

Context-Aware FinCommerce Engine using Vector Memory with Qdrant

January 2026

Contents

1	Introduction	3
2	Problem Statement	3
3	Use Case Description	4
3.1	Representative Scenario	4
4	Project Overview	5
4.1	Input → Processing → Output	5
4.2	Role of Each Component	5
5	System Architecture	6
5.1	High-Level Flow	6
5.2	System Components	7
6	Qdrant Vector Database Design	7
6.1	Why Qdrant	7
6.2	Collections Design	7
6.2.1	Collection 1: products_multimodal	8
6.2.2	Collection 2: user_profiles	9
6.2.3	Collection 3: financial_contexts	9
6.2.4	Collection 4: interaction_memory	10
6.3	Summary Table	10
7	Embedding & Semantic Search	11
7.1	Embedding Model	11
7.2	Vector Semantics	11
7.3	Similarity Search	11
7.4	Performance	12
8	Scoring & Ranking Strategy	12
8.1	Signal 1: Semantic Score (s_{sem})	12
8.2	Signal 2: Affordability Score (s_{aff})	12
8.3	Signal 3: Preference Score (s_{pref})	13
8.4	Signal 4: Collaborative Score (s_{collab})	13
8.5	Signal 5: Popularity Score (s_{pop})	14
8.6	Final Ranking Formula	14
8.7	Rationale for Weights	14
9	Real-Time Interaction Learning	15
9.1	Interaction Types & Weights	15
9.2	Interaction Logging Mechanism	15
9.3	Popularity vs. Collaborative Filtering	15
9.4	Time Decay Rationale	16

10 Explainability Layer	16
10.1 Explanation Framework	16
10.1.1 1. Semantic Relevance Explanation	16
10.1.2 2. Affordability Explanation	16
10.1.3 3. Preference Explanation	16
10.1.4 4. Collaborative Explanation	17
10.1.5 5. Popularity Explanation	17
10.2 Benefits	17
11 User Interface	17
11.1 Interface Sections	17
11.1.1 Sidebar: Context Configuration	17
11.1.2 Main Panel: Search & Results	18
11.2 Real-Time Feedback	18
12 Scalability & Extensibility	18
12.1 Scalability Dimensions	18
12.1.1 Vertical Scaling: More Data	18
12.1.2 Horizontal Scaling: Distributed Deployment	18
12.2 Extensibility	19
12.2.1 Multi-Modal Embeddings	19
12.2.2 Cross-Domain Catalog	19
12.2.3 Advanced Collaboration	19
12.3 Production Readiness	19
13 Conclusion	20

Introduction

Modern e-commerce platforms face a fundamental challenge: delivering recommendations that are simultaneously *semantically relevant*, *financially feasible*, *personally aligned*, and *behavior-informed*. Traditional keyword-based search systems decouple these dimensions, forcing users to navigate fragmented experiences across search, price filters, trending widgets, and recommendation sidebars.

The problem is not one of retrieval—search engines excel at finding relevant items. Rather, it is one of *contextual synthesis*. Users shopping online operate within multiple constraints and motivations simultaneously:

- **Financial Realities:** Available balance, credit limits, debt obligations, installment eligibility
- **Semantic Intent:** Product specifications, use cases, domain language
- **Personal Preferences:** Trusted brands, favorite categories, style preferences
- **Social Proof:** What similar users purchased, trending items in real-time

Vector databases like Qdrant enable a fundamentally different paradigm: *vector-native memory*. Instead of treating vectors as a secondary feature layer, we embed all entities—products, users, financial states, and behavioral interactions—into a shared semantic space. This allows sub-second hybrid retrieval that simultaneously reasons over similarity, affordability, preferences, and trends in a single unified query.

This document describes a complete Context-Aware FinCommerce Engine built on Qdrant, demonstrating how vector memory transforms e-commerce personalization from a post-retrieval ranking problem into an integrated, real-time intelligence system.

Problem Statement

E-commerce personalization today operates as a series of disconnected systems:

- **Search Engine:** Returns semantically relevant items based on keywords and metadata
- **Price Filter:** User manually narrows results by budget
- **Recommendation Widget:** Displays trending or collaborative items separately from search results
- **Personalization Engine:** Applies user preferences offline (not integrated with live search)

This separation creates several failure modes:

1. **Affordability Blindness:** A customer searching for "laptop for machine learning" receives a 4,000\$ laptop despite having a 500\$ balance and no credit. The system never reasons: "This person cannot afford this."
2. **Preference Ignorance:** Recommendations ignore user preferences. A customer loyal to Apple brand receives the same generic results as someone indifferent to brand.

