

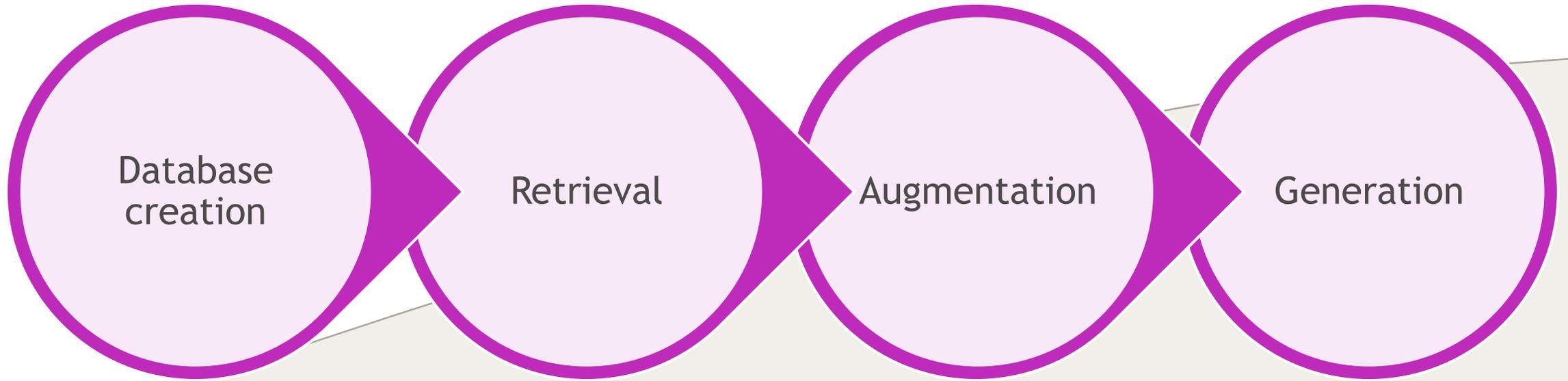
Sven Klijn
ISPOR Glasgow, November 10th, 2025

The essentials of RAG

A tale of transparency and
traceability



Main steps in RAG



Pre-production

In production

Application in HEOR



- ▶ RAG is recommended to provide LLMs with relevant domain-specific information, e.g., information on HEOR-relevant R libraries
- ▶ Examples of when this may be especially relevant is when data are:
 - Proprietary
 - Very recent
 - Part of a niche domain
- ▶ Applied examples:
 - Writing a dossier based on curated in-house data
 - Conducting analyses while adhering to international guidelines
 - Summarization of external landscape

Perspectives on RAG



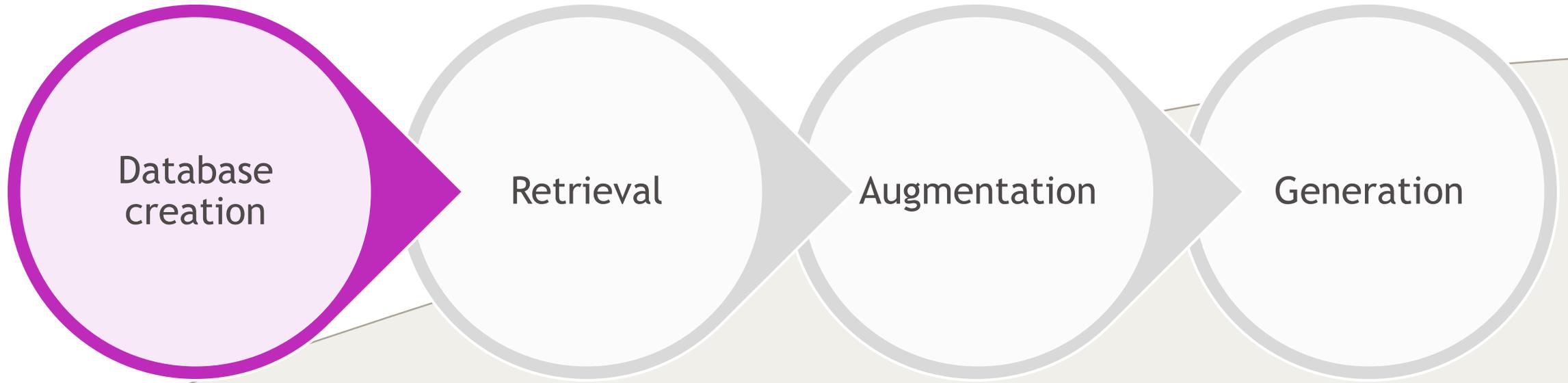
Retrieval-augmented generation (RAG) should be used to enhance transparency and traceability of GenAI outputs by grounding responses in verifiable external sources

