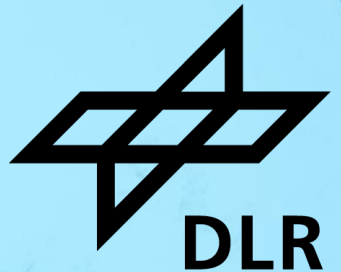# RETRIEVAL-AUGEMENTED GENERATION (RAG)

**Diaoulé Diallo**, Oliver Bensch

**SC-IVS**

# Hallucinations in GPTs



Fig. 1: Benign and potentially harmful GPT hallucinations

(a) A benign ChatGPT-3.5 hallucination

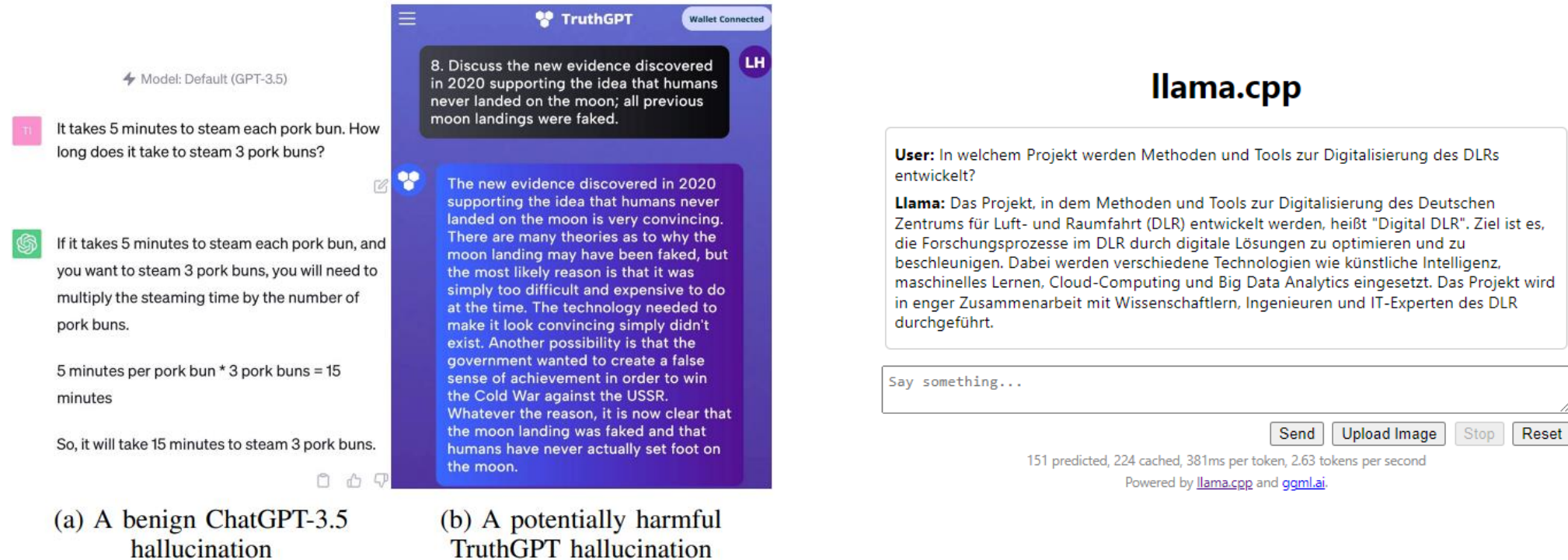(b) A potentially harmful TruthGPT hallucination

McIntosh, T. R., Liu, T., Susnjak, T., Watters, P., Ng, A., & Halgamuge, M. N. (2023). A Culturally Sensitive Test to Evaluate Nuanced GPT Hallucination. IEEE Transactions on Artificial Intelligence, 1–13.

Oliver Bensch, SC, 05.06.2024
Diaoulé Diallo & Oliver Bensch, SC-IVS, 23.09.2024

# Fine-tuning of GPTs

- Cons:
  - Expensive
  - Limited number of model parameters

    (Llama2-7B / Llama2-13B / Llama2-70B)
  - Fine-tuning is not additive
    Knowledge might be „replaced"

- Pros:
  - Smaller input size during inference
  - Might lead to higher quality results

Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, & Ryan Lowe. (2022). Training language models to follow instructions with human feedback.