3. **Trend Isolation:** Trending products are surfaced in a separate widget, disconnected from personalized search. Users miss contextually relevant trending items.
4. **Behavioral Amnesia:** User behavior (clicks, views, purchases) influences recommendations only after hours or days of batch processing. Real-time interaction signals are lost.
5. **Explainability Vacuum:** Users see results but never understand *why*. Trust in recommendations erodes without transparency.

These problems compound in finance-aware e-commerce, where affordability is non-negotiable. A luxury item recommendation to a budget-constrained user damages brand trust and wastes user attention.

Core Challenge: Unify semantic relevance, financial feasibility, personal preferences, collaborative signals, and real-time popularity into a single sub-second recommendation engine with human-readable explanations.

Use Case Description

This project explicitly solves **Use Case 2: Context-Aware FinCommerce Engine**.

Representative Scenario

User Query: “Laptop for machine learning under \$1500.”

User Context:

- Available balance: \$800
- Credit limit: \$1000
- Preferred brands: Apple, Dell
- Preferred categories: Electronics, Computing
- Recent purchases: Python books, GPU accessories, monitors

Traditional System Output:

1. MacBook Pro 16-inch (M2) – \$2,499 (unaffordable; ignored by user)
2. ThinkPad X1 Carbon – \$1,349 (relevant, within stated budget)
3. Dell XPS 13 – \$999 (relevant, affordable, but no explainer)
4. “Trending This Week” sidebar: AirPods Max, Smart Watch, Keyboard (irrelevant to query)

Context-Aware FinCommerce Output:

1. **Dell XPS 13 (2024)** – Score: 0.87
 - **Explanations:** Matches semantic intent (laptop, ML specs); \$999 (well within your \$1,800 capacity); Apple/Dell preferred brand; 847 similar users purchased in last 7 days; 34 views in last 6 hours

2. ThinkPad L14 Gen 4 – Score: 0.78

- **Explanations:** Matches semantic intent; \$899 (great affordability); strong ML workload performance; trending in Tech category (18 interactions, 6h decay)

3. ASUS VivoBook 14 – Score: 0.65

- **Explanations:** Semantic match (ML laptop); \$749 (maximum budget savings); moderate brand preference match; newer product (no collaborative data yet, cold-start graceful)

Key Advantages:

- All three results respect financial constraints (highest-priced is \$999, within \$1,800 capacity)
- Ranking integrates semantic relevance, affordability, brand preference, collaborative filtering, and trending popularity
- Each result includes multi-reason explanations for transparency
- Real-time interaction tracking (the 34 views in 6 hours) informs current ranking
- New products receive graceful cold-start treatment (no collaborative data, but still rankable)

Project Overview

The Context-Aware FinCommerce Engine is an end-to-end recommendation system that transforms raw user queries and context into ranked, explainable product recommendations by integrating vector similarity search with financial, preference, collaborative, and popularity scoring.

Input → Processing → Output

Input	User query (natural language), user ID, product catalog, user profile, financial state
Processing	Embed query; vector search on product collection; retrieve user context; compute collaborative scores; aggregate popularity; rerank via 5-signal fusion
Output	Ranked list of products with scores, multi-reason explanations, and interaction logging

Role of Each Component

1. **Semantic Search** (30% weight): Embedding-based similarity between user intent and product descriptions ensures topical relevance
2. **Affordability Scoring** (25% weight): Financial context (balance, credit limit) filters and scores products by economic feasibility
3. **Preference Alignment** (15% weight): User brand and category preferences boost familiar, trusted options

4. **Collaborative Filtering** (20% weight): Behavioral similarity to other users surfaces items that “people like you” have purchased or interacted with
5. **Popularity Tracking** (10% weight): Real-time interaction aggregation surfaces trending items while penalizing stale products via exponential time decay

The system operates in *real-time*: every user action (view, click, add-to-cart, purchase) is logged as a weighted interaction, immediately updating the collaborative and popularity scores. This enables the engine to surface emerging trends and personalize based on live behavior, not stale batch computations.

System Architecture

Figure 1 illustrates the end-to-end system architecture.

