

Vespa.ai

Simplifying AI Application Development and Deployment in Fintech with Vespa.ai

A Guide For Managers

Vespa.ai

Vespa.ai is a platform for building and running real-time AI-driven applications for search, recommendation, personalization, and RAG. It enables enterprise-wide AI deployment by efficiently managing data, inference, and logic, handling large data volumes and over 100K queries per second. Vespa supports precise hybrid search across vectors, text, and structured metadata. Available as both a managed service and open source, it's trusted by organizations like Spotify, Vinted, Wix, and Yahoo. The platform offers robust APIs, SDKs for integration, comprehensive monitoring metrics, and customizable features for optimized performance.

Document

Simplifying AI Application Development and Deployment in Fintech with Vespa

Piotr Kobziakowski

Senior Principal Solutions Architect

01.01.2025

Contents

➤ Introduction	3
Investment Platforms: A New Age of Accessibility and Control	
Banking and Beyond: The Digital-First Revolution	
Traditional Banking: Not Staying Behind	
➤ GenAI and ML in Growth and Customer Retention	5
Vespa's AI Platform: Enabling a Data-Driven Fintech and Banking Future	
➤ What is Vespa	6
Vespa Use Cases in Banking and Fintech	
➤ Vespa Use Cases in Banking and Fintech	7
Machine Learning Embeddings and Meaning of Vectors (Tensors)	
➤ Key AI Use Cases	9
➤ Deployment Considerations	12
Scalability and Infrastructure	
Data Privacy and Security	
Integration with Existing Systems	
➤ Summary	13

Introduction

The fintech industry has rapidly transformed, reshaping investment and banking by offering streamlined, mobile-first services that put financial control directly into users' hands. Through lowered barriers to entry, innovative features, and technology-driven personalization, fintech platforms have redefined finance across two primary domains: investment and banking. Meanwhile, traditional banks have also been evolving, focusing on internal projects that use automation and augmentation to enhance day-to-day operations. This balance of external innovation and internal modernization allows both fintech firms and established banks to provide more responsive, efficient financial services.

Investment Platforms: A New Age of Accessibility and Control

Investment-focused fintech platforms have democratized access to financial markets, making investing accessible and affordable for a broader audience. These platforms provide professional-grade investment management without the high costs associated with traditional financial advisors through commission-free trading and automated investing with robo-advisors. These platforms deliver tailored insights and actionable advice using AI. AI-driven recommendations help users stay informed about market trends, while personalized investing tips and savings strategies guide users based on their financial habits. This democratization empowers individuals with tools and insights once reserved for institutional investors, making personal investing more inclusive.

Banking and Beyond: The Digital-First Revolution

Digital-first banking platforms have made traditional financial services more efficient, transparent, and customer-centric. These solutions integrate essential banking functions with innovative tools like international currency exchange and budgeting features, automated savings, and many other services that allow users to manage their finances, investments, travel, and insurance with a single app. These platforms thrive by leveraging AI to enhance customer experience and operational efficiency. For example, Traditional ML algorithms improve fraud detection, providing real-time alerts on unusual activity, while GenAI-powered insights enable users to manage spending and savings more effectively. Customer-centric design supported by AI-enabled this platform to offer control, transparency, and financial well-being at users' fingertips.

Traditional Banking: Not staying behind

In traditional banking, institutions are also undergoing significant transformation, focusing on internal projects aimed at modernizing their operations. Automation and GenAI-driven work augmentation play critical roles in enhancing day-to-day functions such as customer service, transaction processing and risk management. Additionally, banks are leveraging advanced technologies for Anti-Money Laundering (AML) and fraud detection, strengthening their ability to identify and mitigate financial crime to stay compliant. Internal knowledge sharing and document digitalization initiatives are streamlining workflows, enabling faster information retrieval and improved decision-making. Optimizing these core activities enables traditional banks to not only improve operational efficiency but also stay competitive in a landscape increasingly shaped by agile fintech platforms. This internal evolution complements their external offerings, ensuring that established banks continue to meet evolving customer expectations while refining their internal processes.

