

---

---

# Scaling AI with Large Language Models (LLMs)

*Meethun Kumar Panda*

*Conf42 Large Language Models (LLM)  
Conference  
06-Mar-2025*

---

---

# Key design principles to enable Gen AI use case at scale utilizing cloud

## Generative AI Characteristics

## Key Design Principles

it is important to be able to experiment with different models and element technologies.

### Flexibility

- **Replaceable modules:** The components of the application are modular and technology-agnostic
- **Flexible storage:** Scalable, platform-agnostic storage solutions
- **Self-healing architecture:** the ability to identify and recover from errors and misconfigurations without human intervention

Focusing on larger-scale unstructured data than conventional AI/ML models

### Scalability

- **Reusable module:** Deploy a repeatable infrastructure stack using "infrastructure as code (IAC)"
- **Coordination:** Control application traffic and other architectural components, such as DNS and load balancing
- **Elastic scaling / hosting:** Ability to automatically add/remove infrastructure based on changes in traffic patterns (such as tuning compute configurations based on CPU, memory, and storage)

Foundation models are 'black boxes' in structure, with low explanatory power and error traceability

### Monitoring

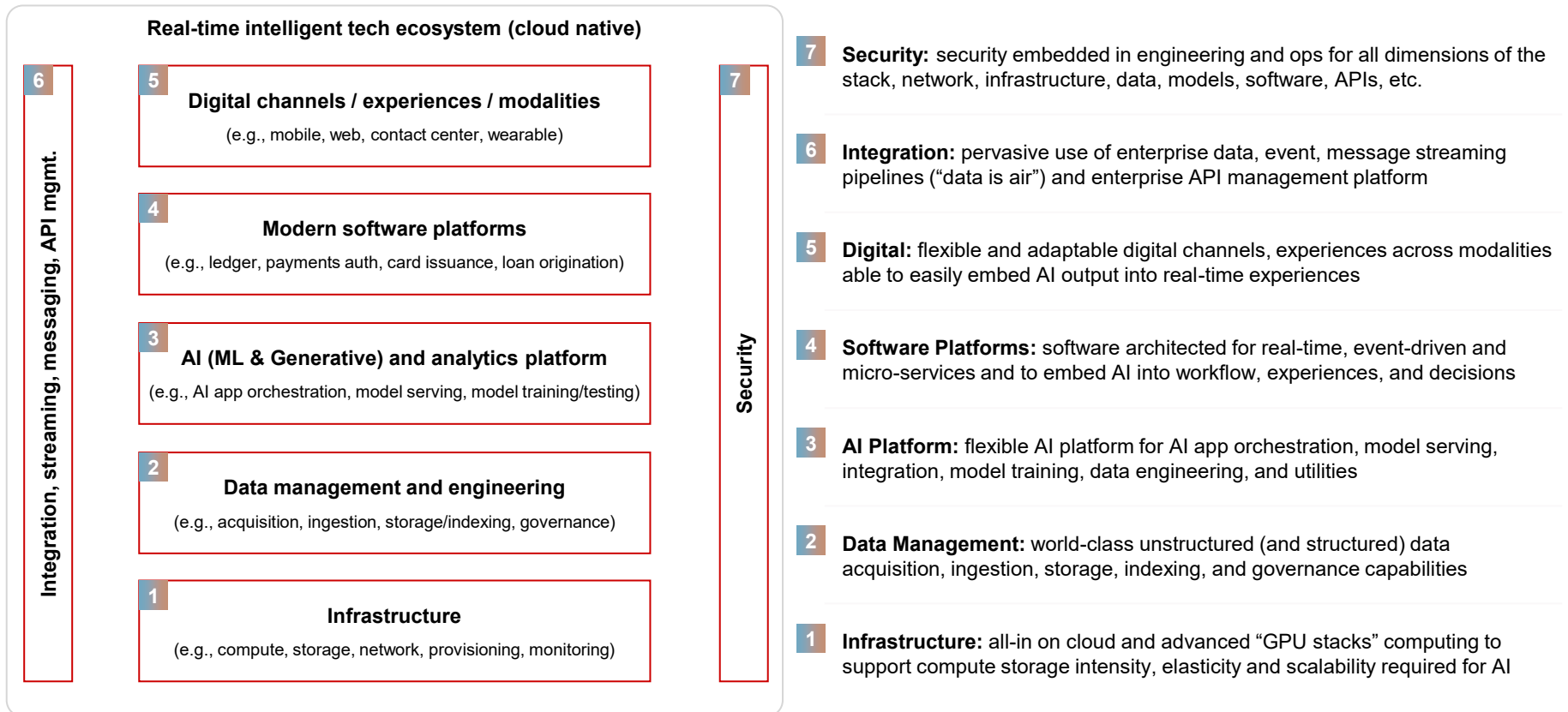
- **Process transparency:** Enables all processes in the architecture to be observable and proactively diagnose issues
- **Drift/Deflection Control:** Implement restrictive measures and mechanisms to monitor/manage drift and bias in the process
- **Continuous improvement:** monitoring for performance benchmarking; Implement a systematic process for incorporating feedback

It is common for foundation models to be connected externally through APIs.