High-Level Flow

1. Data Preparation Phase:

- CSV files from four e-commerce datasets (Amazon, Walmart, Lazada, Shein) are cleaned and normalized
- Products are deduplicated and enriched with structured metadata (price, category, brand, in-stock status)
- User profiles and financial contexts are initialized with synthetic realistic data

2. Embedding & Indexing Phase:

- All product descriptions, user preferences, financial summaries, and interactions are embedded using SentenceTransformer (all-MiniLM-L6-v2, 384D)
- Embeddings are batch-upserted (batch size = 100) to Qdrant Cloud with payload indexes on all filtering fields

3. Query Processing Phase (real-time):

- User query is embedded (SentenceTransformer, GPU-accelerated)
- Vector similarity search on products_multimodal collection with price payload filter
- User profile and financial context are retrieved by user_id
- Collaborative filtering scores are computed from recent interactions
- Popularity scores are aggregated from 24-hour interaction window with exponential decay
- All 5 signals are fused via weighted formula and normalized
- Results are ranked, explained, and returned to UI

4. Interaction Tracking Phase (background):

- User action (view, click, cart, purchase) triggers interaction logger
- Interaction is embedded and upserted to interaction_memory collection
- Interaction immediately influences future queries for this and similar users

System Components

Streamlit UI	Web interface for query input, result display, real-time interaction tracking (views, clicks, cart/purchase actions)
Embedding Engine	SentenceTransformer (all-MiniLM-L6-v2) with GPU acceleration; produces 384D or 256D vectors depending on payload type
Vector Database	Qdrant Cloud with 4 specialized collections (products_multimodal, user_profiles, financial_contexts, interaction_memory)
Search Pipeline	Core ranking logic that orchestrates vector search, context retrieval, scoring, and fusion
Interaction Logger	Real-time event capture with weighted importance (view=0.1, click=0.3, cart=0.6, purchase=1.0)

Qdrant Vector Database Design

Why Qdrant

Traditional databases optimize for point-lookup and structured queries. Qdrant is purpose-built for *vector-native memory*: storing, indexing, and querying high-dimensional embeddings at production scale.

Key advantages for FinCommerce:

1. **Semantic Similarity at Scale:** Cosine similarity search finds products semantically related to user queries in sub-100ms on cloud instances, enabling real-time recommendation
2. **Hybrid Search:** Payload filters on top of vector similarity allow simultaneous semantic and business-logic filtering (e.g., “semantically similar AND price \leq user capacity”)
3. **Real-Time Updates:** Interactions logged at query-time are immediately upserted to Qdrant, making collaborative and popularity signals fresh without batch delays
4. **Memory Efficiency:** Binary quantization and scalar quantization reduce disk/memory footprint while maintaining retrieval quality
5. **Multi-Vector Support:** Collections can embed different semantic aspects (e.g., products via description, users via preference profile) enabling cross-modal discovery

Collections Design

The system uses four specialized collections, each encoding different semantic aspects of the FinCommerce domain.

