



# Modern AI Architecture for Intelligent Banking

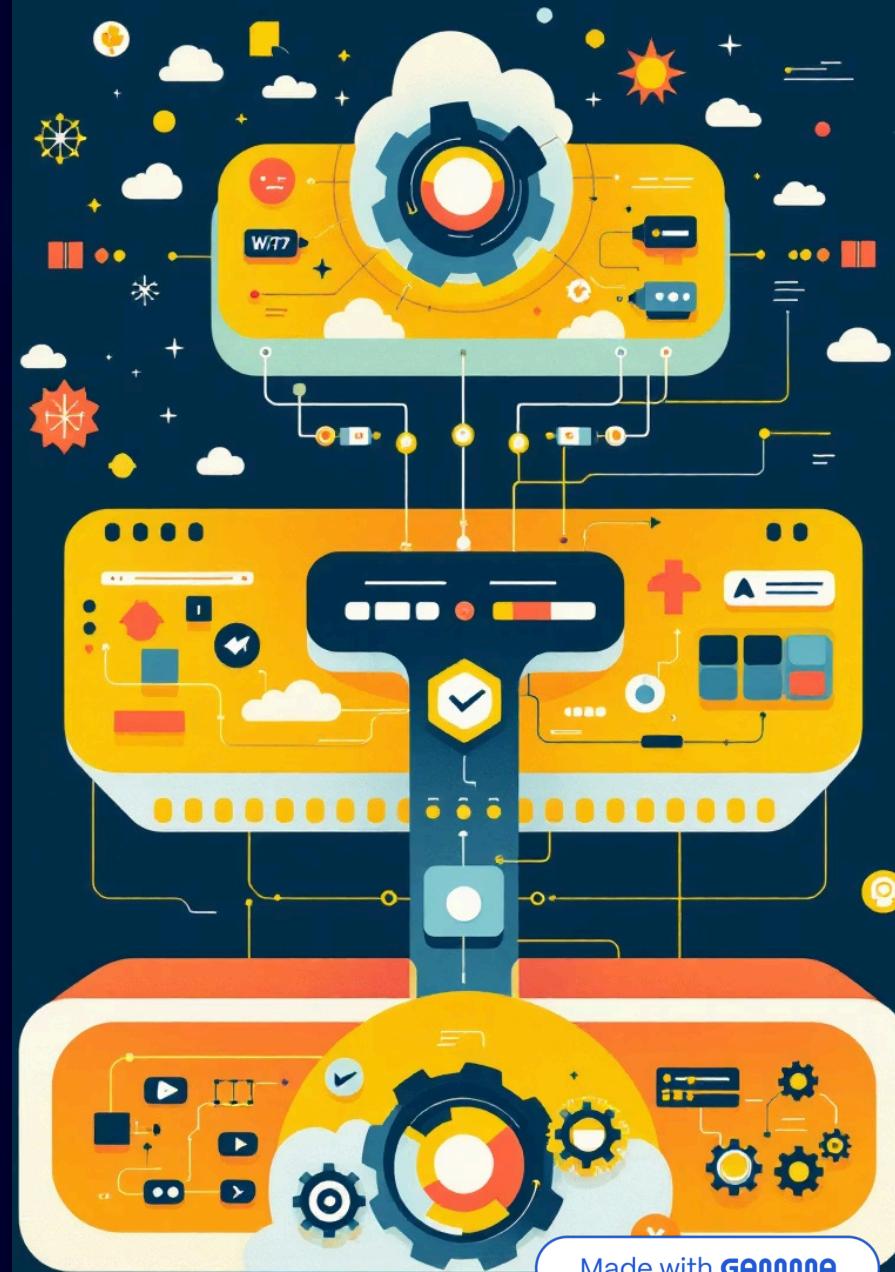
Transforming traditional banking systems into cognitive, learning institutions through layered intelligence

## ARCHITECTURE OVERVIEW

# The Cognitive Banking Architecture

A comprehensive framework for implementing AI-driven intelligence across banking operations. This architecture transforms siloed systems into an integrated cognitive platform that learns, adapts, and delivers actionable insights.