The Role of GenAI and Traditional ML in Growth and Customer Retention

The integration of GenAI and traditional ML technologies is crucial in driving growth, customer loyalty, and operational efficiency across both fintech and traditional banking. For fintech platforms, personalized experiences – such as real-time insights, easy access to data, fraud alerts, and tailored investment advice – build trust and deepen user engagement. These GenAI-powered platforms can anticipate user needs and deliver relevant services, which enhances customer retention and lifetime value.

In traditional banking, GenAI and traditional ML are equally transformative, bringing efficiencies to core functions like transaction processing, regulatory compliance, risk management. Automation in areas like Anti-Money Laundering (AML) and fraud detection allows banks to proactively identify and mitigate risks, while advanced algorithms streamline credit scoring and underwriting, speeding up decision-making without compromising accuracy. Document digitalization and internal knowledge sharing, powered by GenAI, enable faster access to information, reducing manual work and improving internal productivity.

Vespa's AI Platform:

Enabling a Data-Driven Fintech and Banking Future

Vespa's AI platform, with its rich capabilities, significantly enhances fintech operations. Vespa supports high-speed, real-time analytics to generate personalized insights based on individual investment patterns. In banking, Vespa's data processing powers fraud detection and personalized financial recommendations, enhancing both security and user experience. Vespa's ability to manage vast amounts of data from disparate sources and in multiple formats efficiently allows fintech companies to respond dynamically to market conditions and customer needs, ultimately creating a secure, adaptive, and engaging financial ecosystem.

By embracing GenAI and traditional ML with AI platforms like Vespa, fintech firms are well-equipped to meet the expectations of today's digitally savvy consumers, providing a seamless, tailored, and responsive experience that sets them apart in a rapidly evolving financial landscape. Before we dive deeper into Vespa let's step back and talk about Vespa's history.

What is Vespa?

Vespa's origins trace back to the 1990s when it was initially developed within Yahoo to meet the demands of large-scale data processing. It was built to serve millions of users around the world at a time when hardware wasn't as powerful as it is today, making it a robust, highly optimized solution for versatile use cases. By 2017, Vespa was released as an open-source project, and in 2023, a company was formed to focus on its growth and development.

With decades of experience behind it, Vespa.ai is a mature technology that remains highly relevant for modern machine learning and AI applications. Its core strength lies in real-time AI and data serving, enabling it to handle large datasets and deliver insights with minimal latency efficiently. This makes it exceptionally well-suited for the most complex AI-driven use cases of today.

Vespa's initial reach is also impressive. It powers more than 150 different projects within Yahoo. It has been adopted by major organizations outside Yahoo, such as Perplexity.ai, BigData.com, Spotify, Farfetch, Otto, and many more in e-commerce, entertainment and banking space. Vespa.ai is adopted across different projects, including groundbreaking showcases. Its versatility and reliability in handling various challenges, whether powering search engines, recommendation systems, AI solutions or real-time data analytics.

Built from the ground up to be scalable and flexible, Vespa continues to evolve alongside the demands of modern enterprises. Its ability to seamlessly handle high user volumes while processing complex data makes it an ideal platform for businesses looking for a solution that can grow with them.