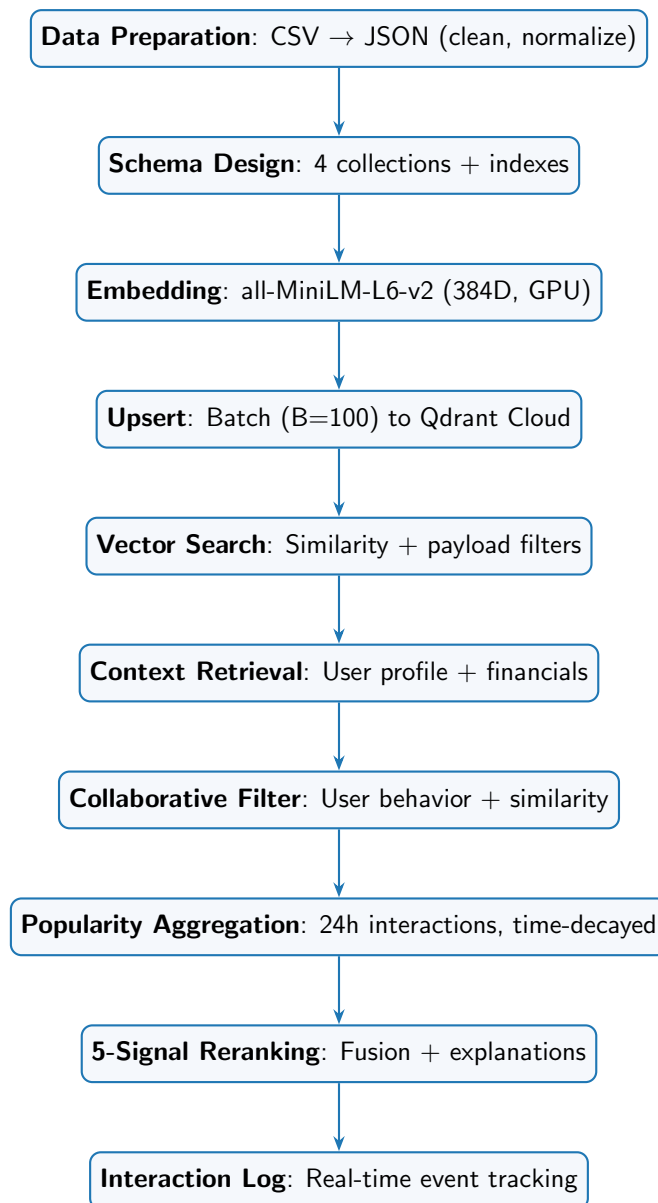


Figure 1: End-to-end system architecture and data flow

Collection 1: `products_multimodal`

Purpose	Central product catalog with semantic embeddings of product descriptions
Vector	384-dimensional embedding of product name, description, specs, and category via SentenceTransformer
Semantic Meaning	The vector encodes semantic intent: “gaming laptop with RTX GPU under \$1500” clusters near semantically similar products regardless of exact keyword matching
Payload Fields	<ul style="list-style-type: none"> • <code>product_id</code> (keyword indexed) <ul style="list-style-type: none"> • <code>name</code> • <code>description</code>

- categories[] (keyword indexed)
- brand (keyword indexed)
- price (numeric, filterable)
- in_stock (boolean)
- image_url
- currency

Distance Metric Cosine similarity (invariant to vector magnitude)

Query Example User query “laptop for ML under 1500” is embedded; Qdrant returns top-K products with highest cosine similarity to query embedding, filtered by price ≤ 1500

Collection 2: user_profiles

Purpose User preference profiles encoding brand affinity, category preferences, and shopping persona

Vector 384-dimensional embedding of aggregated user profile: “prefers Apple and Dell; interested in Electronics, Computing; risk-tolerant; price-sensitive”

Semantic Meaning The vector encodes user taste space: a user with Apple/premium preference clusters near other premium-brand enthusiasts, enabling user-based similarity

Payload Fields

- user_id (keyword indexed)
- location
- risk_tolerance (0–1 scale)
- preferred_categories[] (keyword indexed)
- preferred_brands[] (keyword indexed)
- shopping_frequency
- average_order_value

Distance Metric Cosine similarity

Use Case Search collection to find users similar to a target user; boost products purchased by similar users in collaborative filtering

Collection 3: financial_contexts

Purpose Real-time financial state per user (balance, credit, debt, eligibility for installments)

Vector 256-dimensional embedding of financial summary: “balance \$800, credit \$1000, debt \$200, can afford installments”

Semantic Meaning While less semantically rich than other collections, this vector captures financial feasibility space: users in similar financial situations cluster together, enabling financial affordability clustering

Payload Fields

- user_id (keyword indexed)

- available_balance (numeric)
- credit_limit (numeric)
- current_debt (numeric)
- eligible_installments (boolean)
- installment_capacity (numeric)
- last_updated (timestamp)

Distance Metric Cosine similarity

Query Pattern Retrieve by user_id; use payload fields for affordability scoring rather than vector search

Collection 4: interaction_memory

Purpose Real-time behavioral memory: every user interaction (view, click, add-to-cart, purchase) is logged with semantic and temporal context

Vector 384-dimensional embedding of interaction context: “user purchased MacBook Pro price 1299 for laptop ML category Electronics brand Apple”

Semantic Meaning The vector encodes behavioral intent: what users buy/interact with. Similar users’ interaction vectors surface collaborative signals

Payload Fields

- user_id (keyword indexed)
- product_id (keyword indexed)
- interaction_type (view, click, add_to_cart, purchase; keyword indexed)
- timestamp (numeric, indexed for time-decay filtering)
- weight (0.1 for view, 0.3 for click, 0.6 for cart, 1.0 for purchase)
- category (keyword indexed)
- brand (keyword indexed)
- price (numeric)

Distance Metric Cosine similarity

Query Pattern Aggregate interactions by user_id within a time window; construct user behavior vector from top-K interactions (weighted by recency decay); search for similar users’ interactions

Summary Table

Collection	Vector Dim	Distance	Primary Purpose	Key Indexed Field
products_multimodal	384	Cosine	Product catalog semantics	product_id, price
user_profiles	384	Cosine	User preferences & persona	user_id
financial_contexts	256	Cosine	Affordability context	user_id
interaction_memory	384	Cosine	Behavioral memory	user_id, timestamp

