# Safe harbor statement

The following is intended to outline our general product direction. It is intended for information purposes only, and may not be incorporated into any contract. It is not a commitment to deliver any material, code, or functionality, and should not be relied upon in making purchasing decisions. The development, release, timing, and pricing of any features or functionality described for Oracle's products may change and remains at the sole discretion of Oracle Corporation.

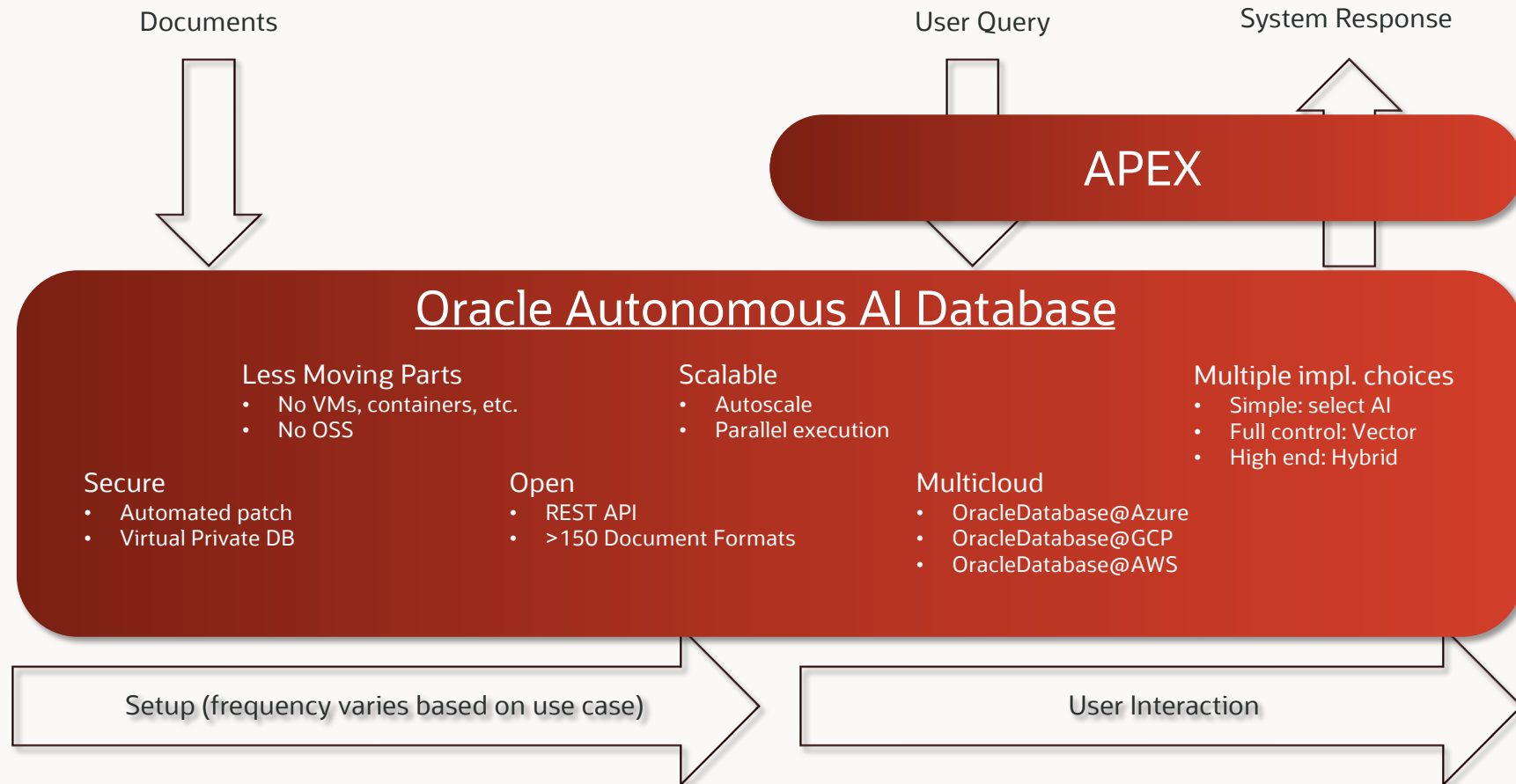# Retrieval Augmented Generation (RAG)