Six distinct layers work in harmony: from foundational systems of record through intelligent reasoning engines to customer-facing applications. Each layer builds upon the previous, creating a complete intelligence stack.



# Core Banking & Enterprise Platforms

Systems of Record — The Foundation



## Core Banking System

Central ledger and account management infrastructure



## AML Transaction Monitoring

Real-time surveillance for suspicious activities



## KYC & Onboarding

Customer identity verification and due diligence



## CRM & Channels

Customer relationship and omnichannel management



## Payments & Trade Finance

Transaction processing and trade operations

"Systems that know what happened – excellent at recording transactions, but not interpreting meaning"

# Real-Time & Batch Data Pipeline

The Nervous System of Modern Banking



## Kafka / Event Streaming

Captures every transaction, customer interaction, and system event as it occurs, enabling real-time decision making

## ETL / Data Lake

Consolidates historical data from disparate sources into a unified repository for analysis and pattern detection

## CDC Pipelines

Change Data Capture ensures seamless synchronization between operational systems and analytical platforms

- This layer moves data from operational silos to centralized intelligence, transforming static records into flowing, actionable information streams.

# Embedding & Feature Engineering Layer

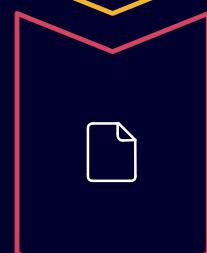
## Converting Raw Data into Numerical Meaning

The transformation layer where traditional banking data becomes AI-ready intelligence. Machine learning models convert complex transactions, documents, and behaviors into dense vector representations that capture semantic meaning.



### Transaction Behavior Embeddings

Payment patterns, frequency, amounts, and counterparty relationships encoded as vectors