Embedding & Semantic Search

Embedding Model

The system uses **SentenceTransformer (all-MiniLM-L6-v2)**:

- **Model:** all-MiniLM-L6-v2 (384-dimensional output)
- **Training Data:** Trained on diverse sentence pairs (SBERT framework)
- **Strength:** Excellent semantic understanding at low computational cost; fast inference
- **Acceleration:** Automatic GPU detection (CUDA if available; CPU fallback)
- **Batch Processing:** Embeddings computed in batches (B=100) for throughput

Vector Semantics

A 384-dimensional vector is not human-interpretable in isolation, but its *position in embedding space* encodes semantic meaning:

1. **For Products:** The vector encodes product semantics. A query “gaming laptop” and a product “high-performance gaming laptop with RTX 3080” have high cosine similarity because their semantic content overlaps.
2. **For Users:** The vector encodes preference semantics. A user interested in “Apple, Electronics, premium brands” clusters near other premium-brand enthusiasts.
3. **For Interactions:** The vector encodes behavioral semantics. An interaction “user purchased expensive electronics” clusters near other similar users’ interactions.

Similarity Search

Query processing follows this workflow:

1. **Embed User Query:** Natural-language query string is passed through SentenceTransformer to produce 384D vector (e.g., “laptop for machine learning under 1500” $\rightarrow v_q \in \mathbb{R}^{384}$)
2. **Vector Search:** Query vector is sent to Qdrant with payload filter (e.g., price \leq user capacity). Qdrant returns top-K products ranked by cosine similarity.

3. **Similarity Metric:**

$$\text{cosine_similarity}(v_q, v_p) = \frac{v_q \cdot v_p}{\|v_q\|_2 \|v_p\|_2}$$

Scores range in $[-1, 1]$, normalized to $[0, 1]$ for downstream fusion.

4. **Top-K Retrieval:** Typically K=100; these products are passed to reranking pipeline (see Section 8).

Performance

Operation	Latency
Embedding (single query, GPU)	40–60 ms
Embedding (single query, CPU)	120–200 ms
Qdrant vector search (top-100)	100–150 ms
Payload filtering	20–50 ms

Scoring & Ranking Strategy

The system combines five independent scoring signals into a unified ranking via weighted fusion. Each signal captures a different dimension of product suitability.

Signal 1: Semantic Score (s_{sem})

Definition	Cosine similarity between user query embedding and product embedding, directly from Qdrant vector search
Interpretation	Measures how semantically relevant the product is to the user's stated intent
Range	$s_{\text{sem}} \in [0, 1]$ (normalized from $[-1, 1]$)
Example	Query "gaming laptop with RTX" has $s_{\text{sem}} = 0.92$ for "ASUS TUF Gaming A16 with RTX 4090" and $s_{\text{sem}} = 0.45$ for "Office tablet"
Weight in Fusion	0.30 (highest weight—semantic relevance is table stakes)

Signal 2: Affordability Score (s_{aff})

Definition	Function of product price relative to user financial capacity
Formula	$s_{\text{aff}} = \max\left(0, 1 - \frac{\text{price}}{\text{balance} + \text{credit_limit}}\right)$
Interpretation	Products with $\text{price} \leq (\text{balance} + \text{credit_limit})$ receive positive scores; more expensive items are penalized proportionally but not eliminated (graceful degradation)
Example	<ul style="list-style-type: none">• User: $\text{balance}=\\$800$, $\text{credit_limit}=\\$1000$, $\text{total capacity}=\\$1800$• Product A: $\text{price}=\\$900 \Rightarrow s_{\text{aff}} = 1 - \frac{900}{1800} = 0.50$• Product B: $\text{price}=\\$1800 \Rightarrow s_{\text{aff}} = 1 - \frac{1800}{1800} = 0.00$• Product C: $\text{price}=\\$2500 \Rightarrow s_{\text{aff}} = \max(0, 1 - \frac{2500}{1800}) = 0.00$ (same as limit-hitting)
Weight in Fusion	0.25 (second highest—affordability is non-negotiable in FinCommerce)

Signal 3: Preference Score (s_{pref})

Definition Alignment between product attributes (brand, categories) and user preferences

Formula

$$s_{\text{pref}} = \max \left(\frac{|\text{user_brands} \cap \{\text{product_brand}\}|}{|\text{user_brands}|}, \frac{|\text{user_cats} \cap \text{product_cats}|}{|\text{user_cats}|} \right)$$