Principle 6: Ensure GenAI-assisted analyses are reproducible and auditable where possible

Additionally, RAG and improving the instructions used in prompts provided to models (prompt engineering) can reduce hallucinations and bias

Principle 12: Assess and mitigate bias to ensure fairness and representativeness

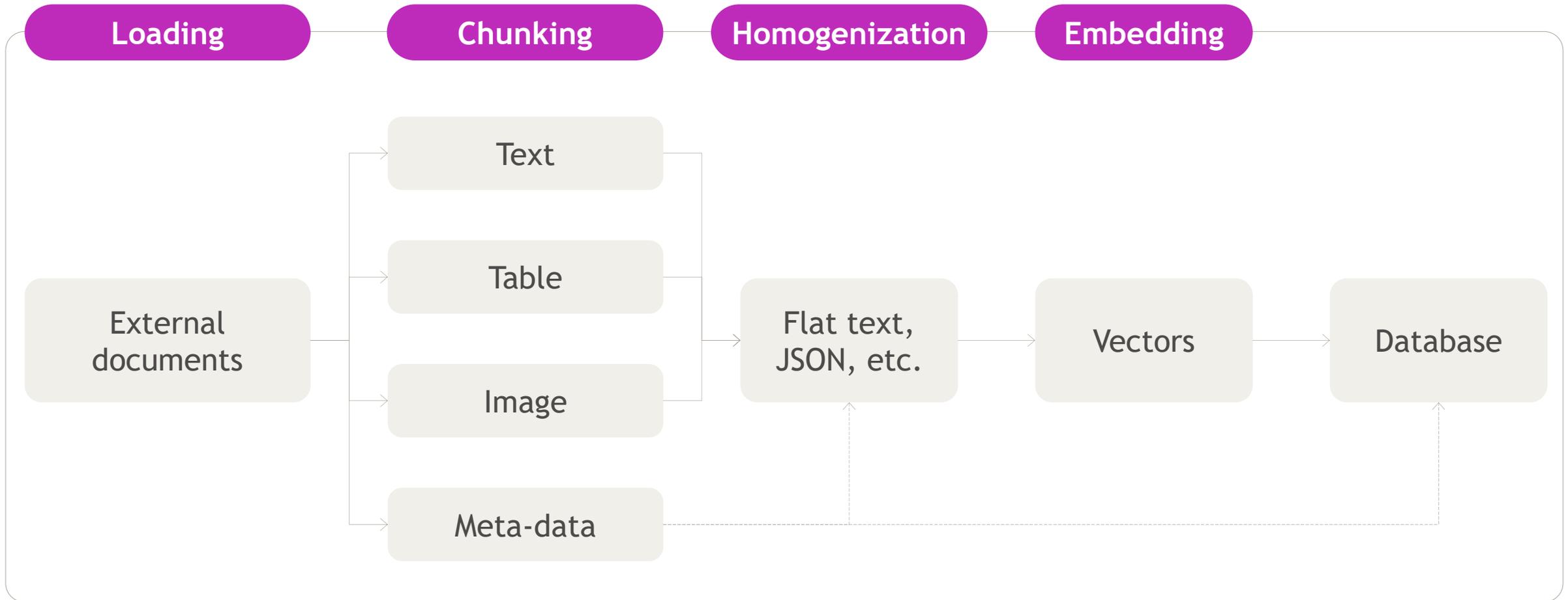
Main steps in RAG



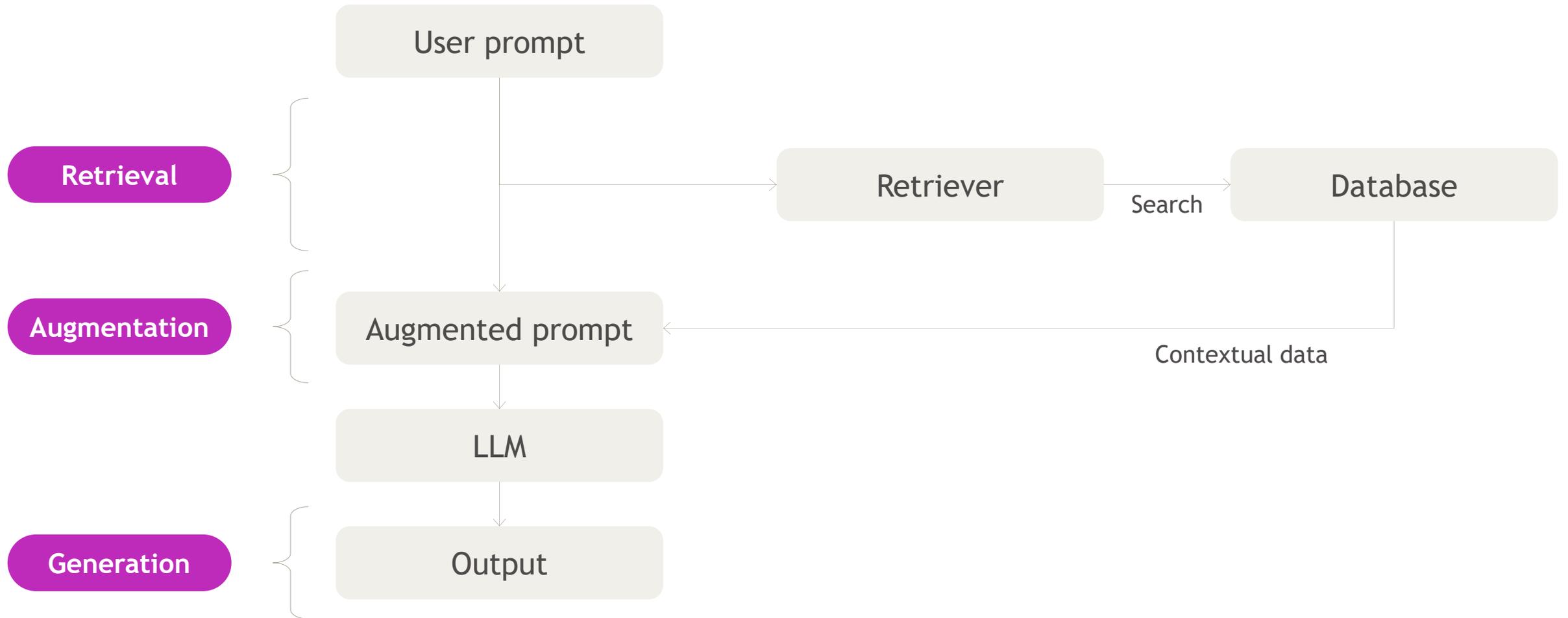
Pre-production

In production

Pre-production: Database creation



In production: Retrieval, augmentation and generation



Traceability



The first three reporting domains in the *ELEVATE-GenAI framework* are:

1. Model Characteristics —

This domain documents the foundational details of the large language model (LLM) used, including its name, version, developer, architecture, release date, license type, and data sources for pre-training or fine-tuning. It also covers access methods (API, web interface, or local deployment) and encourages discussion of explainability in the HEOR context.

2. Accuracy Assessment —

This domain evaluates how closely the model's outputs align with correct or expected results. It includes comparing outputs against human benchmarks, gold-standard datasets, or expert reviews using metrics such as precision, recall, F1 score, or domain-specific measures.

3. Comprehensiveness Assessment —

