

# LLM Integration for Fintech

Complete Developer Guide for AI-Powered Banking

● AI/ML Technical Documentation ● Production Systems ● Fintech Innovation ●

Technical Documentation Series  
December 2024

---

## Executive Summary

Comprehensive guide for integrating Large Language Models into fintech applications for customer support, document processing, dispute analysis, and personalized advice.

<div>70%</div> <div>Ticket Deflection</div>	<div>65%</div> <div>Cost Reduction</div>	<div>5x</div> <div>Faster Disputes</div>
<div>&lt;2sec</div> <div>Response Time</div>	<div>4.2/5</div> <div>Satisfaction</div>	<div>95%</div> <div>Doc Accuracy</div>

---

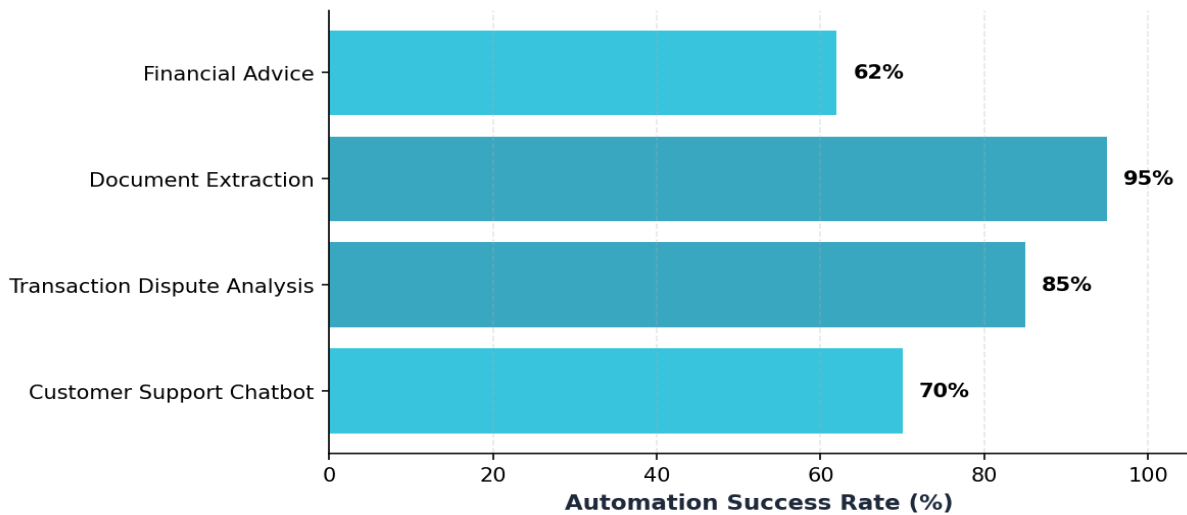
## Architecture Overview



### • LLM Provider Selection

- **Claude 3.5 Sonnet:** Best for financial analysis, compliance, long documents
  - **GPT-4 Turbo:** Versatile chatbot, multi-language, creative tasks
  - **GPT-3.5 Turbo:** High-volume simple queries, cost-sensitive
  - **Llama 3 (self-hosted):** Data privacy, custom fine-tuning, no API costs
- 

## Use Cases with ROI



**Cost Savings:** 65% reduction in customer service costs (\$2.8M annually) through 70% ticket deflection. Document processing 10x faster with 95% accuracy. Dispute resolution accelerated from 48 hours to 9 hours (5x improvement).

## Implementation Patterns

### • 1. Customer Support Chatbot

- RAG (Retrieval-Augmented Generation) over knowledge base
- Vector database (Pinecone) for semantic search
- Context window management for conversation history
- Escalation triggers for human handoff

### • 2. Document Processing

- Multimodal models for PDF/image processing
- Structured extraction with JSON schema validation
- Cross-document entity resolution
- Confidence scoring for quality control

### • 3. Security & Compliance

- PII filtering before LLM processing
- Audit logging for all interactions
- Access control and API key management
- GDPR compliance with data retention policies

## Cost Optimization

- **Caching:** 65% hit rate, \$12K/month savings
- **Prompt Compression:** Reduce context by 40%
- **Model Routing:** GPT-3.5 for simple, Claude for complex
- **Batch Processing:** Group similar requests