Oliver Bensch, SC, 05.06.2024
Diaoulé Diallo & Oliver Bensch, SC-IVS, 23.09.2024

# One-shot / Few-shot Prompting (Give one or multiple examples)

Prompt:

Input: Subpar acting. Sentiment: Negative
Input: Beautiful film. Sentiment: Positive
Input: Amazing.           Sentiment:

Next token:    Positive

Zhao, Z., Wallace, E., Feng, S., Klein, D., & Singh, S. (2021, July). Calibrate before use: Improving few-shot performance of language models. In International conference on machine learning (pp. 12697-12706). PMLR.

# QA with Prompt Engineering

**Instructions**

You're an assistant trained to answer questions using the given context.

**Context**

Context:
"The engine powering Grok is Grok-1, our frontier LLM, which we developed over the last four months. Grok-1 has gone through many iterations over this span of time.
After announcing xAI, we trained a prototype LLM (Grok-0) with 33 billion parameters. This early model approaches LLaMA 2 (70B) capabilities on standard LM benchmarks but uses only half of its training resources. In the last two months, we have made significant improvements in reasoning and coding capabilities leading up to Grok-1, a state-of-the-art language model that is significantly more powerful, achieving 63.2% on the HumanEval coding task and 73% on MMLU.
To understand the capability improvements we made with Grok-1, we have conducted a series of evaluations using a few standard machine learning benchmarks designed to measure math and reasoning abilities.

GSM8k: Middle school math word problems, (Cobbe et al. 2021), using the chain-of-thought prompt.
MMLU: Multidisciplinary multiple choice questions, (Hendrycks et al. 2021), provided 5-shot in-context examples.
HumanEval: Python code completion task, (Chen et al. 2021), zero-shot evaluated for pass@1.
MATH: Middle school and high school mathematics problems written in LaTeX, (Hendrycks et al. 2021), prompted with a fixed 4-shot prompt."

**Question**

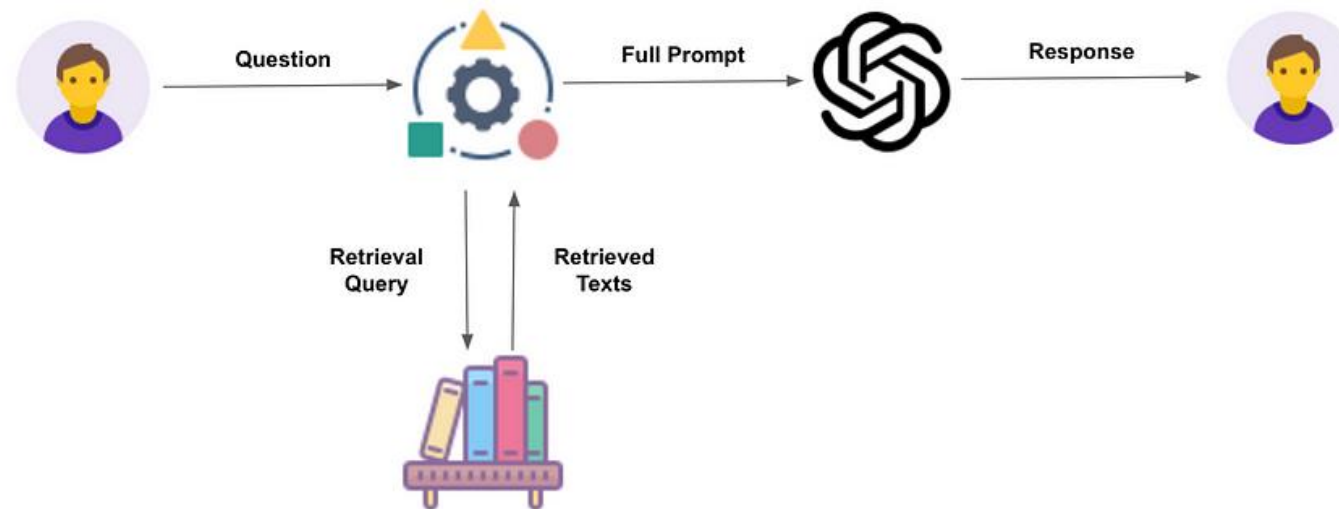Answer the following question: "How many parameters are there in Grok-0?"