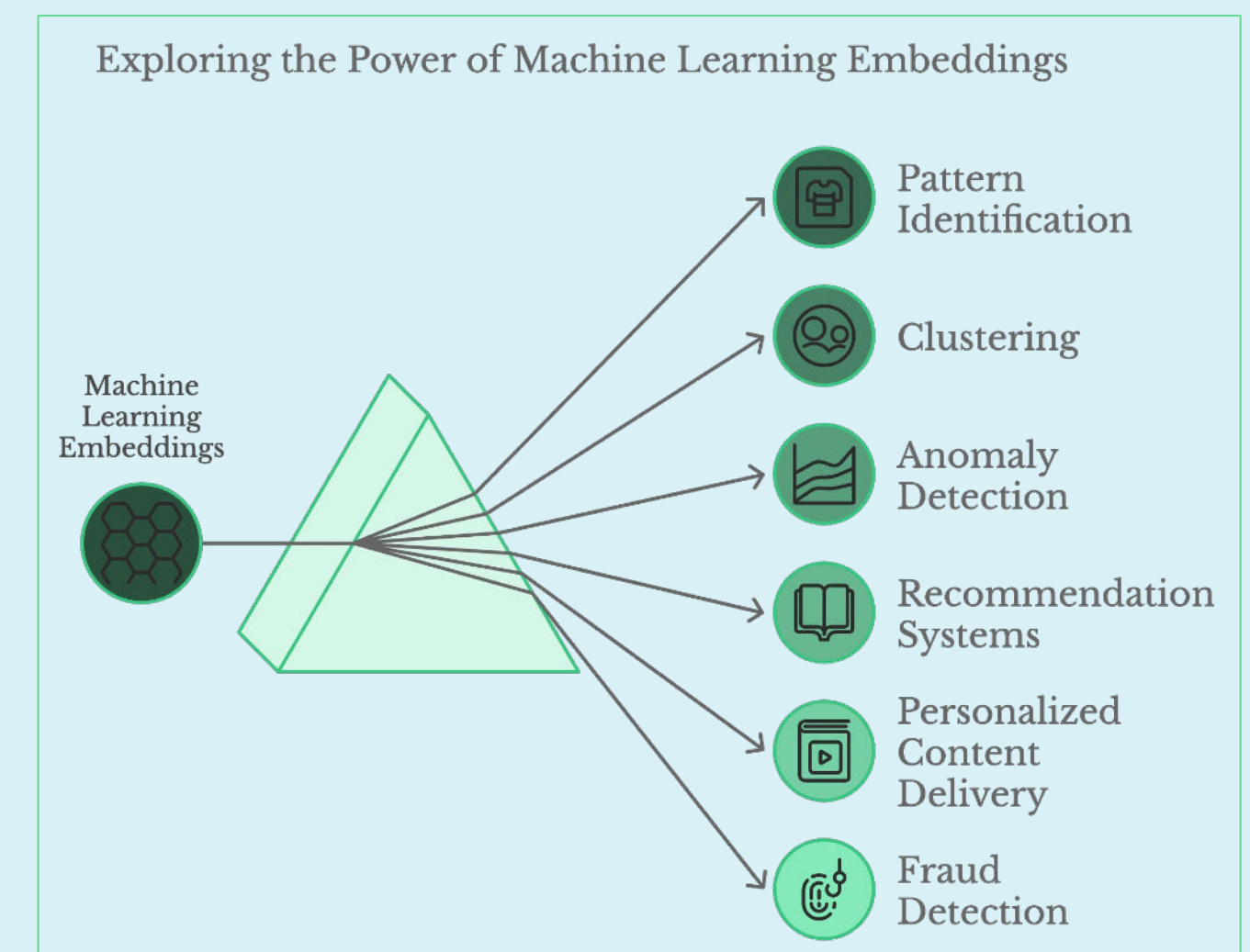
Vespa Use Cases in Banking and Fintech

In the fintech and banking sectors, Vespa's advanced vector search and tensor operations are revolutionizing applications ranging from enhancing customer interactions to strengthening security measures. While Vespa excels in delivering personalized search and recommendation systems that elevate customer experiences, its capabilities extend beyond that to sophisticated data analytics tasks, especially when combined with different available ML models. Vespa's diverse toolset can be used in data analytics tasks that are little less known. Let's start by explaining the base concepts.

Machine Learning Embeddings and Meaning of Vectors (Tensors)

Machine learning model embeddings transform complex data into high-dimensional vector space that represent concepts, objects, and relationships in a way that algorithms can efficiently process. By encoding the essence of textual, visual, or other data types into these embeddings, the Vespa.ai engine can perform operations in vector space to identify patterns, relationships, and similarities. This capability extends beyond basic document retrieval to more advanced applications like clustering similar items, detecting outliers, and powering recommendation systems by identifying related content. Embeddings also enable personalized content delivery, fraud detection, and intent recognition, making them versatile tools across industries such as e-commerce, finance, and healthcare. Their strength lies in capturing contextual features and enables similarity calculation, measuring how close or far apart

different concepts are while maintaining computational efficiency. This allows to solve complex problems by leveraging the power of vector representations.