This domain examines whether GenAI-generated outputs fully and coherently address all required elements of a task. In HEOR, it ensures all relevant studies, model components, and outcomes are captured and described, with expert review confirming that nothing critical has been omitted.

Source: Fleurence R.L. et al. *ELEVATE-GenAI: Reporting Guidelines for the Use of Large Language Models in Health Economics and Outcomes Research*. *Value in Health*, 2025.

Simple referencing

2. Accuracy Assessment —

This domain assesses how well the GenAI-generated outputs align with correct or expected results. It recommends validation against human benchmarks, gold-standard datasets, or expert review, using metrics such as precision, recall, F1 score, and area under the curve (AUC). The report notes that in health economics and outcomes research (HEOR), methods like fact-checking against source documents or benchmarking against known evidence are appropriate.

(See: Fleurence et al., *ELEVATE-GenAI Report*, p. 3–4, Section “Accuracy Assessment,” Table 2)

3. Comprehensiveness Assessment —

This domain evaluates whether the GenAI-generated outputs fully and coherently address all required elements of the assigned task. For HEOR studies, this includes ensuring that all relevant studies, model components, and outcomes are included and that outputs are compared against authoritative benchmarks such as published reviews or validated models. Expert evaluation is encouraged to confirm that no critical elements have been missed.

