# Context-Aware FinCommerce Engine using Vector Memory with Qdrant

## Project Overview

This project presents a **Context-Aware FinCommerce Engine** that enables intelligent product discovery and personalized recommendations by combining **semantic understanding** with **user-specific financial context**.

Instead of relying solely on keyword-based search or static recommendation rules, the platform uses **vector embeddings stored in Qdrant** to capture semantic meaning across products, user profiles, financial constraints, and interaction history.

At query time, the system retrieves semantically relevant products and then **re-ranks them using affordability constraints and personal preferences**, ensuring that recommendations are not only relevant but also financially realistic for each user.

The solution demonstrates how **Qdrant can act as a vector-native memory layer** for real-time, context-aware decision-making in FinCommerce applications.

## Problem Statement

Modern e-commerce and fintech platforms struggle to deliver recommendations that are both **semantically relevant** and **financially appropriate** for individual users.

Traditional recommendation systems typically:

- Rely on keyword-based search or collaborative filtering

- Ignore real-time financial constraints such as account balance or credit limits

- Treat personalization, affordability, and user behavior as disconnected systems

This results in unrealistic product suggestions, poor user experience, and reduced conversion rates.

There is a growing need for a **unified, vector-native system** capable of reasoning simultaneously over **products, user preferences, financial context, and interaction history**, while supporting fast retrieval and filtering at scale.

## Use Case Solved

The project addresses the **Context-Aware FinCommerce Recommendation** use case.

When a user submits a natural-language query (e.g., *"Laptop for machine learning under 1500"*), the system must return products that:

- Are semantically relevant to the user's intent

- Match the user's personal preferences (brands, categories)

- Respect financial constraints such as available balance or credit limits

- Adapt based on historical interactions

Two users issuing the same query may receive different recommendations depending on their financial situation and preferences. This enables **real-time, personalized, and financially realistic discovery**, which is not achievable with traditional search engines.

# How Qdrant Is Used

Qdrant serves as the **central vector memory and retrieval engine** of the system.

## Vector Collections

The system uses multiple Qdrant collections, each representing a different type of memory:

- **Product Collection**: Stores semantic embeddings of product descriptions along with structured metadata such as price, category, brand, and availability.

- **User Profile Collection**: Represents user preferences and personalization signals.

- **Financial Context Collection**: Stores financial constraints such as balances, income, and credit limits.

- **Interaction Memory Collection**: Captures historical user interactions and queries.

## Query-Time Workflow

At query time, the system performs the following steps:

1. The user query is embedded into a vector representation.

2. Qdrant retrieves the most semantically similar products using vector similarity search.

3. Payload filters enforce structured constraints such as price limits and availability.

4. Results are re-ranked using affordability and preference-based scoring.

## Why Qdrant

Qdrant enables:

- High-dimensional semantic search over multimodal data

- Hybrid retrieval combining vector similarity and structured filtering

- Low-latency performance suitable for real-time personalization

- A scalable and explainable architecture for FinCommerce intelligence

Rather than acting as a passive database, Qdrant functions as an **active intelligence layer** powering context-aware recommendations.

## Alignment with Hackathon Use Case

This project directly addresses **Use Case 2: Context-Aware FinCommerce Engine**, show-casing how Qdrant can enable personalized, constraint-aware discovery and recommendation using vector-native memory at scale.