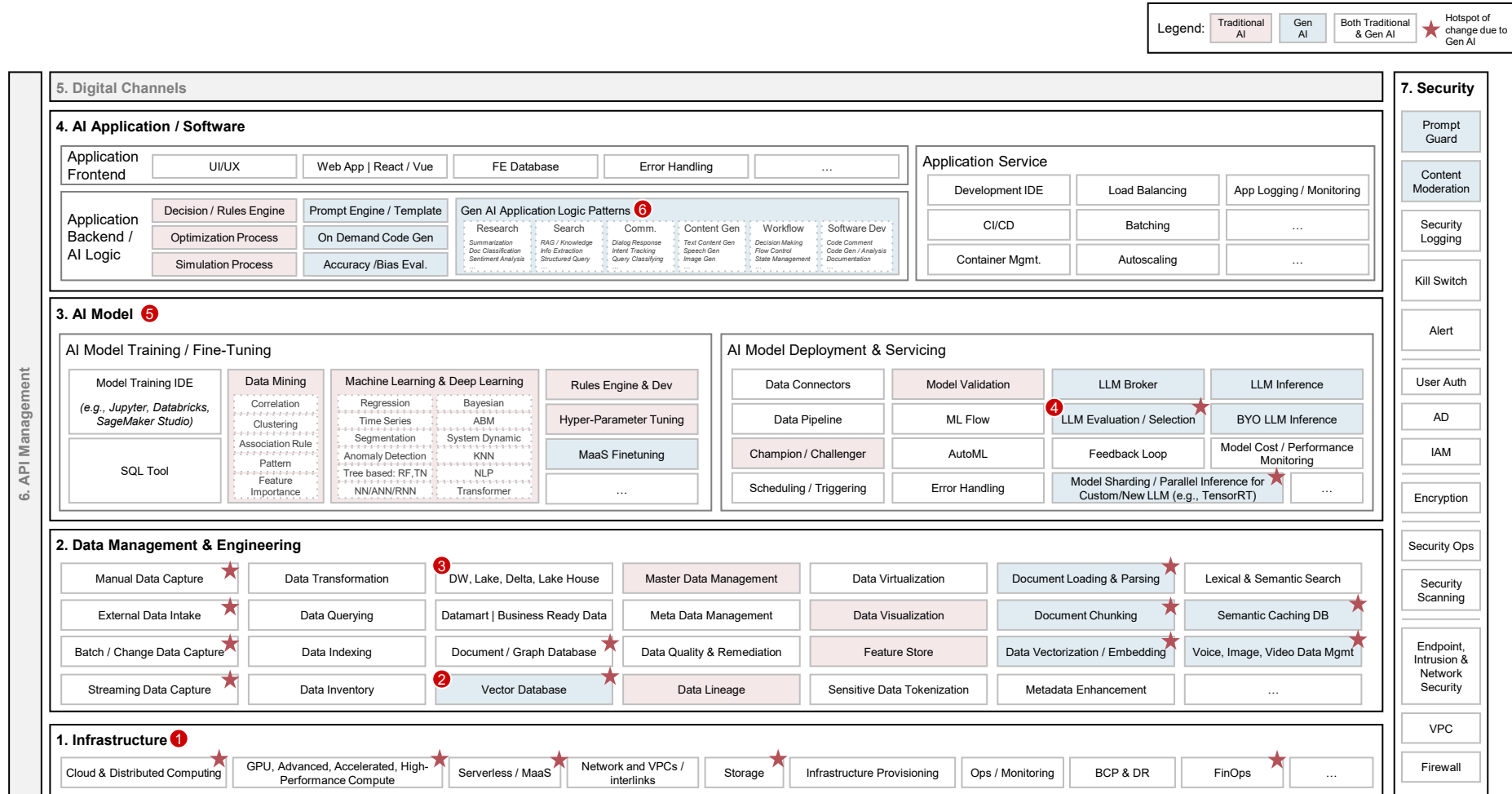
### Security

- **Data management:** Secure management of sensitive data (such as PII)
- **Access management:** Enforce access management to restrict access to infrastructure and authenticate users of applications
- **API key and secret management:** (1) centralized key vault, (2) dynamic secret, (3) automatic secret injection, and (4) monitoring/auditing for secret activities