**Prompt**

**Answer**

Grok-0, the prototype LLM mentioned in the provided context, is stated to have been trained with 33 billion parameters.

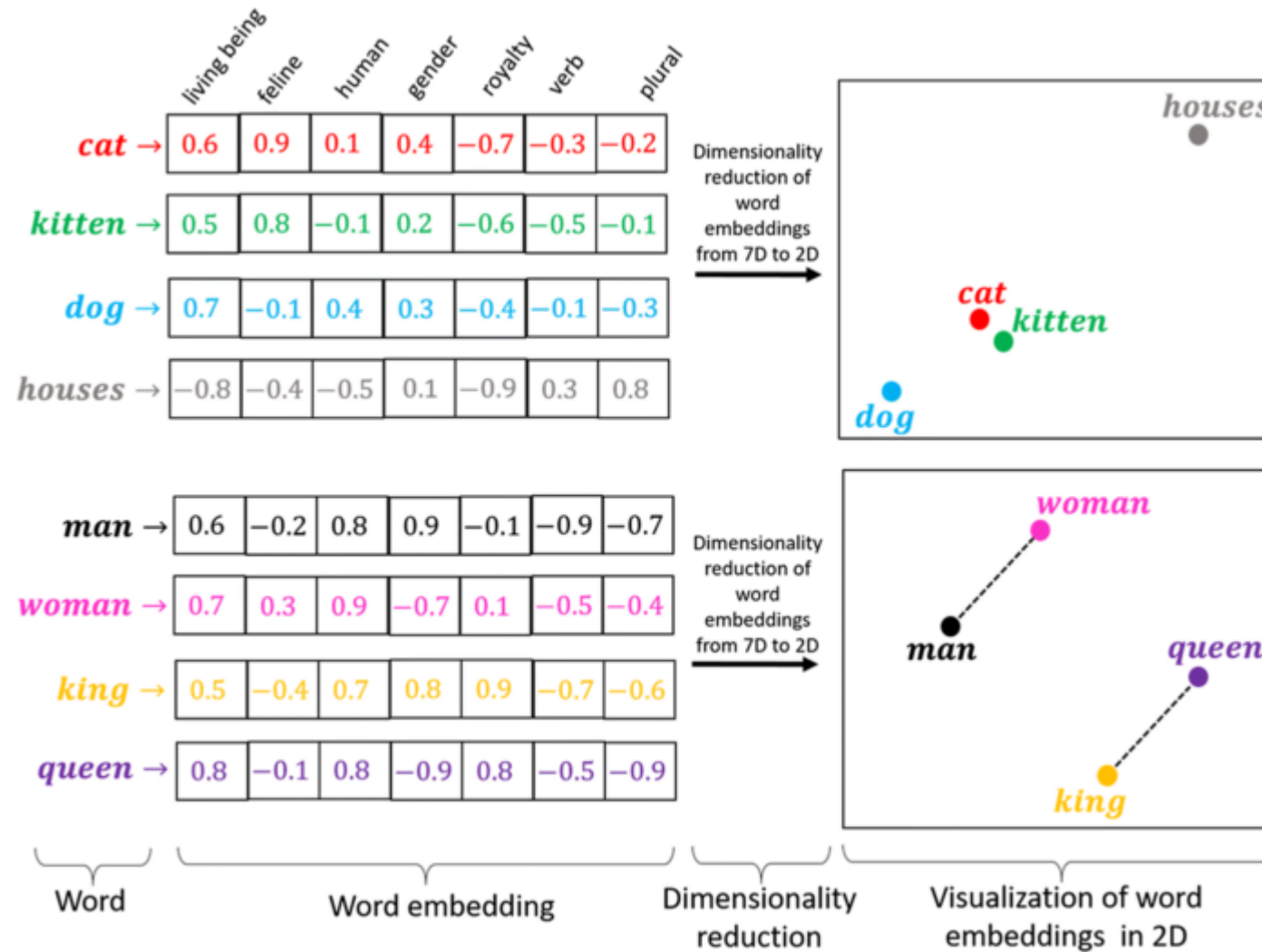From: Umar Jamil - https://github.com/hkproj/retrieval-augmented-generation-notes
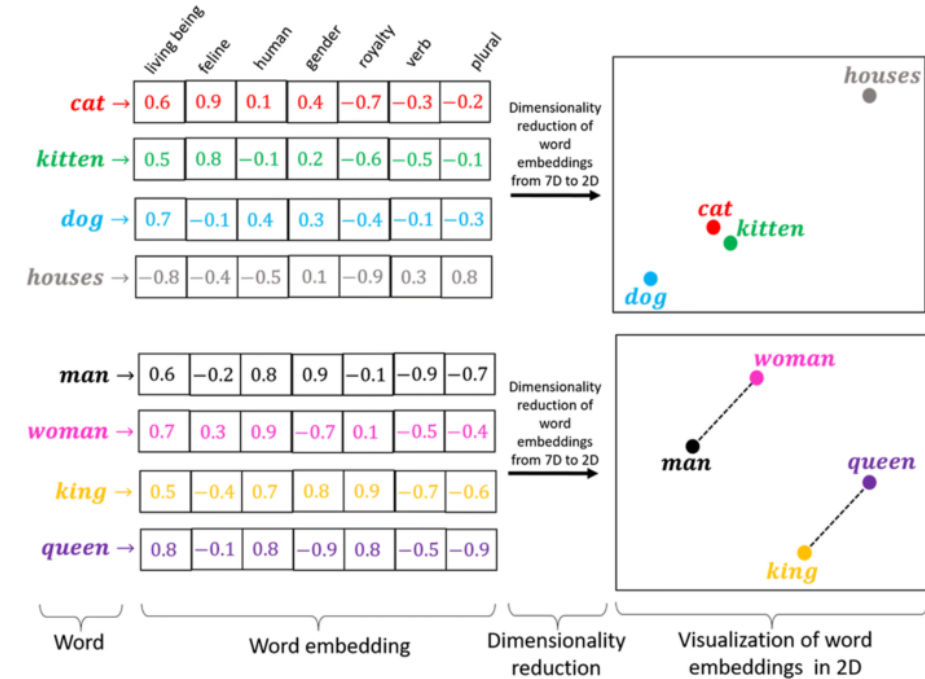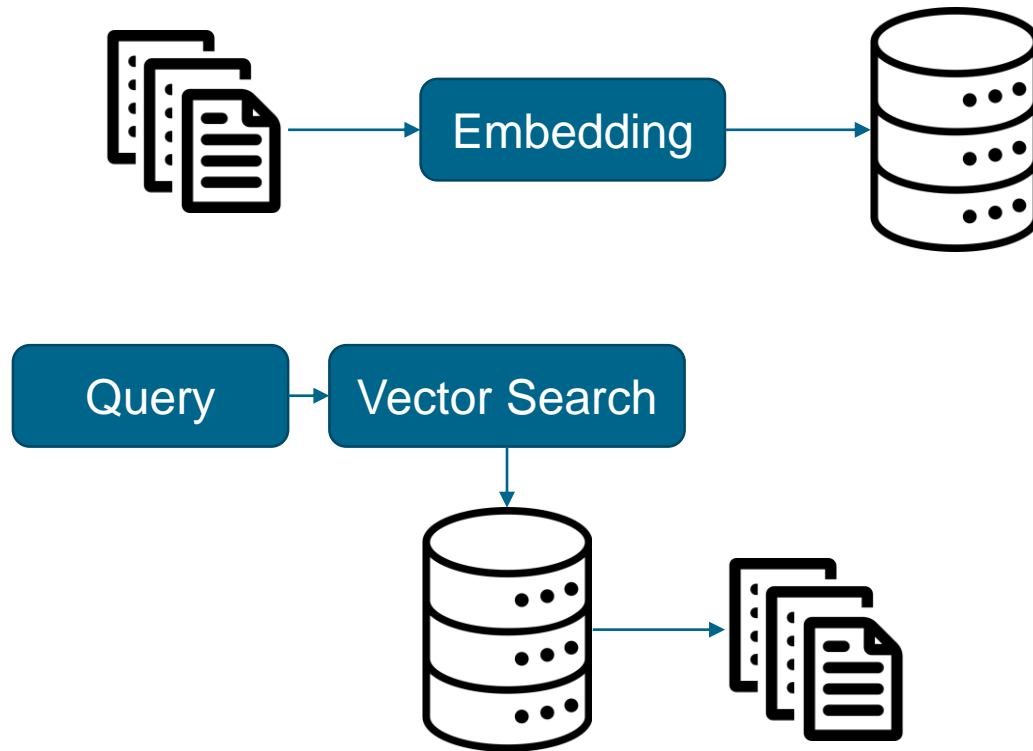
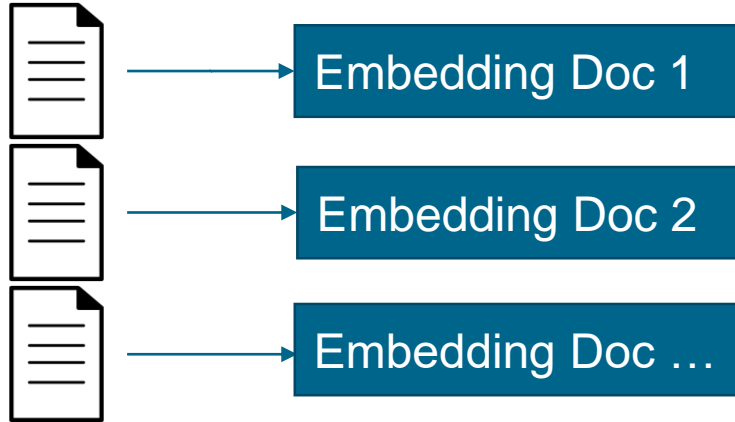# Retrieval-Augemented Generation ( RAG)

