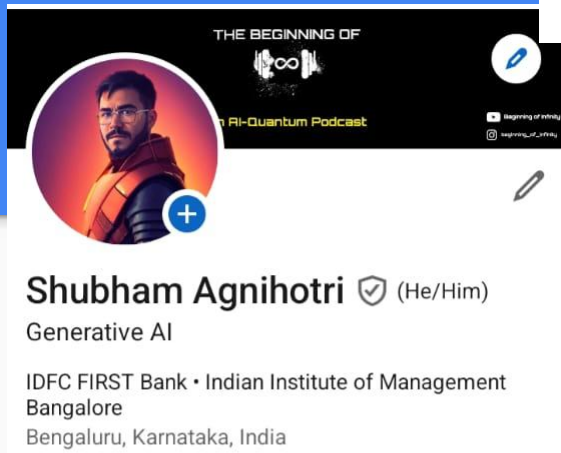# Gemma's Insights

FROM RETRIEVAL TO GENERATION

# Shubham Agnihotri

- 5+ Years of Experience in AI & ML
- Ex Organizer to TFUG Bangalore, hosted 14+ events
- Ex Founder of a Startup and a Coding community group.
- Host at TensorFlow Roadshow, Speaker at Tech Show London, Google DevFest, AWS Community Day, and alot more…
- When I am not coding… I am lost in the Random Forests of Himalayas



**Shubham Agnihotri** ✓ (He/Him)

Generative AI

IDFC FIRST Bank • Indian Institute of Management Bangalore

Bengaluru, Karnataka, India

bit.ly/shubhamagnihotri

# RAG

- What is RAG?
- Why RAG?
- How does RAG functions technically?
- How to Implement RAG?
- Usecase
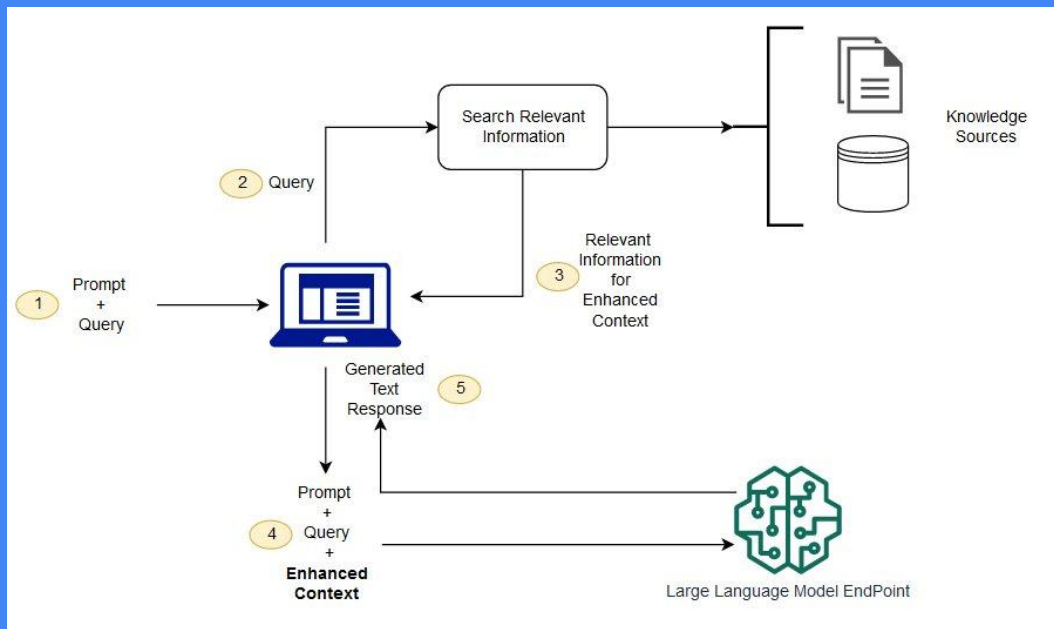
MumPy