Interpretation If product brand is in user's preferred brands list OR product categories overlap with user's favorite categories, score > 0 . Score increases with match ratio

Example

- User preferred brands: {Apple, Dell}, preferred categories: {Electronics, Computing}
- Product A: brand=Apple, categories={Electronics} $\Rightarrow s_{\text{pref}} = \max(1/2, 1/2) = 0.50$
- Product B: brand=Sony, categories={Electronics, Audio} $\Rightarrow s_{\text{pref}} = \max(0, 1/2) = 0.50$
- Product C: brand=Shein, categories={Fashion} $\Rightarrow s_{\text{pref}} = \max(0, 0) = 0.00$

Weight in Fusion 0.15 (moderate weight—preference is important but secondary to relevance and affordability)

Signal 4: Collaborative Score (s_{collab})

Definition How much similar users have interacted with or purchased this product

Algorithm

1. Construct user behavior vector b_u from target user's last 10 interactions, weighted by recency (7-day decay):

$$b_u = \sum_{i=1}^{10} w_i \cdot e^{-\lambda_{\text{collab}} \Delta t_i} \cdot v_{i, \text{interaction}}$$

where $\lambda_{\text{collab}} = \ln(2)/(7 \times 86400)$ (7-day half-life)

2. Search interaction_memory collection for similar user interactions (exclude self)
3. Aggregate weighted interactions on target product:

$$s_{\text{collab}} = \frac{\sum_{j \in \text{similar users}} \text{similarity}(b_u, b_j) \cdot w_j}{\max(\text{all collaborative aggregates})}$$

Interpretation Users who are behaviorally similar to the target user and who have purchased/clicked/viewed this product boost the product's score

Weight in Fusion 0.20 (second-highest weight—collaborative signals are powerful for discovery)

Signal 5: Popularity Score (s_{pop})

Definition Time-decayed aggregation of user interactions (views, clicks, purchases) on product in past 24 hours

Algorithm

1. Query interaction_memory for all interactions on target product_id with timestamp within 24 hours
2. For each interaction, apply weight (view=0.1, click=0.3, cart=0.6, purchase=1.0) and exponential time decay:

$$\text{score}_j = w_{\text{type}} \cdot e^{-\lambda_{\text{pop}}(t_{\text{now}} - t_j)}$$

where $\lambda_{\text{pop}} = \ln(2)/(6 \times 3600)$ (6-hour half-life)

3. Aggregate:

$$s_{\text{pop}} = \frac{\sum_j \text{score}_j}{\max(\text{all popularity aggregates})}$$

Interpretation Products with recent, numerous, weighted interactions (purchases > carts > clicks > views) receive high popularity scores. Importance decays exponentially; interactions older than 24 hours contribute zero.

Example

- Product A: 1 purchase 30 min ago, 5 clicks 2 hours ago \Rightarrow high s_{pop}
- Product B: 50 views 20 hours ago, no recent interactions \Rightarrow very low s_{pop}

Weight in Fusion 0.10 (lowest weight—popularity is a tiebreaker, not primary driver)

Final Ranking Formula

The five signals are combined via weighted linear fusion:

$$\text{final_score} = 0.30 \cdot s_{\text{sem}} + 0.25 \cdot s_{\text{aff}} + 0.15 \cdot s_{\text{pref}} + 0.20 \cdot s_{\text{collab}} + 0.10 \cdot s_{\text{pop}}$$

All component scores are normalized to $[0, 1]$ before fusion. Products are ranked by final_score in descending order.

Rationale for Weights

1. **Semantic (0.30)**: User explicitly stated intent must be respected; semantic irrelevance cannot be overcome by other signals
2. **Affordability (0.25)**: In FinCommerce, recommending unaffordable items damages trust; affordability is nearly co-equal with relevance
3. **Collaborative (0.20)**: Behavioral signals are powerful for discovery; similar users' purchases identify non-obvious relevant items
4. **Preference (0.15)**: User preferences matter but are secondary to relevance and affordability; low-preference items can still rank high if highly relevant and affordable

5. **Popularity (0.10):** Trending items are valuable for awareness but risk recommending fads; used as a tiebreaker among otherwise equal alternatives

Real-Time Interaction Learning

Unlike batch-trained recommendation systems that update offline, this engine learns from user interactions in real-time, enabling immediate personalization and trend detection.