Key AI Use Cases

Vespa.ai's advanced vector search and tensor operations enable financial institutions to deliver personalized services and make informed decisions

Fraud Detection and Prevention

Advanced Search Capabilities	<p>Vespa.ai supports late interaction models such as colPali. It enables search beyond the traditional approach and is able to navigate charts, tables, images such as signatures or handwritten text. Finding all documents signed by a specific person, finding insider trading by identifying images captured by a camera, or finding exact documents using text search for specific content, including image description, table content, chart values or direct image search is now possible.</p>
Recommendation Engines for Banking Products	<p>Converting user profiles, transaction histories, and product features into vector representation will enable Vespa.ai to perform real-time matching of customers with suitable financial products through nearest neighbor search. Its distributed architecture ensures scalability for large user bases and extensive product catalogs.</p>
Personalized Investment Recommendations	<p>Vespa.ai can vectorize customer profiles, risk tolerance and investment history to dynamically align investment options with individual customer vectors. This system adapts recommendations in real time as customer data and market conditions evolve.</p>

User Experience

Dynamic Credit Scoring	<p>Incorporating non-traditional data sources into customer embeddings, Vespa.ai assesses the similarity between customer vectors and profiles of known credit risks. This approach provides up-to-date credit assessments with low latency.</p>
-------------------------------	--

Key AI Use Cases

Payment processing and Optimization

Fraud Detection in Payments	By vectorizing payment attributes, Vespa.ai can identify anomalies through comparison with typical transaction vectors, enabling real-time risk mitigation by blocking or flagging suspicious payments.
Dynamic Merchant Scoring	Vespa.ai can create embeddings using different models based on merchants' transaction histories and behaviors, assessing transaction legitimacy through vector similarity. This proactive approach helps reduce chargebacks by identifying high-risk merchants.

Fraud Detection and Prevention

Outlier Detection in Transactions	By representing transactions as vectors that capture features like amount, location, merchant, and time, Vespa.ai can establish a baseline of normal transactions and identify deviations from these patterns using distance metrics, effectively enabling it to flag potential anomalies.
Synthetic Identity Fraud Detection	By vectorizing user behaviors across multiple sessions, Vespa.ai can identify inconsistencies indicative of synthetic identities. It clusters similar behaviors to detect outliers, thereby uncovering fraudulent accounts.
Transaction Clustering	Vespa.ai can cluster transactions based on vector similarities. This can be used to group transactions in the spending type or even identify coordinated fraudulent activities.

Key AI Use Cases

Risk Assessment and Management

Risk Profile Analysis	By creating vector embeddings from customers' financial behaviors, Vespa.ai enables continuous assessment of risk profiles through similarity searches. This dynamic analysis allows for personalized financial advice tailored to individual risk profiles.
Creditworthiness Evaluation	Vespa.ai integrates both traditional and alternative data into multidimensional embeddings, facilitating similarity searches that compare applicant vectors with profiles of various credit outcomes. This approach provides underwriters with data-driven insights, supporting informed decision-making.

Compliance and Regulatory Reporting

Anti-Money Laundering (AML)	By vectorizing transactional data, Vespa.ai can identify patterns indicative of money laundering. Its rapid querying capabilities enable real-time monitoring and prompt detection of anomalies, ensuring efficient maintenance and updating of compliance checks.
Know Your Customer (KYC) Optimization	Vespa.ai can store customer documents as vector embeddings, allowing for automatic flagging of discrepancies by comparing document vectors to profile vectors. This process streamlines verification efforts, reducing the need for manual intervention.

Customer Service and Engagement

Automated Financial Assistance through Chatbot	Vespa.ai can be used for a highly scalable backbone of any chatbot. Its capabilities enable organizations to quickly build chatbots that can deliver accurate responses through advanced ranking and retrieval techniques implemented using agentic architecture.
---	---

Key AI Use Cases

Credit and Loan Management

Loan Default Prediction	Behavioral embeddings from borrower actions and external factors can be transformed into elector representation, enabling continuous monitoring by comparing current behavior vectors to default patterns. This approach allows for preventive actions when increased risk is detected.
Portfolio Segmentation	By vectorizing loan repayment data, Vespa.ai can be used to cluster borrowers with similar behaviors, facilitating targeted strategies and interventions for specific segments.