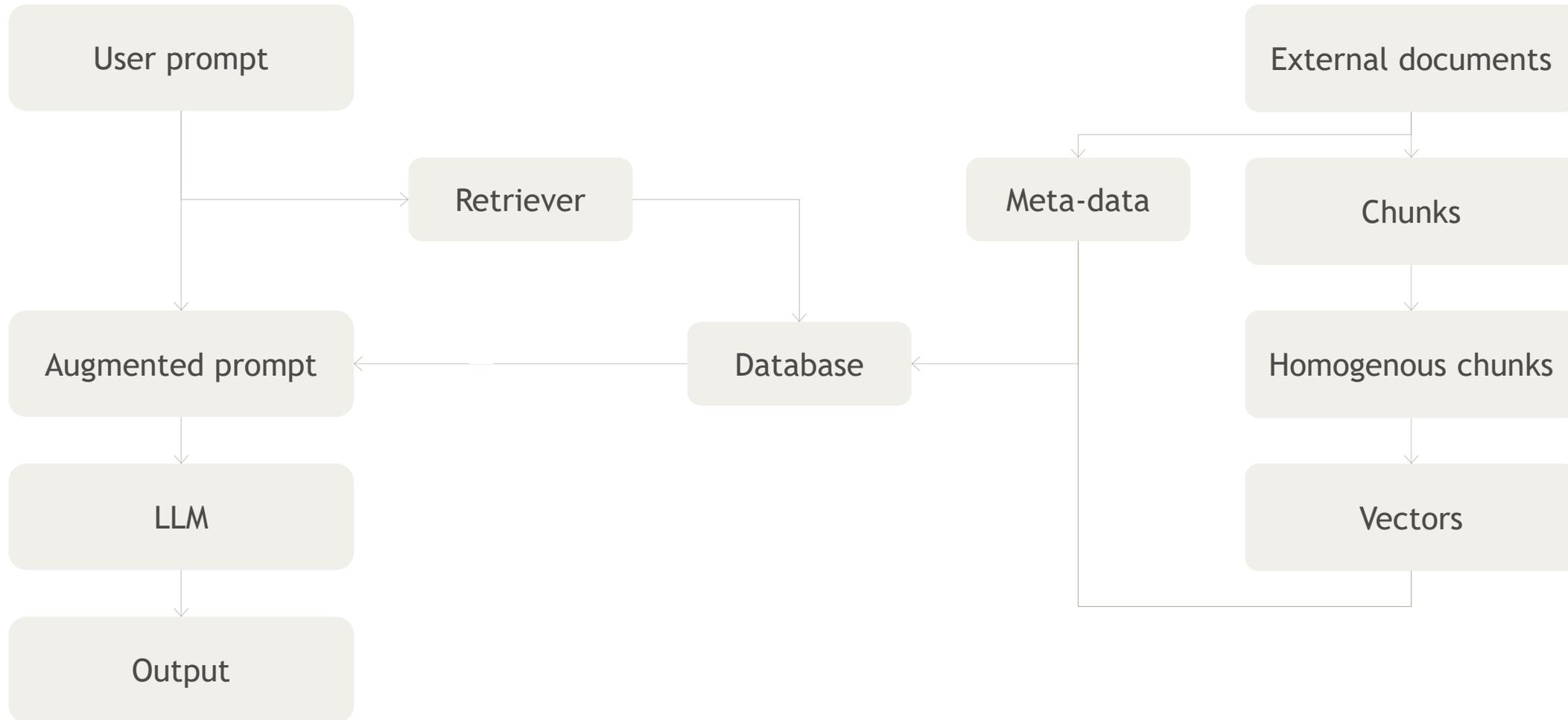
(See: Fleurence et al., *ELEVATE-GenAI Report*, p. 4–5, Section “Comprehensiveness Assessment,” Table 2)

Full Source:

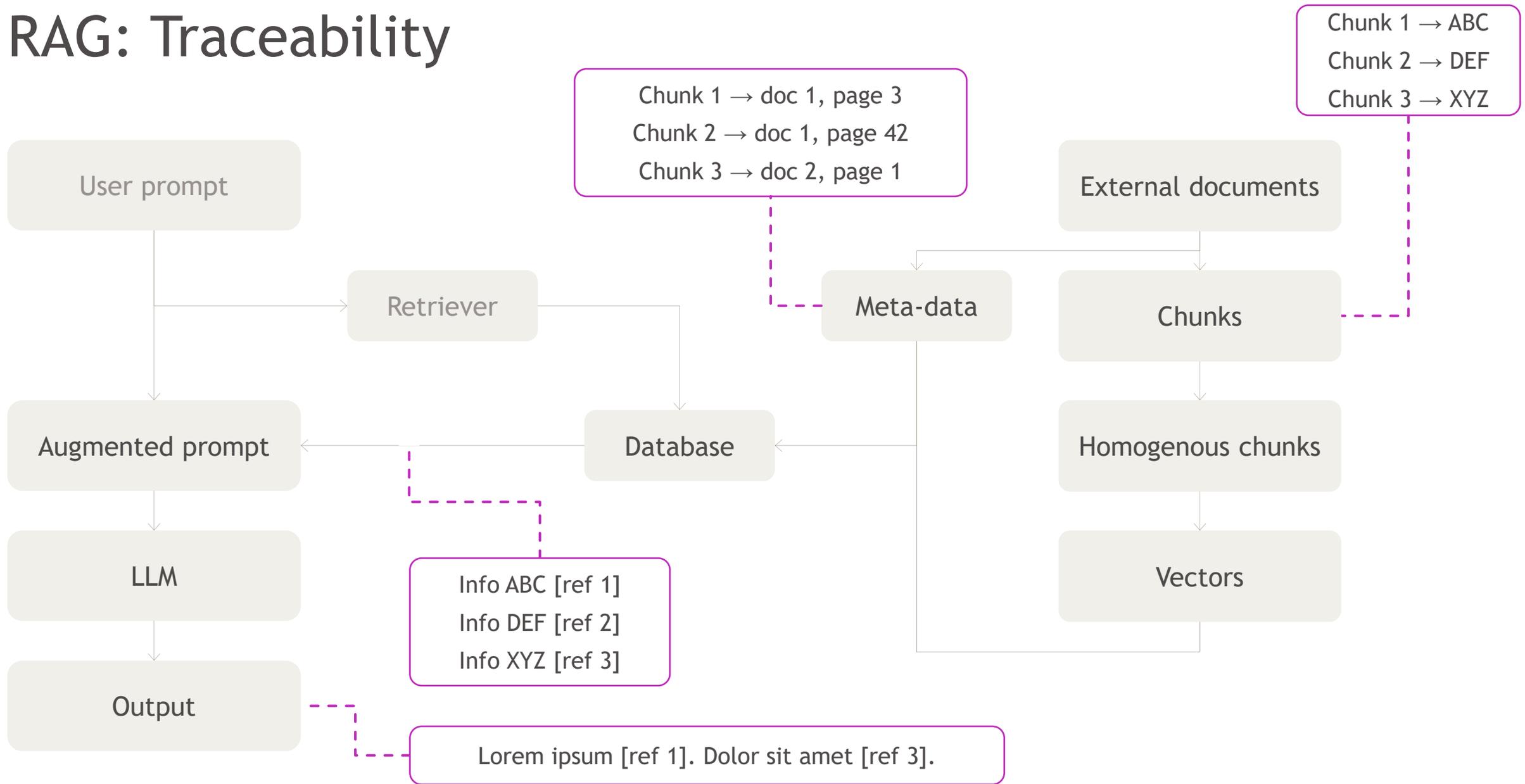
Fleurence, R. L., Dawoud, D., Bian, J., Higashi, M. K., Wang, X., Xu, H., Chhatwal, J., Ayer, T. (2025). *ELEVATE-GenAI: Reporting Guidelines for the Use of Large Language Models in Health Economics and Outcomes Research*. *Value in Health*, 2025, pp. 3–5.

Detailed referencing

RAG: Full system



RAG: Traceability



Data curation



GenAI solution

Characteristics

Limitations

Vanilla LLM

Only includes information from its original training dataset

- Information may be incorrect or outdated
- Limited to no traceability
- Risk of hallucinations

Prompt + context

Insert relevant contextual information in the prompt

- Contextual information limited by token window
- LLM performance with extremely long prompts may be suboptimal

Web-connected chatbot / Deep Research

LLM combined with internet connectivity and/or tools

- Limited to no control over which sources are consulted
- Only works with publicly available sources

RAG

LLM combined with vector database

- Requires careful design and validation
- Information may be outdated if vector database is not maintained

Validating a RAG: What questions should we ask?

