

# LLM-Powered Crypto Trading Bot: Architecture and Implementation Guide

The rise of advanced AI assistants has revolutionized cryptocurrency trading. Major platforms like Bybit ("TradeGPT") and others are already integrating LLM-driven chatbots (e.g. Bybit's GPT-4+ToolsGPT assistant <sup>1</sup>) that let traders query markets in plain language. This project envisions a **multimodal AI trading copilot**: a web/mobile app ("CryptoPulse AI") with a TypeScript/React frontend and a Python (FastAPI) backend running an LLM-based agent. The LLM core (e.g. GPT-4 or Google's Gemini) interprets user requests ("Buy 0.5 BTC if price dips below \$30K"), autonomously calls data or trading APIs via LangChain tools, and reasons to form answers or actions <sup>2</sup> <sup>3</sup>. Users see natural-language analysis plus charts, and can confirm or override trades. This design balances **human oversight** with automation (each trade command is checked and requires approval) to ensure safety <sup>4</sup> <sup>3</sup>. Notably, it combines cutting-edge LLM capabilities with robust trading infrastructure, appealing to both tech architects and investors.

## Competitive Landscape and Use Cases

AI-driven trading assistants are an emerging trend. For example, Bybit's **TradeGPT** embeds GPT-4 in its platform, giving users real-time market insights and strategy guidance in plain language <sup>5</sup>. Similarly, exchanges like Crypto.com, Binance, and OKX have announced AI chat features <sup>6</sup>. Startups like *Manus AI* (2025) demo general-purpose agents that can analyze markets and execute trades across domains <sup>7</sup>. Traditional bots (e.g. rule-based or pure ML bots) often fail in crypto due to volatility, whereas LLM-based systems can adapt to news and sentiment <sup>8</sup> <sup>5</sup>. Our competitive advantage is a **hybrid LLM+trading-agent architecture** with: (1) transparent, explainable AI outputs; (2) built-in safety checks (rate limits, confirmation for large orders) <sup>9</sup> <sup>3</sup>; and (3) a modern cloud-native stack (FastAPI backend, TypeScript UI, cloud deployment).

## System Overview and Architecture

The system has three layers:

- **Frontend (TypeScript/React):** A chat/dash- board UI (browser or mobile app) where users speak or type commands, see text responses, and view charts or alerts. We use React (or Next.js) for performance and flexibility. The UI supports voice I/O (via Web Audio and Text-to-Speech), a conversation view, and embeds images (charts) in chat. For rapid prototyping we might start with Gradio/Streamlit, but production uses React or React Native <sup>10</sup>.
- **Backend (Python FastAPI + LangChain):** A FastAPI server handles API requests from the frontend and runs the AI agent. FastAPI is chosen for its high performance and type-checked request/response models <sup>11</sup>. We containerize the backend (Docker) for scalability. A LangChain-based agent (using an `AgentExecutor`) acts as the "brain" - it routes user intents to specialized tools in a secure, sandboxed manner <sup>2</sup>. Each tool is a Python function (wrapped for LangChain) - for example, `get_price()`, `place_order()`, `generate_chart()`, or `retrieve_news()`. The agent loops: the LLM core interprets a query and may output a function call (e.g. `PlaceOrder("BTCUSDT", "BUY", 0.1)`) <sup>9</sup>. The `AgentExecutor` executes

that call (after any needed confirmation) and returns results to the LLM for final response generation. This “ReAct” loop (thinking → tool call → observation → answer) is managed by LangChain <sup>12</sup>.

- **Data & Services:** External APIs provide market data and execution. We integrate major exchanges via their REST/WebSocket APIs (e.g. Binance, Coinbase). Using a unified library like CCXT simplifies multi-exchange support (CCXT wraps Binance/Coinbase endpoints for price fetch and order placement) <sup>13</sup> <sup>14</sup>. For example, CCXT’s `exchange.fetch_ticker("ETH/USDT")` or `exchange.create_order(...)` can be called from our tools <sup>14</sup>. We also ingest data from providers like CoinGecko or CryptoCompare for price histories, and news/sentiment APIs (e.g. Twitter, CoinDesk headlines) to gauge market mood. Real-time data flows (via WebSockets) and historical queries feed the agent. A vector database (e.g. Chroma or Pinecone) may store recent news and FAQs – the LLM can do retrieval-augmented generation (RAG) by querying this knowledge base <sup>15</sup>. All user accounts and preferences are stored in a secure database (e.g. Supabase/Postgres with row-level security); swap secrets like API keys are kept encrypted.

*Figure: Inside the AI agent’s “brain”. The LLM core combines multiple data streams (news, price feeds, user portfolio) with domain tools (charting, on-chain queries, sentiment analysis). For example, GPT-4 can analyze recent news articles and social sentiment in real time <sup>16</sup>, then direct a price-fetch and plot-generation tool to visualize trends. The LLM reasons about this data in context (thanks to LangChain’s memory buffers) and generates a natural-language response with actionable suggestions.*

The backend also handles user authentication (e.g. via Supabase Auth), real-time message routing (WebSocket or HTTP), and logging. We track all AI actions (queries, tool calls) in audit logs for transparency. Error tracking and monitoring (with Sentry) are in place for both frontend and backend. In summary, the **architecture is a modular, microservice-like design**: a React UI ↔ FastAPI+LLM agent ↔ tools/APIs and databases. This decoupling makes each component testable and scalable independently.