Enables LLMs to use business data to produce better and more contextually relevant answers to user questions while keeping business data secure
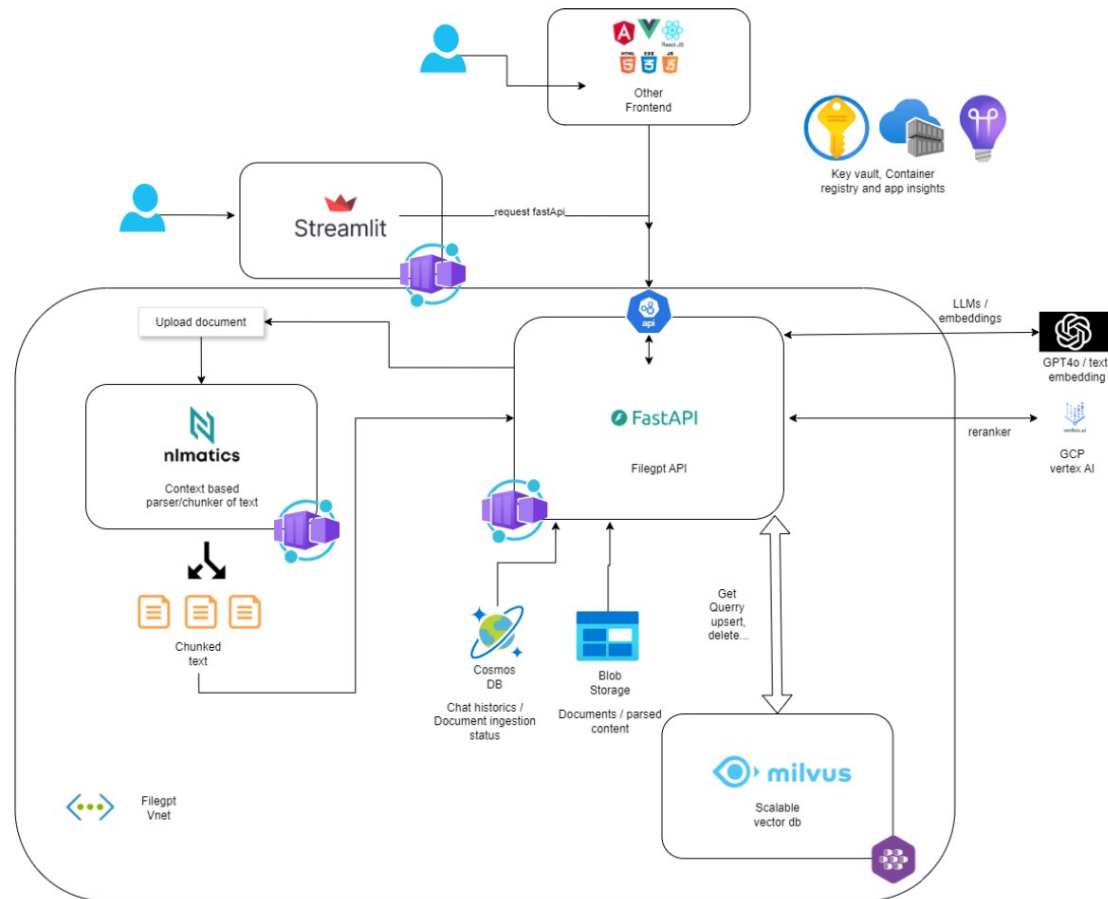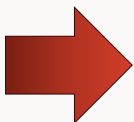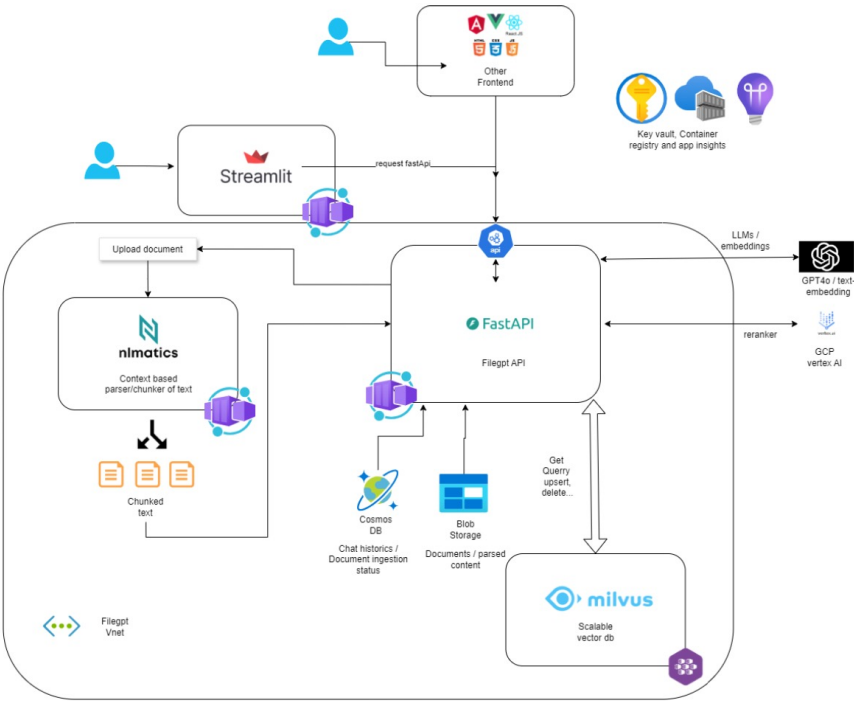
# Typical RAG Pipeline