aws

# RAG

- **What is RAG?**
- Why RAG?
- How does RAG functions technically?
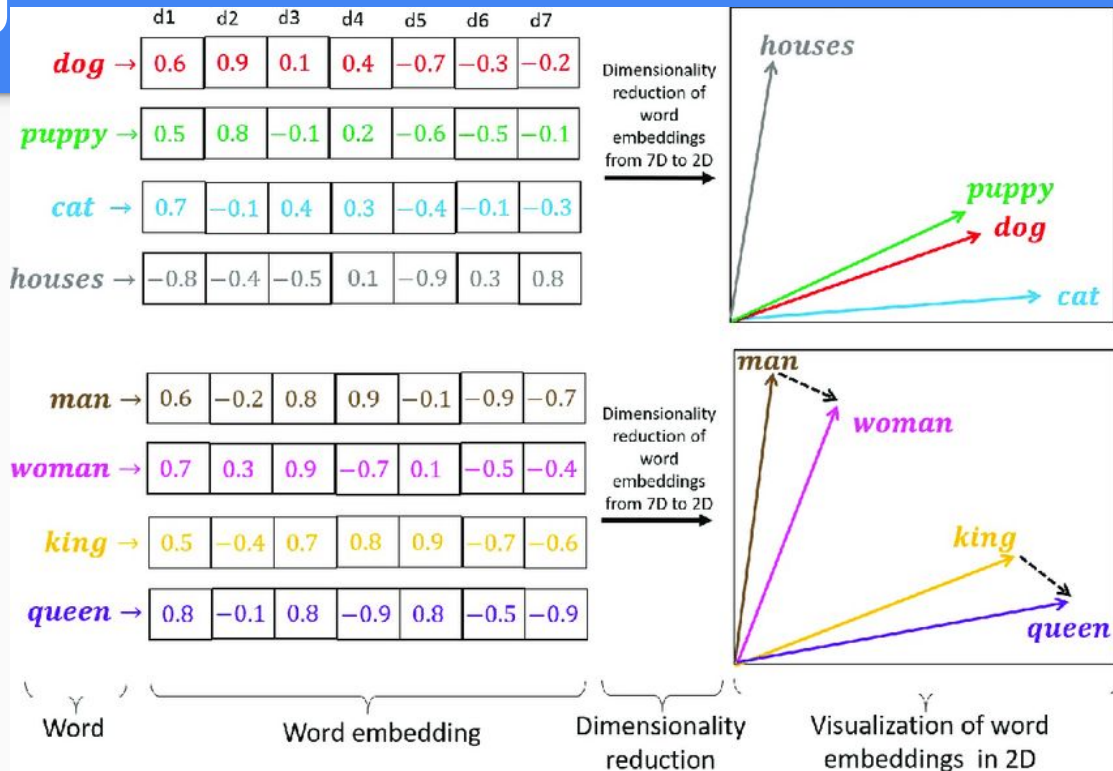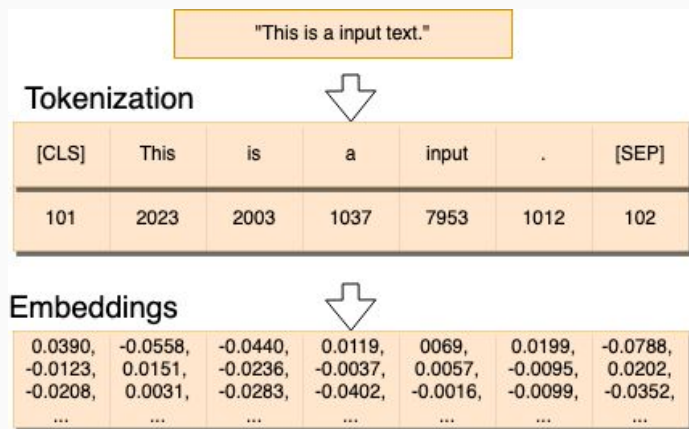- How to Implement RAG?
- Usecase

# What is RAG?

A Process Optimizing outputs of LLMs so it references knowledge base outside training data sources