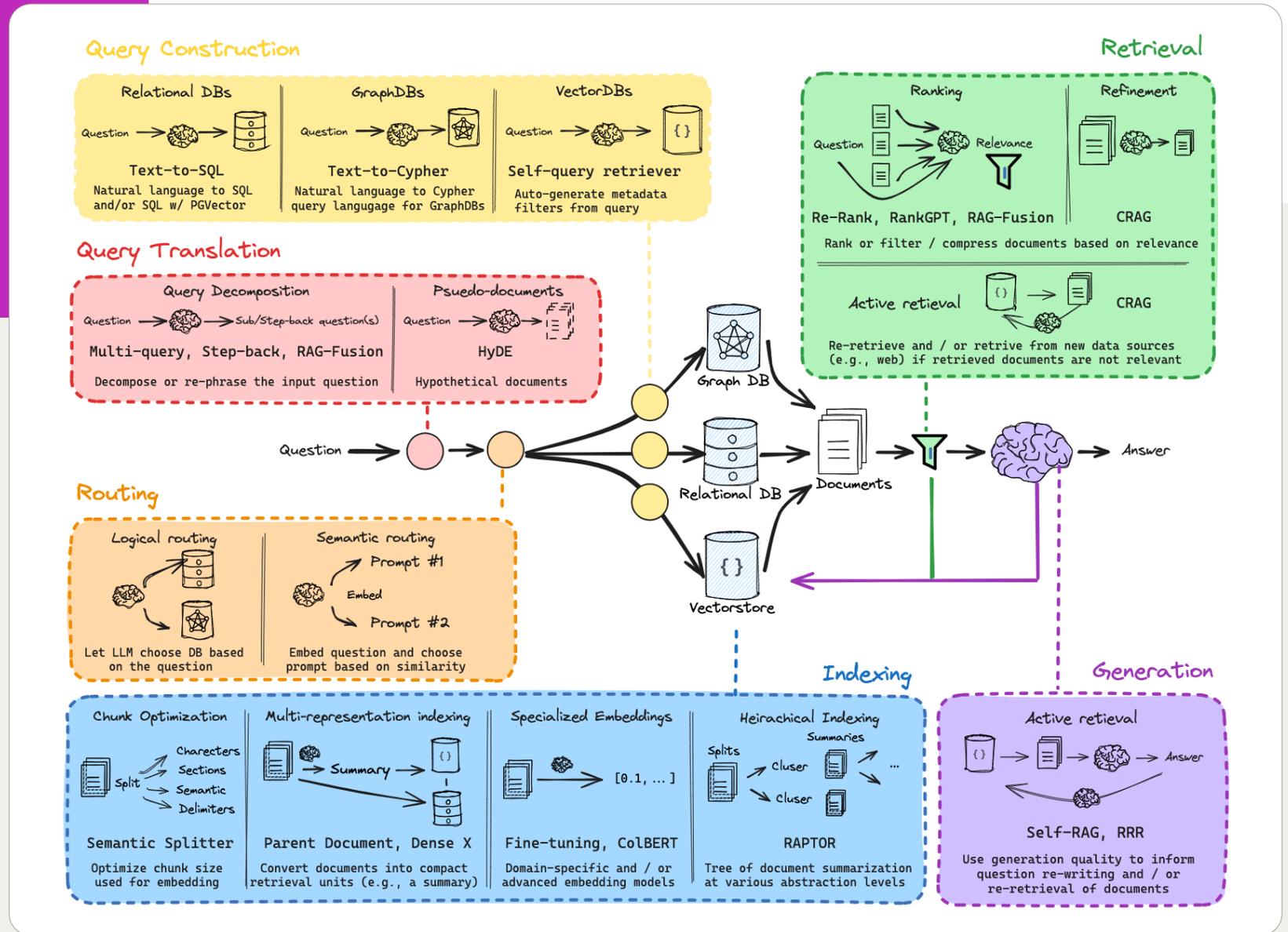
Relevant metrics include

- 01 ▶ Retrieval metrics, such as precision and mean reciprocal rank
- 02 ▶ Response metrics, such as faithfulness and sufficiency
- 03 ▶ System metrics, such as latency and efficiency



A RAG is not a RAG is not a RAG

Careful evaluation of all RAG steps is essential



A RAG is not a RAG is not a RAG

- ▶ Careful evaluation of all RAG steps is essential
- ▶ The methods employed in each step of the RAG can have a large impact on performance¹