Data Analytics and Business Intelligence

Customer Segmentation and Profiling	By embedding demographic and behavioral attributes into vectors, Vespa.ai clusters customers for targeted marketing, enhancing personalization through tailored offers.
Churn Prediction	Vespa.ai represents customer engagement and satisfaction levels as vectors, detecting at-risk customers through pattern analysis and facilitating retention strategies to reduce churn.

Please note: the vectorization process requires using ML models built for the purpose.

Deployment Considerations

Scalability and Infrastructure

Vespa Cloud is designed to handle large-scale applications efficiently, offering features like automatic scaling to adjust resources based on real-time demand. This ensures optimal performance and cost-effectiveness by allocating resources as needed. Its distributed architecture supports seamless scaling, allowing applications to manage increasing data volumes and user requests without compromising performance. Additionally, Vespa Cloud provides continuous deployment and upgrades, enabling applications to evolve and scale smoothly.

Data Privacy and Security

Vespa Cloud prioritizes data privacy and security through several key measures:

- **Mutual TLS (mTLS) Authentication:** All communication between services is protected using mTLS, ensuring that only authenticated clients can access endpoints and that services communicate with trusted sources. (Vespa Blog)
- **Data Encryption at Rest:** All customer data is encrypted at rest using the cloud provider's native encryption capabilities (AWS KMS or Google Cloud KMS). (Vespa Cloud)
- **Access Control and Service Isolation:** Vespa Cloud employs strict access control mechanisms and service isolation to prevent unauthorized interactions between different applications and services. (Vespa Cloud)

These measures collectively ensure that Vespa Cloud maintains a secure environment for data processing and storage.

Data Privacy and Security

Vespa is designed to easily integrate existing technology stacks through its flexible APIs. Businesses can adopt Vespa incrementally without disrupting ongoing operations or requiring a complete re-platforming effort.

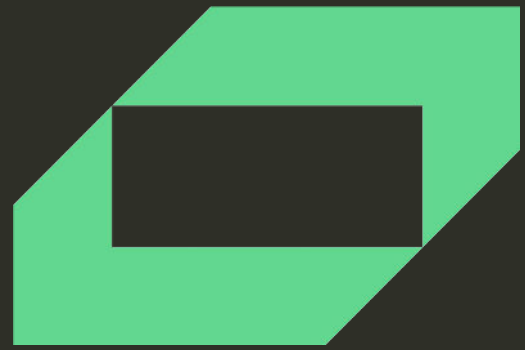
Summary

The fintech industry is transforming rapidly, blending innovative digital-first solutions and traditional modernization to deliver seamless, personalized financial services. Fintech platforms are democratizing investment and banking, offering tools like robo-advisors, AI-driven insights, and fraud detection. At the same time, traditional banks are leveraging AI and automation to enhance operational efficiency and compliance, balancing external competition with internal transformation.

Vespa.ai can play a pivotal role in this evolution, providing a robust AI-powered platform that supports advanced vector search, real-time analytics, and personalized recommendations. Its ability to handle vast datasets efficiently enables applications like fraud detection, customer profiling, and credit risk management. Vespa's versatile architecture integrates seamlessly with existing systems, ensuring scalability, data security, and cost-efficiency.

Vespa.ai's transformative potential in fintech and banking, is proven by its history, core features, and groundbreaking use cases, showcasing how it powers innovation and drives success in today's competitive financial landscape.

13



About Vespa.ai

Vespa.ai is a platform for building and running real-time AI-driven applications for search, recommendation, personalization, and RAG. It enables enterprise-wide AI deployment by efficiently managing data, inference, and logic, handling large data volumes and over 100K queries per second. Vespa supports precise hybrid search across vectors, text, and structured metadata. Available as both a managed service and open source, it's trusted by organizations like Spotify, Vinted, Wix, and Yahoo. The platform offers robust APIs, SDKs for integration, comprehensive monitoring metrics, and customizable features for optimized performance.

Interested to learn more? We have many different resources and information available through our social platforms

[GitHub](#)

[Twitter](#)

[LinkedIn](#)

[YouTube](#)