# What is RAG?

A Process Optimizing outputs of LLMs so it references knowledge base outside training data sources

# Input Embedding

# RAG

- **What is RAG?**
- Why RAG?
- How does RAG functions technically?
- How to Implement RAG?
- Usecase

# RAG

- What is RAG?
- Why RAG?
- How does RAG functions technically?
- How to Implement RAG?
- Usecase

MumPy

aws

# Why RAG?

# Why RAG?

- Cost Effective Implementation
  - Don't need to train or retrain the model
  - Domain Specific information rich replies

# Why RAG?

- Cost Effective Implementation
  - Don't need to train or retrain the model
  - Domain Specific information rich replies
- Up to date with latest information
  - Any new update can be added in few seconds.

# Why RAG?

- Cost Effective Implementation
  - Don't need to train or retrain the model
  - Domain Specific information rich replies
- Up to date with latest information
  - Any new update can be added in few seconds.
- Enhanced User Trust
  - Get Source of the data

# Why RAG?

- Cost Effective Implementation
  - Don't need to train or retrain the model
  - Domain Specific information rich replies
- Up to date with latest information
  - Any new update can be added in few seconds.
- Enhanced User Trust
  - Get Source of the data
- More Developer Control
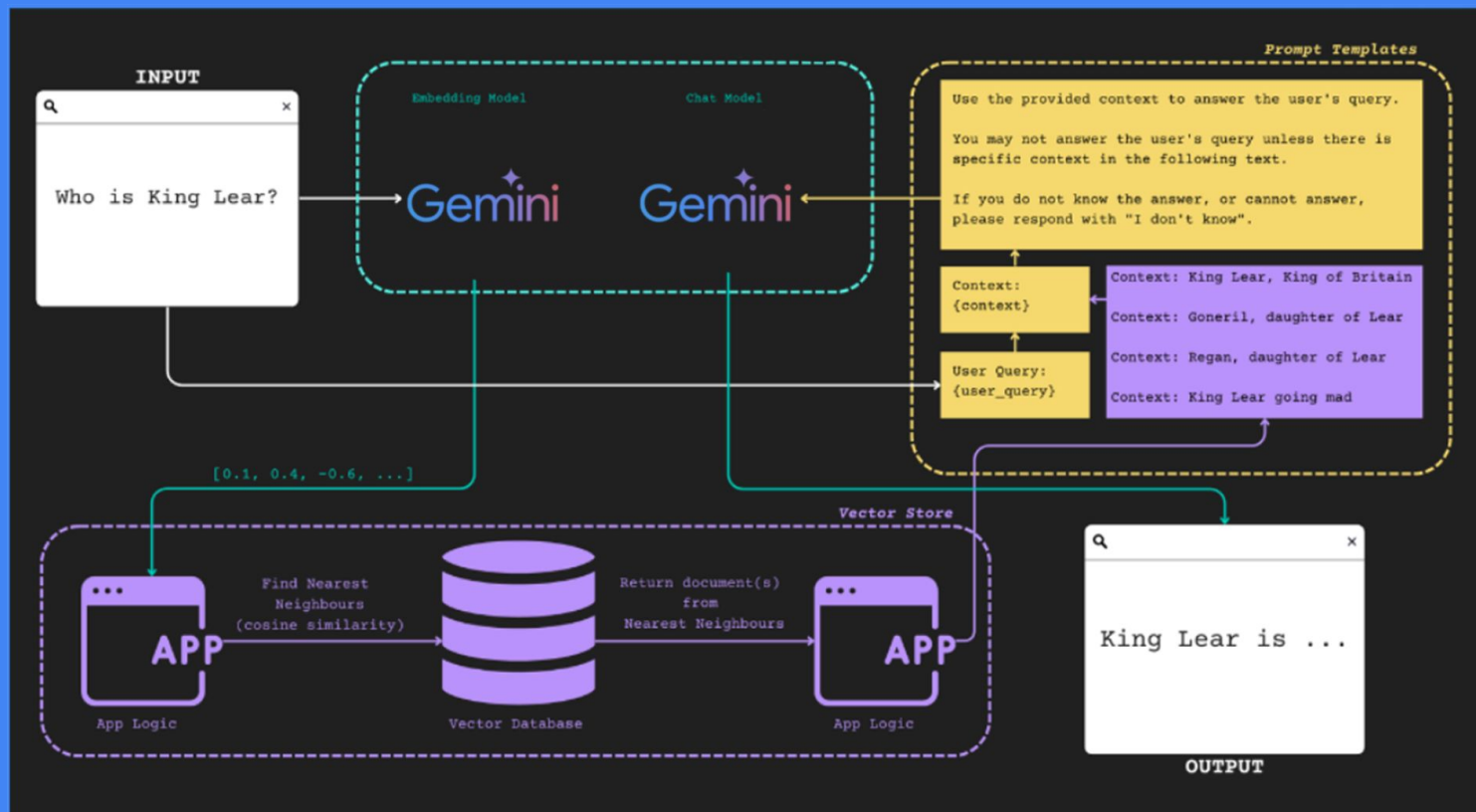  - Easier to Debug
  - Easier to Fix outputs