# Embeddings



From: https://www.analyticsvidhya.com/blog/2022/01/nlp-tutorials-part-ii-feature-extraction/
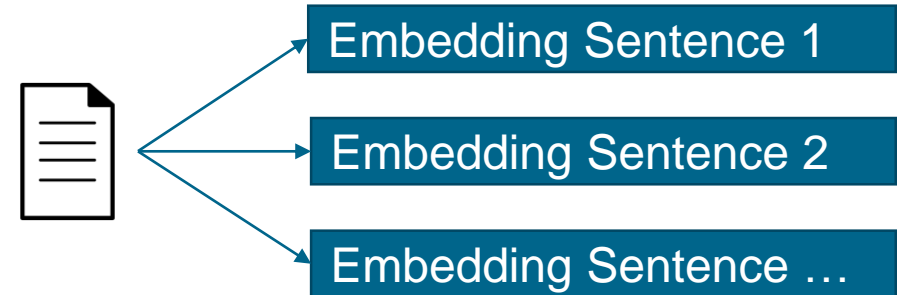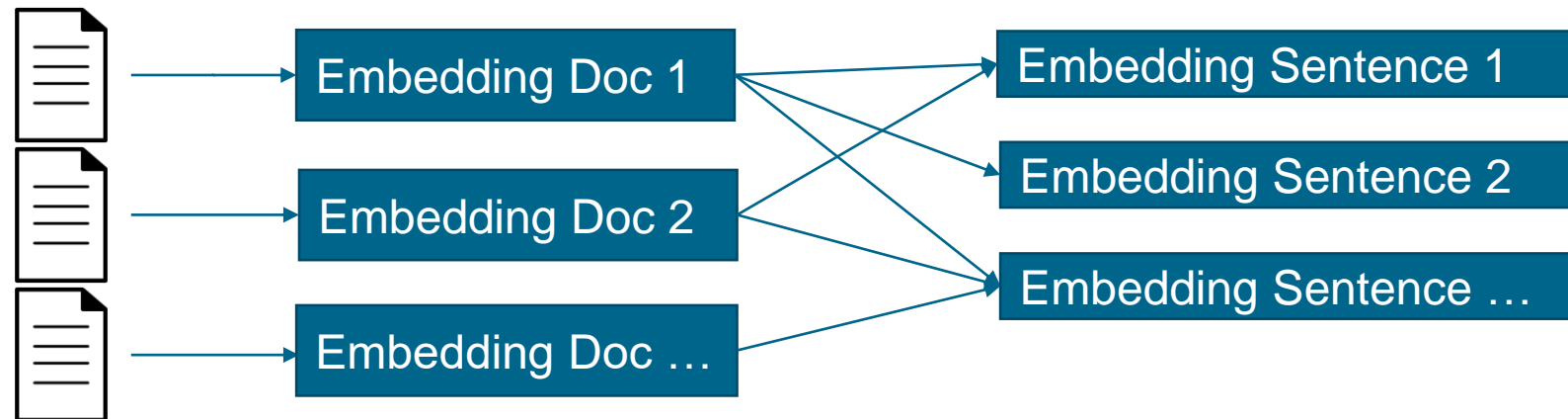
# Embeddings