# Architecting AI-Driven Ecosystems: A Bottom-Up and Destination-Back Approach



# AI platform reference architecture (traditional & Gen AI)



# Strategic foundational choices on Generative AI building blocks

①	<b>Infrastructure</b>	Single cloud	Multi-cloud	SaaS foundation model API	Self-hosted hybrid/on-premise
		GPU full-stack AI platform (NVIDIA)	Cloud vendor specialized AI chips	Niche provider AI chips	
②	<b>Vector database</b>	Enterprise-level vector storage		Application-level vector storage	
		Vector indexing added to existing data store		Specialized vector data store	
③	<b>Data architecture</b>	Single integrated data lake		Multiple federated data lakes / data mesh	
		Integrated structured and unstructured storage (e.g. Delta Lake)	Separate structured and unstructured/vector data stores	Separate structured, unstructured, and vector stores	
④	<b>LLM(s) selection</b>	Large foundation model	Customized large foundation model	Multiple smaller specialized FMs	
		Open source/proprietary		Commercial	
⑤	<b>AI / ML Ops platforms</b>	Open source	Niche best of breed	Hyperscale cloud vendor (e.g., AWS SageMaker)	GPU full-stack AI platform (e.g., NVIDIA)
⑥	<b>Gen AI app multi-cloud architecture</b>	All on same cloud as LLM (ingest, vector db, AI app FE/BE, LLM)		AI app front-end on current cloud, all else on same cloud as LLM	All on current cloud; LLM hosted on relevant LLM cloud

# Unstructured data readiness for “AI everywhere” is a critical pre-requisite

## Three key data management challenges



**Organizing, cleaning, & labeling and catalog unstructured data**

### Data requirements for Gen AI

Data must be **structured, accurate, and contextually “sound”** before it’s fed into foundation models for Gen AI solutions

### Why it’s challenging

**Most data is unstructured** (e.g., live in PDFs or images) and must be **organized, cleaned, and labeled** prior to being used in a data pipeline

### Ecosystem perspective



*Finding and labeling data for a given use case has been the singular bottleneck - You’re set up for success when you can properly label the data for the specific use case at scale.*



**Storing data to be effectively used by Gen AI solutions**

Data needs to be stored in a way that’s **compatible with being “fed” into foundation models** (e.g., in vector databases)

Company **with limited or poorly maintained** data storage practices will have a **vast amount of data that is not ready for ingestion** into foundation models



*Model output is directly related to how good of an input you can put into the model – the benefit of working with a database of real data is that it’s generally really good input.*



**Extracting, transforming, & loading data into foundation models**

Foundation models **continuously learn and adapt**, which requires **constant managing of data pipelines** to ensure compliance, accuracy, and sustaining of data privacy

Companies must establish **appropriate data pipelines** and integration processes in line with overall **data governance strategies**



*It’s not only about cleaning data at rest. It’s about being able to exchange or move that data around... If we can have data in a state that allows us to move it elsewhere and plug it into [FMs] it makes a huge difference.*

# What is hard about building and operationalizing Generative AI beyond typical AI/ML project challenges

## TECHNICAL CHALLENGES



- Decreasing **hallucinatory response** and **improving quality of output** requires **new prompt engineering skills** and **new technical skills** (e.g., grounding, fine-tuning, chaining)
- **Vastly different MLOps** with many tools nascent or non-existing (e.g., vector DBs, p-stores, chaining, failovers, guardrails)
- Creation of **grounding databases** (e.g., sources, chunking)
- **New types of cyber-security risks** (e.g., chat prompt injection)
- Increasing **scalability and deployment** efficiency (incl. caching)
- More and **diverse data improves performance** but requires **extensive integration**
- **Balancing data privacy and security** with cloud LLMAaaS

## OPERATIONAL CHALLENGES



- Picking smart (build vs buy) & **building a roadmap** that balances today forward with future-back
- Research **user experience and unmet needs** to develop hypotheses on agent automation
- Need for a different type of **UI innovation**
- **System to identify issues** with response pathways to feed back into development
- Estimating value and implementing **systems to track value and prove real value**
- Workforce **engagement and training** at scale
- Planning for **radically redesigned processes** and **potential organizational changes**