# RAG

- What is RAG?
- Why RAG?
- How does RAG functions technically?
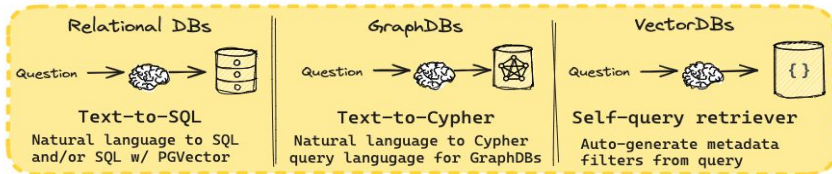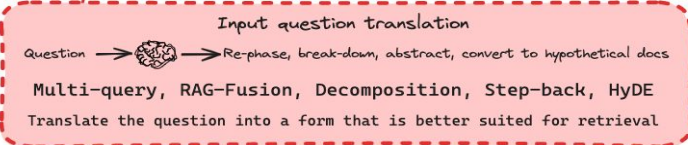- How to Implement RAG?
- Usecase

MumPy

aws

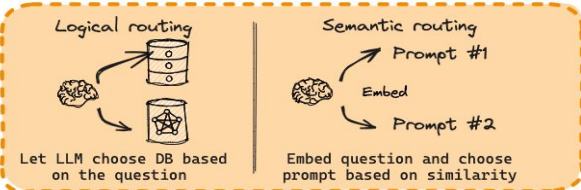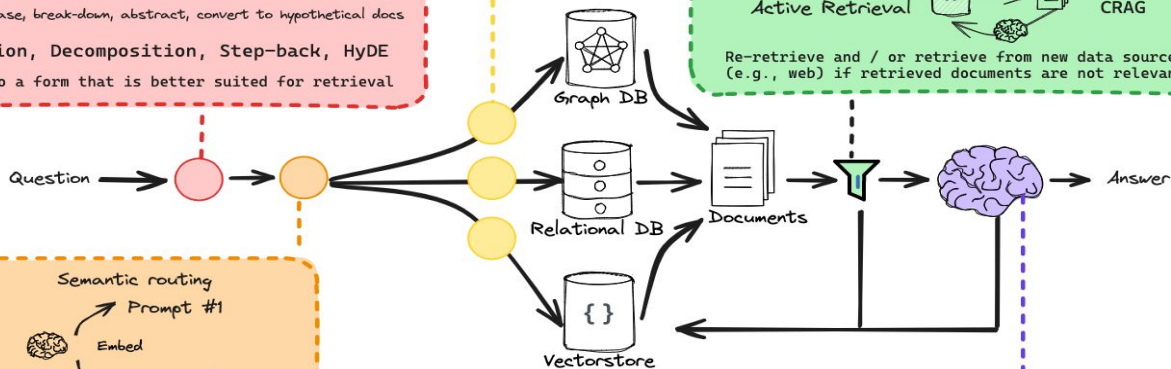# How does RAG functions technically?