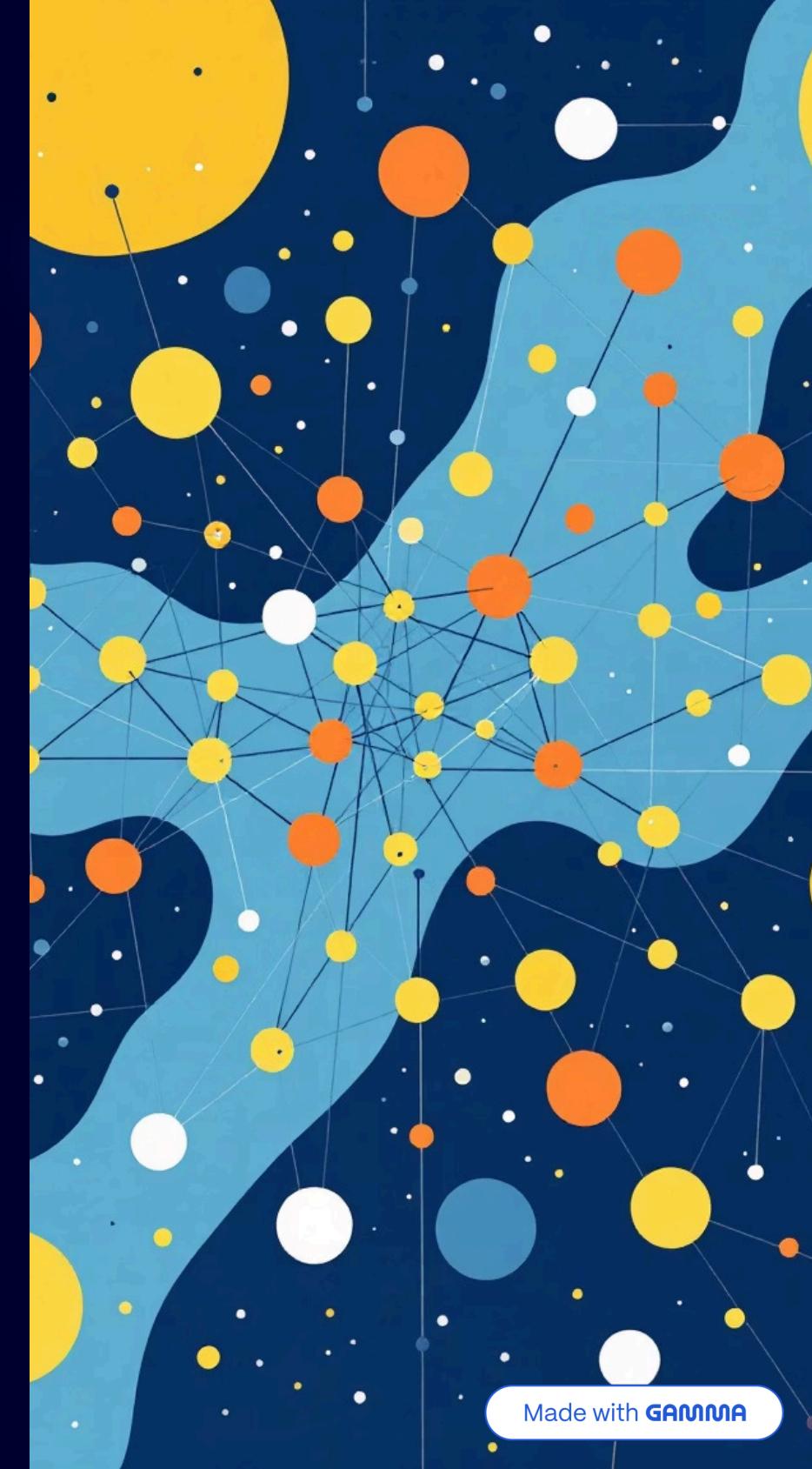
### KYC Document Embeddings

Identity documents, beneficial ownership structures, and compliance records transformed into searchable formats



### Fraud Pattern Embeddings

Historical fraud cases and suspicious activity patterns captured as reference vectors



# Enterprise Vector Database

The Institutional Memory of the Bank

## Customer Behavior Vectors

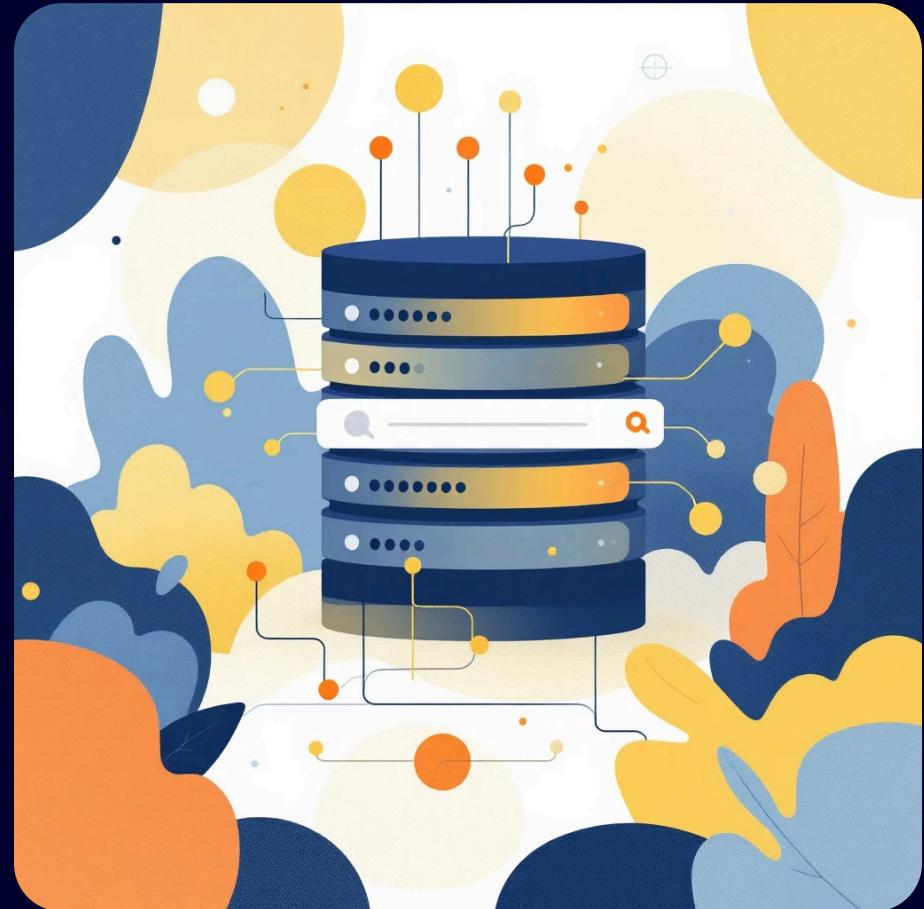
- Transaction patterns across time
- Product usage profiles
- Risk indicators and anomalies