# GenAI tech architecture is structured around six key elements

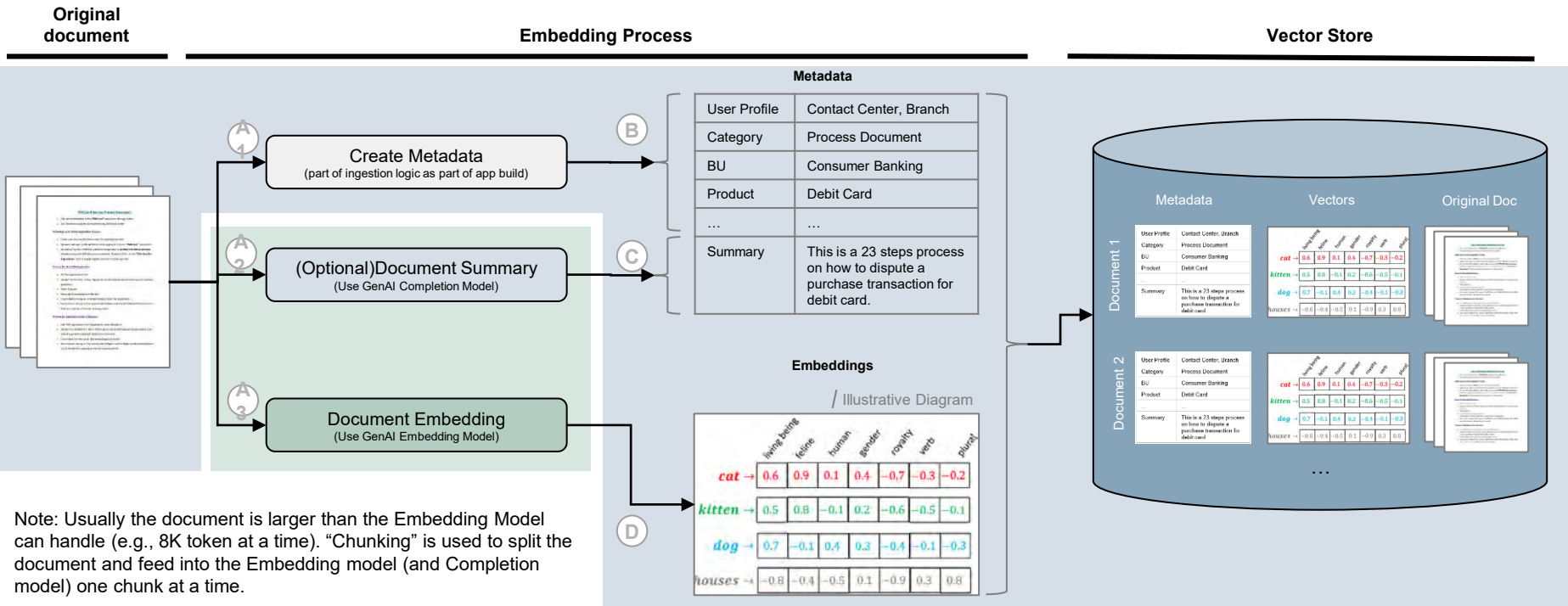
GenAI architecture elements		Description	Considerations for decisions
GENAI TECH STACK	LARGE LANGUAGE MODELS	<ul style="list-style-type: none"> <li>Model trained on <b>vast amounts of data</b> to understand the structure of natural language</li> <li>Enables generation of <b>novel content</b> and response to questions</li> </ul>	<ul style="list-style-type: none"> <li>Capabilities</li> <li>Data security</li> <li>Ease of use and deployment</li> <li>Multilingual capabilities</li> </ul>
	EMBEDDINGS	<ul style="list-style-type: none"> <li>Representation of content and meaning that captures <b>semantic and contextual relationships</b></li> <li>Enables comparison of user questions with data in knowledge base <b>to determine relevance</b></li> </ul>	
	VECTOR DATABASE	<ul style="list-style-type: none"> <li><b>Storage of data</b> with their embeddings</li> <li>Enables <b>retrieval of relevant context</b> for specific queries based on semantic and contextual relationships</li> </ul>	
DATA & INFRASTRUCTURE	PLATFORM	<ul style="list-style-type: none"> <li><b>Consolidates data from various sources</b> and provides tools for analysis, application development, and visualizations</li> <li><b>Accelerates AI application</b> and insights delivery</li> </ul>	<ul style="list-style-type: none"> <li>Readiness of the options</li> <li>Compliance with security norms</li> <li>Data access availability</li> <li>Scalability to accommodate new use cases</li> </ul>
	HOSTING	<ul style="list-style-type: none"> <li><b>Computational resources</b> and storage to run applications and store application data</li> <li><b>Enables end users to utilize application</b> in line with data security requirements</li> </ul>	
	DATA STORE	<ul style="list-style-type: none"> <li><b>Data utilized for each use case</b>, stored in business applications and IT systems</li> <li>Data needs to be provided and pre-processed <b>in compliance with data security and privacy requirements</b></li> </ul>	