**Query Construction**

Relational DBs — GraphDBs — VectorDBs

Text-to-SQL: Natural language to SQL and/or SQL w/ PGVector

Text-to-Cypher: Natural language to Cypher query language for GraphDBs

Self-query retriever: Auto-generate metadata filters from query

**Query Translation**

Input question translation: Re-phase, break-down, abstract, convert to hypothetical docs

Multi-query, RAG-Fusion, Decomposition, Step-back, HyDE: Translate the question into a form that is better suited for retrieval

**Routing**

Logical routing — Semantic routing

Prompt #1, Prompt #2

Let LLM choose DB based on the question

Embed question and choose prompt based on similarity

**Retrieval**

Ranking — Refinement

Re-Rank, RankGPT, RAG-Fusion: Rank or filter / compress documents based on relevance

CRAG

Active Retrieval — CRAG: Re-retrieve and / or retrieve from new data sources (e.g., web) if retrieved documents are not relevant

**Indexing**

Chunk Optimization: Characters, Sections, Semantic, Delimiters

Semantic Splitter: Optimize chunk size used for embedding

Multi-representation indexing: Summary

Parent Document, Dense X: Convert documents into compact retrieval units (e.g., a summary)

Specialized Embeddings: [0.1, ...]

Fine-tuning, ColBERT: Domain-specific and / or advanced embedding models

Hierarchical Indexing: Splits, Summaries, Cluser

RAPTOR: Tree of document summarization at various abstraction levels

**Generation**

Active Retrieval: Answer

Self-RAG, RRR: Use generation quality to inform question re-writing and / or re-retrieval of documents

Graph DB, Relational DB, Vectorstore, Documents, Question, Answer

Source

# How does RAG functions technically?

# Recommendation system



COLLABORATIVE FILTERING

Read by both users

Similar users

Read by her, recommended to him!