Documents

User Query

System Response

Chunking

Embeddings

Retrieval

Answer

Python running OSS frameworks, e.g., langchain, nlmatics, etc.

LLM invocation and storage in a Vector DB, e.g. Pinecone, Milvus, etc.

LLM invocation to create embeddings from the user question and compare vector distance from Vector DB info. Eventual rerank.

LLM invocation with local Context.

Chunked text + metadata

Local knowledge as context

Setup (frequency varies based on use case)

User Interaction

# Typical RAG Pipeline

Documents

User Query

System Response

**APEX**

## Oracle Autonomous AI Database

**Less Moving Parts**
- No VMs, containers, etc.
- No OSS

**Scalable**
- Autoscale
- Parallel execution

**Multiple impl. choices**
- Simple: select AI
- Full control: Vector
- High end: Hybrid

**Secure**
- Automated patch
- Virtual Private DB

**Open**
- REST API
- >150 Document Formats

**Multicloud**
- OracleDatabase@Azure
- OracleDatabase@GCP
- OracleDatabase@AWS

Setup (frequency varies based on use case)

User Interaction

# A real architecture as example…


System Architecture

# A real architecture as example…

# The "Three Ways of RAG" using Autonomous AI Database

**Write your own server code**
this is the recommended way when you want total control

**No code / use embedded capabilities**
this is the recommended way when you are focused on getting a result fast

**Agentic approach**
this is the recommended way when RAG is to be one part in a larger set of operations

# The "Three Ways of RAG" using Autonomous AI Database

*Write your own server code – project assumptions*

- Domain Knowledge stored as documents in an object store bucket

- Don't want/need to load the documents into the DB first (duplication/cost)

- No front-end dependencies (once I get the answer how it gets visualized is independent from it)

# The "Three Ways of RAG" using Autonomous AI Database

*Write your own server code – technology choices*

- Select AI
  - Essentially just two calls
  - Limited control on process settings
  - SQL interface

- Vector Index
  - Definitely more than two calls
  - Full control on chunking, embedding, results filtering, and answer creation

- Hybrid Vector Index
  - Complex but very powerful query syntax
  - Less control on the settings
  - Huge optimization (avoid expensive reranking!) through union of Text and Vector searches

# The "Three Ways of RAG" using Autonomous AI Database
*No code approach – just load your documents*

# The "Three Ways of RAG" using Autonomous AI Database
*Agentic RAG – build your own RAG agent*



Copyright

# Recap & Take Aways

**1** **Oracle Autonomous AI Database helps reduce complexity**
It is multicloud, secure and scalable

**2** **Multiple technology paths**
Freedom of choice between code, no-code, and agentic approach

**3** **Focus on the business case of the RAG solution**
Technology makes it simple, but be very clear on what business problem you are addressing

RAG in a Box
Blog Article

Select AI

Vector Search

Hybrid Vector Search

Docs

**Your feedback is important.**

**Scan this QR Code or use the Mobile App to share your thoughts on this session.**

To contact us

**Massimo Castelli**

**Doug Hood**

# Thank you!

ORACLE