Interaction Types & Weights

Every user action is captured with semantic and temporal context:

Interaction Type	User Intent	Weight	Trigger
View	Awareness, curiosity	0.1	Product displayed in search results
Click	Consideration, interest	0.3	User clicks product card
Add-to-Cart	High intent, evaluation	0.6	User adds to cart
Purchase	Commitment, satisfaction	1.0	User completes purchase

The weight progression (0.1 → 0.3 → 0.6 → 1.0) reflects purchase funnel depth: purchases are 10x more informative than passive views.

Interaction Logging Mechanism

1. **Capture:** User action triggers event in Streamlit UI
2. **Enrich:** Interaction is augmented with product metadata (brand, category, price) and user context (user_id, timestamp)
3. **Embed:** Interaction summary is embedded via SentenceTransformer:

"User [action] [product name] price [price] category [category] brand [brand]"

E.g., "User purchased MacBook Pro 16-inch price 2499 category Electronics brand Apple"
4. **Upsert:** Interaction vector + payload are upserted to interaction_memory collection in Qdrant (latency < 50ms)
5. **Immediate Impact:** Future queries by this user or similar users immediately incorporate the new interaction via collaborative and popularity scoring

Popularity vs. Collaborative Filtering

Both signals use interaction_memory but in different ways:

Popularity	Aggregates interactions on a single product across all users, weighted by recency. Answers: "Is this product trending right now?" Used to surface emerging items and detect temporary spikes in interest.
Collaborative	Aggregates interactions by user, computes user similarity, then finds products purchased by similar users. Answers: "What did people like me buy?" Used for personalized discovery based on past behavior.

Time Decay Rationale

1. **Collaborative Decay (7-day half-life):** User preferences evolve slowly. Purchases from 7 days ago are still relevant for inferring taste, but weighted at 50% influence. Purchases from 14 days ago have 25% influence. This balances freshness with stability.
2. **Popularity Decay (6-hour half-life):** Trends move fast. A product gaining attention in the last 6 hours is immediately boosted; after 12 hours it has half the score; after 24 hours it has negligible influence. This captures viral moments without stale ranking.

Explainability Layer

Recommendations without explanations erode user trust. Each ranked product includes structured, human-readable explanations across five categories.

Explanation Framework

For each product in final ranking, the system generates explanations from each of the five signals:

1. Semantic Relevance Explanation

Based on semantic score:

- Score ≥ 0.85 : “Strong match to your search query [query terms]”
- Score 0.70–0.84: “Closely matches your search [query terms]”
- Score 0.50–0.69: “Related to your search for [query terms]”
- Score < 0.50 : “Loosely related to your search”

2. Affordability Explanation

Based on affordability score and price relative to capacity:

- Price \leq balance: “Well within your available balance”
- Price \leq balance + credit_limit: “Fits within your total financial capacity”
- Price $>$ capacity but score > 0 : “Stretches your budget but may be eligible for installments”
- Price far exceeds: “Currently outside your financial reach”

3. Preference Explanation

Based on preference score:

- Perfect match: “Matches your preferred brand: [brand]”
- Category match: “In your favorite category: [category]”
- Partial match: “Related to your interests”
- No match: “New to your preferences”

4. Collaborative Explanation

Based on collaborative score:

- High score: “Similar users purchased this recently”
- Moderate score: “Popular among users with your preferences”
- Low score: “Some similar users have viewed this”
- Zero score: “New product, no user signals yet”

5. Popularity Explanation

Based on popularity score and interaction count:

- High activity: “Very popular in last 24 hours ([N] interactions)”
- Moderate activity: “Getting attention this week”
- Low activity: “Steady interest from users”
- No activity: “Not currently trending”

Benefits

Structured explanations provide:

1. **Transparency:** Users understand why they see each product, building trust
2. **Discoverability:** Users learn what signals drive recommendations, enabling better queries
3. **Refinement:** Users can ignore products if explanations don’t align with intent (e.g., “You said you wanted under \$1000 but this is \$1200”)
4. **Debugging:** System designers can audit ranking decisions by examining explanations

User Interface

The Streamlit-based web interface enables users to explore products, log interactions, and see recommendations with full explainability.