## Semantic knowledge search use case example

## Step 1: Preparation Knowledge Vector Store

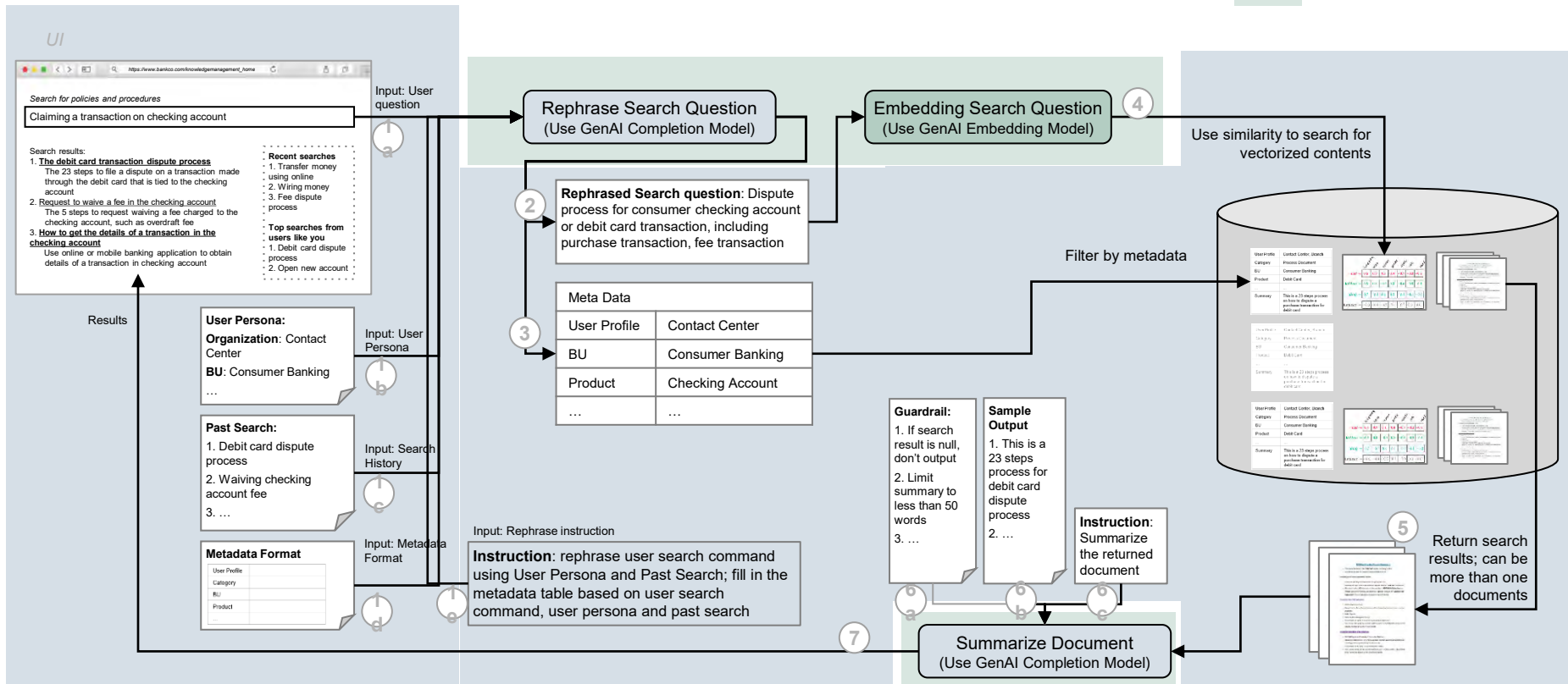
A key step is to ingest the document and prepare the vector store. Both metadata and document embeddings are prepared. The more precise and complete the metadata, the better the search result.



# Semantic knowledge search use case example

## Step 2: Core app function – semantic search of knowledge base

AWS  
Azure



# Use case delivery requires multiple technical roles and resources

## Use case delivery requires multiple technical roles and requirements

