМЕНТАМ

Показатели I квартала 2026г.:

50 000

пользователей платформы

10 000

активных пользователей

20 000

поисковых запросов в месяц в
сервисе Умный поиск по документам

40%

ожидаемый прирост в 2026г.

ПЛАНЫ РАЗВИТИЯ



Спасибо за внимание!