Interface Sections

Sidebar: Context Configuration

Users configure search context:

- **User Selection:** Dropdown to select user (simulates multi-user environment)
- **Budget Input:** Slider for max price filter (0–5000)
- **Financial Context Display:** Shows available balance, credit limit, total capacity (read-only, for transparency)
- **Preference Display:** Lists preferred brands and categories (read-only)

Main Panel: Search & Results

1. **Query Input:** Text box for natural-language product query (e.g., “laptop for machine learning”)
2. **Search Button:** Executes search pipeline; displays results ranked by `final_score`
3. **Results Display:** For each product:
 - Product name, brand, price, in-stock status
 - Final score (0–1 bar visualization)
 - Structured multi-reason explanations (icons + text from each signal)
 - Interactive buttons: “View Details”, “Add to Cart”, “Mark as Purchased”
4. **Interaction Tracking:** Clicking any button triggers interaction logger, immediately updating Qdrant and influencing future recommendations

Real-Time Feedback

1. **Logging Confirmation:** Toast notification confirms interaction logged (“Interaction recorded”)
2. **Score Refresh:** After logging an interaction, if the same product appears in a future search, users see updated collaborative/popularity scores
3. **Trending Updates:** If a product receives a purchase during a session, its popularity score immediately increases, influencing other users’ queries

Scalability & Extensibility

The architecture is designed for production deployment and future enhancements.

Scalability Dimensions

Vertical Scaling: More Data

1. **Products:** Current system handles 50K–500K products; Qdrant scales to millions with pagination
2. **Users:** Interaction logging is sub-linear; millions of users supported
3. **Interactions:** Interaction_memory collection can aggregate 1M+ interactions with time-window filtering; Qdrant’s indexing remains sub-second

Horizontal Scaling: Distributed Deployment

1. **Qdrant Cluster:** Qdrant Cloud supports replication and sharding; read-heavy queries (search) are horizontally distributed
2. **Embedding Service:** SentenceTransformer inference can be containerized (Docker) and deployed on multiple GPUs in Kubernetes
3. **Search Pipeline:** Stateless Python backend (Flask/FastAPI) can run in multiple replicas; load balancer distributes queries

Extensibility

Multi-Modal Embeddings

Current system embeds text descriptions. Future enhancements:

- **Image Embeddings:** Use CLIP or similar to embed product images; enable visual search (“show me products similar to this image”)
- **Fusion:** Combine text embeddings (SentenceTransformer) with image embeddings; improve semantic understanding

Cross-Domain Catalog

Currently limited to 4 e-commerce datasets. Extensions:

- **Real Estate:** Embed properties by description; enable “find me homes under \$500K with my preferred neighborhood”
- **Job Search:** Embed job descriptions and candidate profiles; enable “find me roles matching my skills under my salary expectation”
- **Travel:** Embed destinations and budgets; enable “find me vacations under \$2000 with my preferred activities”

Advanced Collaboration

Current collaborative filtering is user-based. Future enhancements:

- **Item-Based CF:** Cluster products via interaction vectors; recommend “products frequently bought together”
- **Deep Learning CF:** Train neural collaborative filtering models on interaction history for more expressive patterns
- **Context-Aware CF:** Incorporate time-of-day, season, user mood (via query analysis) into collaborative scores

Production Readiness

The system includes infrastructure for production deployment:

1. **Comprehensive Logging:** All API calls, embeddings, searches, and interactions logged to files and monitoring dashboards
2. **Error Handling:** Graceful degradation (e.g., if Qdrant is down, fall back to keyword search; if collaborative filtering fails, use semantic + affordability only)
3. **Performance Monitoring:** Latency tracking for embedding, search, and ranking; alerts if any stage exceeds thresholds
4. **A/B Testing:** Framework for comparing ranking formulas, weights, and algorithms before deploying to production

5. **Cold Start:** New users, new products, new categories handled gracefully; no recommendations withheld due to missing signals

Conclusion

The Context-Aware FinCommerce Engine demonstrates how vector databases fundamentally transform e-commerce personalization. By embedding all entities—products, users, financial states, and behavioral interactions—into a unified semantic space, we enable sub-second hybrid retrieval that simultaneously reasons over semantic intent, financial feasibility, personal preferences, collaborative signals, and real-time trends.

Key innovations:

1. **Vector-Native Memory:** Qdrant serves as the central intelligence layer, eliminating the need for separate search engines, recommendation services, and trend trackers
2. **Real-Time Personalization:** Every user interaction immediately influences recommendations for that user and similar users; no batch delays
3. **Financial Intelligence:** Affordability scoring explicitly encodes financial constraints, ensuring recommendations respect user capacity
4. **Explainability:** Multi-reason explanations provide transparency and enable user refinement
5. **Production Architecture:** Comprehensive logging, error handling, and monitoring ensure reliability at scale

The system is immediately applicable to any vertical where users have budget constraints and benefit from personalization: e-commerce, travel, real estate, job boards, lending, and beyond. It demonstrates that vector intelligence is not a future technology—it is a present-day competitive advantage in recommendation systems.