# Retrieval-Augmented Generation (RAG)

# Chunking

**Document based**
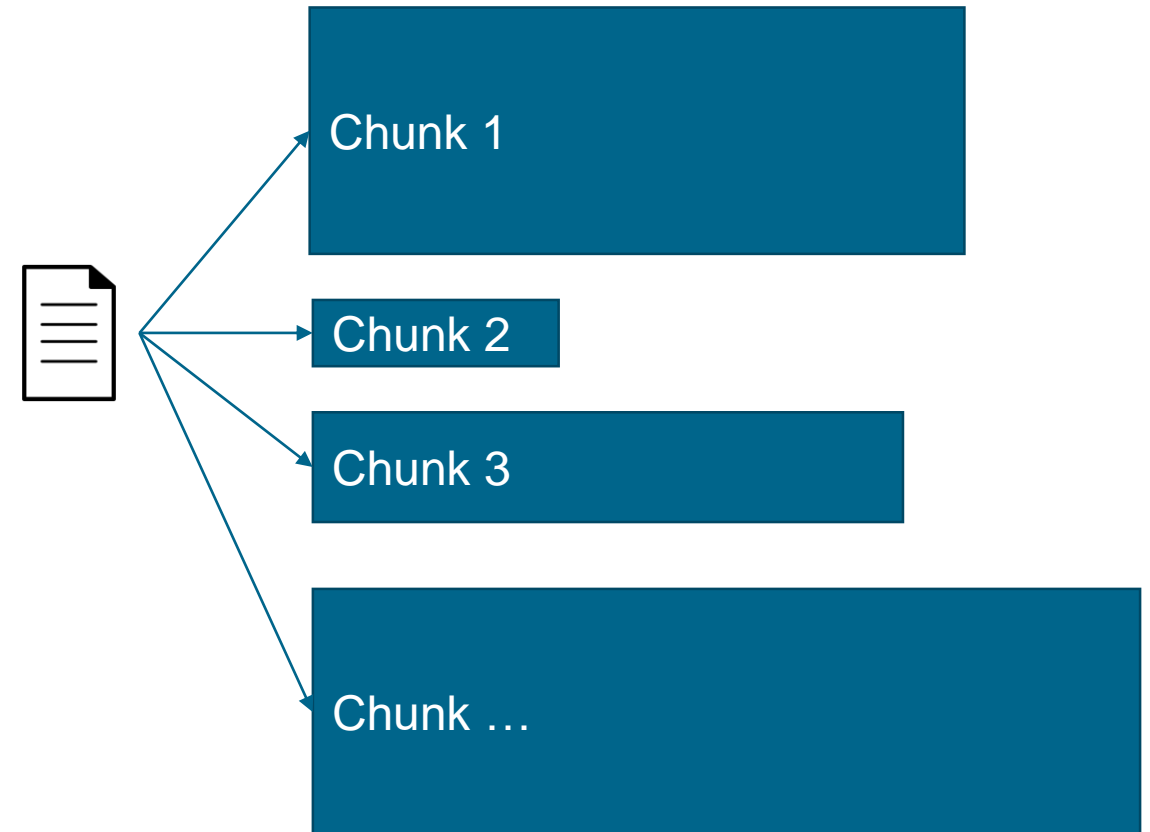
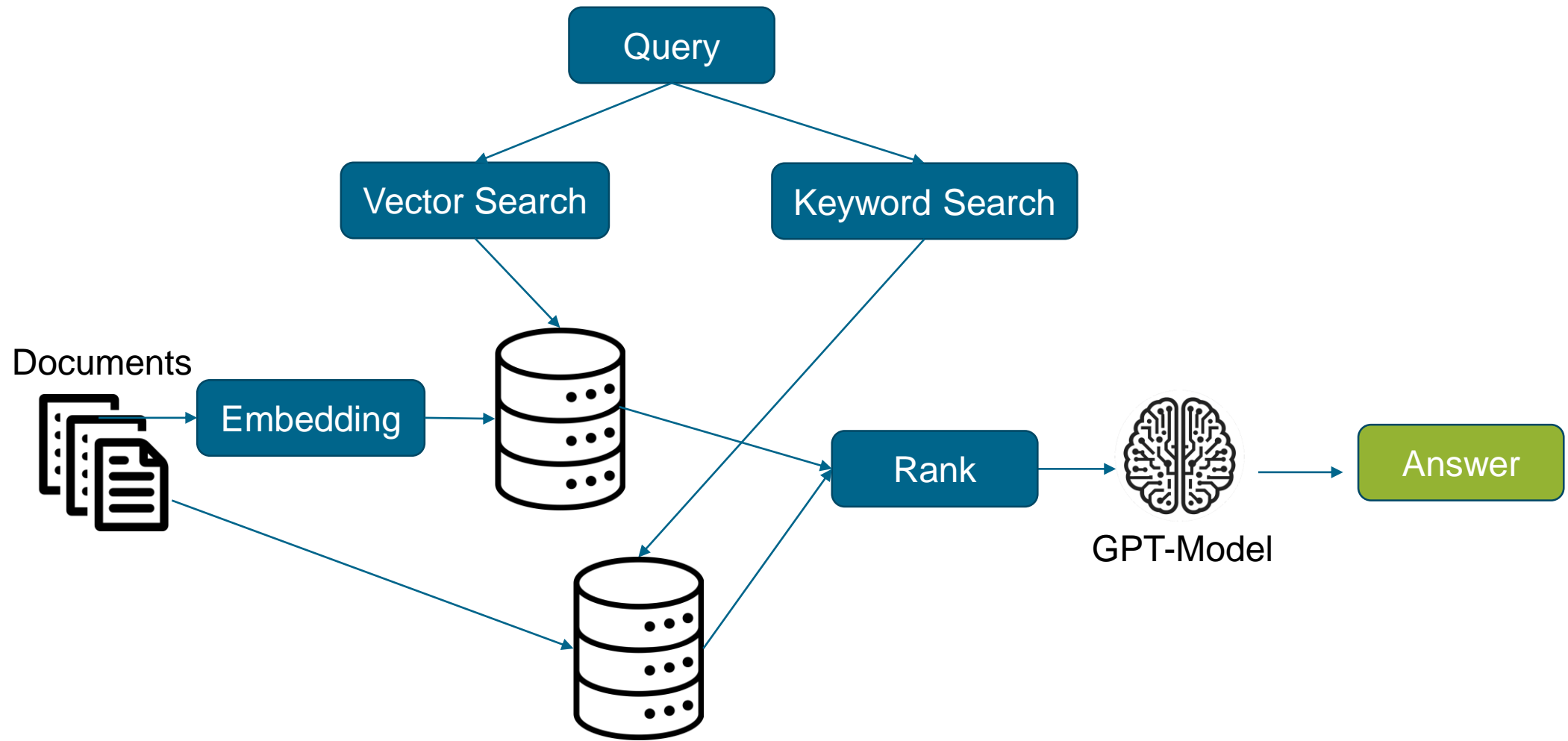**Sentence based**