## Case History Vectors

- Investigation outcomes and resolutions
- Regulatory filing precedents
- Decision rationales and approvals

## Fraud Pattern Vectors

- Known fraud typologies and schemes
- Emerging threat signatures
- Cross-channel attack patterns



## Advanced Indexing Technology

**HNSW** (Hierarchical Navigable Small World): Graph-based indexing for sub-millisecond similarity search at scale

**IVF + PQ** (Inverted File + Product Quantization): Compressed storage enabling billions of vectors with minimal memory footprint

# LLM Reasoning & Orchestration Layer

The Brain — Turning Memory into Intelligence

01

## Retrieval-Augmented Generation (RAG)

Queries the vector database to find relevant historical cases, similar transactions, and applicable regulations before generating responses

02

## Context Reasoning

Synthesizes retrieved information with real-time data, applying domain knowledge and regulatory requirements to form comprehensive understanding

03

## Explanation Engine

Generates human-readable justifications for decisions, linking conclusions to evidence, policies, and precedents for full auditability

"The vector database remembers every case and pattern. The LLM reasons across that memory to generate insights."



Made with GAMMA

# AI-Powered Banking Applications

## Where Intelligence Becomes Business Value

The culmination of the cognitive architecture: user-facing applications that empower banking professionals with AI-augmented capabilities. Each application leverages the full intelligence stack to deliver contextualized insights.



### AML Investigator Copilot

Accelerates case review by surfacing similar historical investigations, suggesting next steps, and auto-generating preliminary reports



### Fraud Explanation Engine

Provides real-time explanations for fraud alerts, connecting suspicious patterns to known typologies with confidence scores



### KYC Analyst Assistant

Streamlines customer due diligence by extracting key information from documents and flagging potential risks automatically



### Relationship Manager AI

Equips RMs with comprehensive customer intelligence, product recommendations, and proactive risk alerts during client interactions

# The Intelligence Flow

From Transaction to Action in Milliseconds



## Transaction

Customer activity triggers the flow



## Vector

Embedded as semantic representation



## Retrieval

Similar patterns found in memory



## Reasoning

LLM synthesizes context and insight



## Action

Intelligent response delivered

---

This horizontal flow represents the complete journey of information through the cognitive banking stack. Each transaction becomes an opportunity for learning, pattern recognition, and intelligent intervention.

# Evolving from Transaction Processors to Learning Systems



**At the bottom**, we have core banking and AML systems – excellent systems of record that know what happened, but not what it means.

**The data pipeline** serves as the nervous system, streaming transactions and documents in real time to create a continuous flow of information.

**The embedding layer** converts raw transactions, KYC documents, and behaviors into numerical meaning that machines can process and compare.

**The vector database** becomes the bank's semantic memory, using HNSW to instantly find similar past cases and patterns across billions of records.

**The LLM reasoning engine** thinks across this memory, synthesizing insights that would take humans hours or days to uncover.

**Finally**, this intelligence appears in copilots for investigators, fraud teams, and KYC analysts – making every banking professional more effective.

- **This is the future of banking:** Institutions that don't just process transactions, but learn from every interaction, remember every pattern, and continuously evolve their intelligence.