CONTENT-BASED FILTERING

Read by user

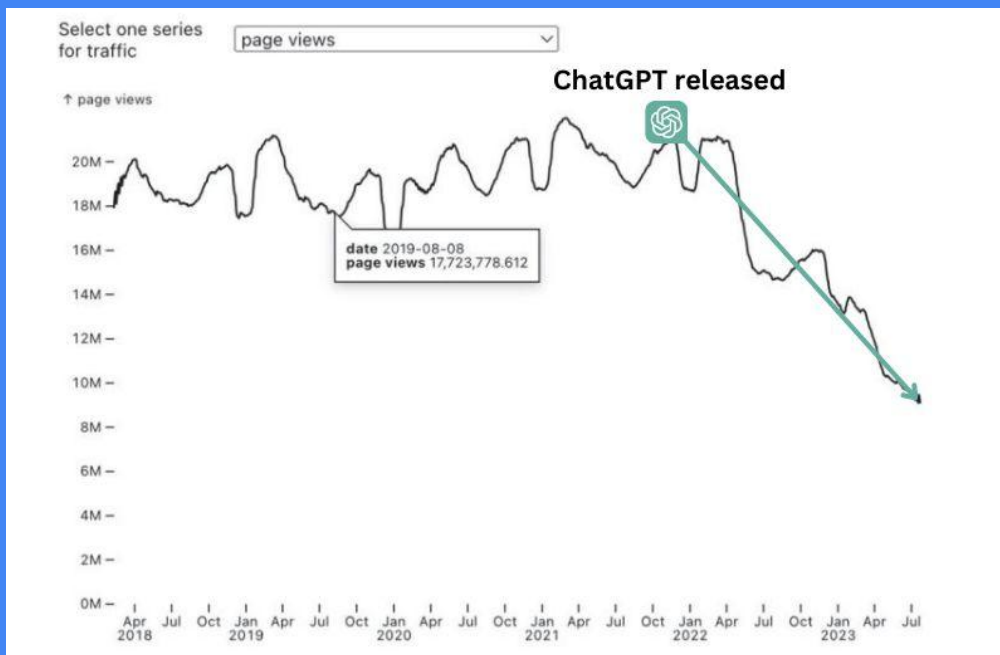Similar articles

Recommended to user
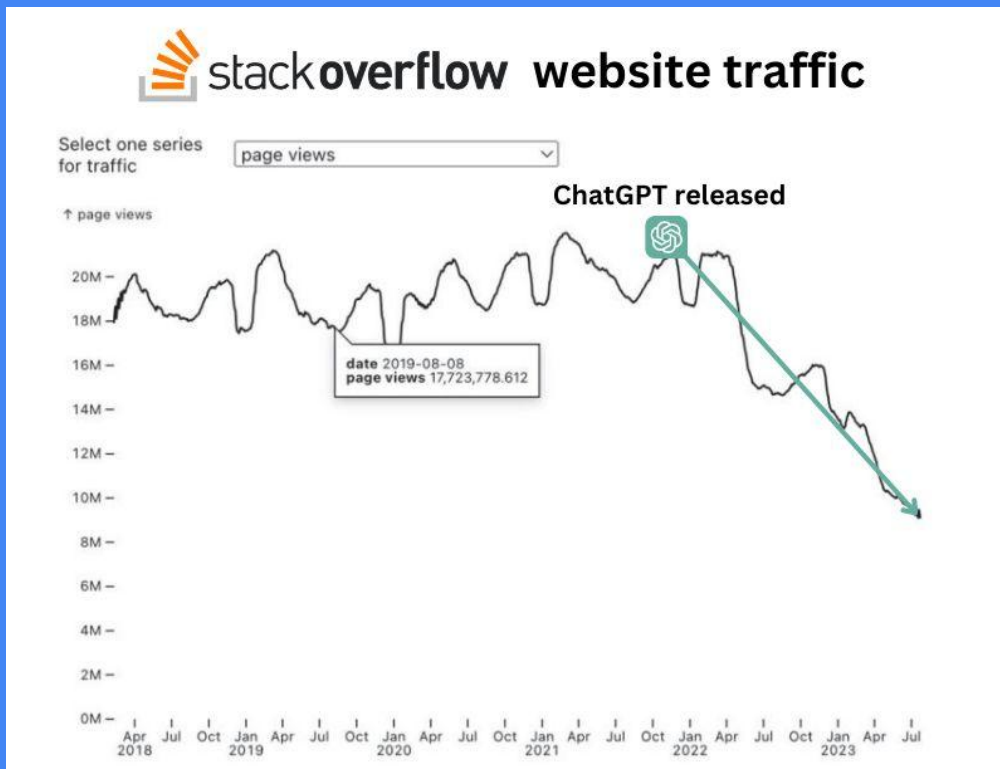
# How to Implement RAG?

# RAG

- What is RAG?
- Why RAG?
- How does RAG functions technically?
- How to Implement RAG?
- Usecase

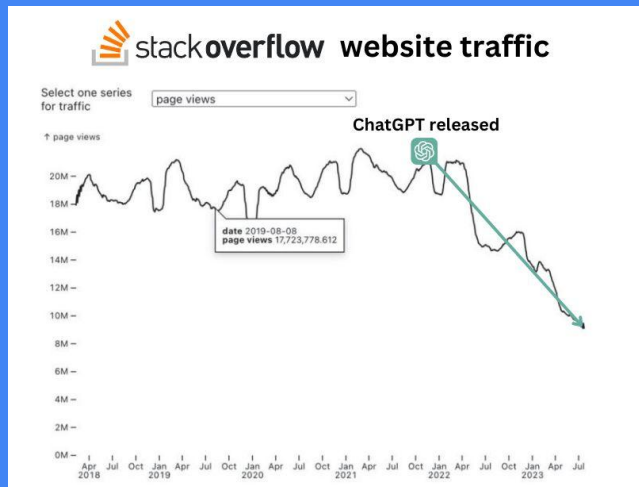# Use Case

# Use Case

# Usecase

# Thanks!

Shubham Agnihotri
Senior Manager - Generative AI



bit.ly/shubhamagnihotri