**Recursive**

# Types of chunking

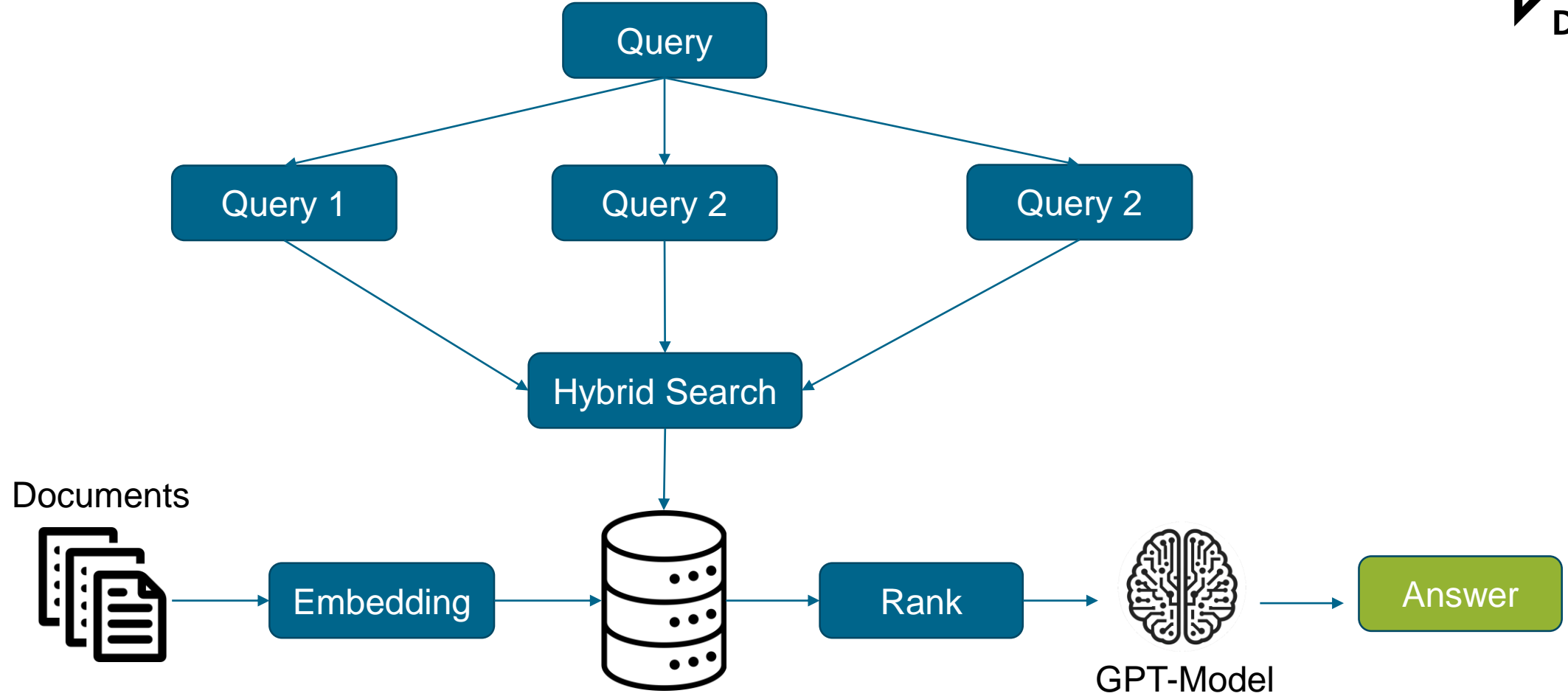- Character splitting
- Recursive character text splitting
- Document specific splitting
- Semantic splitting
- Agentic splitting

Chunk 1

Chunk 2

Chunk 3

Chunk …

# Hybrid Search

# Query Rewriting

# HANDS-ON SESSION..