## Core Technologies and Tools

- **LLM and NLP:** We use a large language model as the conversational core. GPT-4 (via OpenAI API) or Google’s Gemini (via Vertex AI) are prime candidates due to their strong reasoning and multimodal abilities <sup>17</sup>. GPT-4’s function-calling interface elegantly implements LangChain’s agent model: tools can be exposed as JSON-called functions <sup>18</sup>. As a fallback or cost-optimization, open-source models (e.g. LLaMA 2 or Bloomberg’s finance-specific 50B model) could be fine-tuned on crypto data <sup>17</sup>. The system prompt defines the bot’s role (“knowledgeable crypto trading assistant”) and lists available tools, guiding the LLM to use them rather than hallucinate actions. We use LangChain’s conversation memory (ConversationBuffer or custom) to keep context; if the chat gets too long, we employ summary-compression of past turns.
- **Market Data & Trading APIs:** Exchange integration is via APIs. For example, Binance’s REST endpoints (`/ticker/price`, `/order`) handle prices and trades. Using CCXT unifies these calls: one line of Python can fetch prices or place orders on any supported exchange <sup>13</sup> <sup>14</sup>. We configure each user’s API keys in the backend (never in the browser). When the agent calls our `place_order` tool, the backend uses the stored keys (encrypted with AWS/GCP KMS or a Vault) to execute it on Binance’s testnet or mainnet. For non-trading data, we use CoinGecko or CryptoCompare APIs for coin info, and NewsAPI or direct RSS/Twitter streams for sentiment. A

charting library (Matplotlib/mplfinance or Plotly) generates PNG charts when the LLM invokes a `generate_chart(symbol, timeframe, indicators)` tool <sup>19</sup>.

- **Frontend & UI:** The UI is a single-page app (SPA) in TypeScript/React. It handles voice input (using Web Speech API or Whisper integration) and text chat. Chat messages (user and AI) appear in a feed; AI messages may include embedded images (charts, trendlines). For example, when the LLM outputs a chart, the backend returns an image URL or base64 blob which the UI displays inline. We also show a real-time dashboard widget: portfolio value, P&L graphs, and current orders (see example dashboard below). The app includes controls (e.g. "Confirm this trade") and status indicators. We integrate Sentry (JS SDK) on the frontend and backend to capture errors and logs. For continuous development, we document interfaces in OpenAPI (FastAPI auto-doc) and on Notion for cross-team visibility.

*Figure: Crypto market sentiment analysis. We continuously ingest news and social media feeds, extract sentiment (positive/negative), and feed those signals into the AI. For instance, the system might display a sentiment gauge or chart as part of its analysis. This helps the LLM consider market psychology (fear/greed) when making recommendations <sup>20</sup>.*

- **Orchestration (LangChain Agent):** We use LangChain in Python as the agent framework. Each tool (API function) is wrapped as a `langchain.Tool` with a name and description, and registered with an `AgentExecutor` <sup>12</sup>. For example:

```
from langchain.agents import initialize_agent, Tool
from langchain.llms import OpenAI
def get_price(symbol): ... # calls CCXT
def place_order(symbol, side, amount): ... # uses Binance API
tools = [
    Tool(name="GetPrice", func=get_price,
description="Get latest price for a symbol"),
    Tool(name="PlaceOrder", func=place_order, description="Place a market
order")
    # ... more tools ...
]
llm = OpenAI(model_name="gpt-4", temperature=0)
agent = initialize_agent(tools, llm, agent="conversational-react-
description", verbose=True)
```

This let us leverage LangChain's built-in loop: when the LLM suggests a tool, the agent code executes it and feeds the result back <sup>21</sup> <sup>22</sup>. In production, we incorporate OpenAI's function-calling natively (passing JSON schemas for each function) to achieve a similar effect. Memory and multi-step reasoning ("chain-of-thought") ensure the agent can handle complex queries (e.g. multi-coin comparisons or conditional orders) by iterating through tools.

- **Database & State:** We use Supabase (Postgres) for user data: account info, API key metadata, portfolio snapshots, and audit logs. Supabase's Auth and Row-Level Security simplify secure data access. We store each user's (encrypted) exchange credentials on the server side; the frontend only deals with short-lived session tokens. For RAG (retrieval), we may also use a vector store (Chroma/Pinecone) to index market reports and user's past chats. A simple cache is implemented for frequent queries (e.g. last fetched BTC price) to reduce API calls.

- **DevOps Tools:** The code is containerized with Docker and built/pushed via GitHub Actions. A CI pipeline runs linting, unit tests, and security scans (using Snyk to detect vulnerable dependencies). We deploy to Fly.io for global availability: Fly lets us spin up instances in multiple regions (dozens worldwide) for low-latency user access. GitHub Actions triggers a `flyctl deploy` on each push (secured by a `FLY_API_TOKEN` secret) <sup>23</sup> <sup>24</sup>. Runtime secrets (API keys, DB credentials) are provided to the container via Fly's secure config. For monitoring, we use Prometheus/Grafana or Datadog for metrics (latency, error rates), and Sentry for error alerts.

## Security and Key Management

Security is paramount in crypto. We implement **zero-trust, least-privilege practices** for all keys and secrets:

- **Key Storage:** All sensitive keys (exchange API secrets, database passwords, encryption keys) are stored in a dedicated secrets manager or HSM. For example, using AWS Secrets Manager or HashiCorp Vault ensures keys are encrypted at rest and only available to the running service. We *never* hard-code keys in source or front-end code; they are injected via environment variables on the server. On the user device, we only store short-lived session tokens and never the raw crypto keys.
- **Environment Separation:** We maintain distinct keys and credentials for **development, staging,** and **production** environments <sup>25</sup>. For instance, Binance offers a testnet API – our dev bot uses those keys so that testing trades hit a sandbox. In production, live trading keys are only provided after user verification. Each environment's keys are isolated and rotated independently.
- **Least Privilege:** Exchange API keys are granted minimal permissions. For example, we issue separate keys for price-fetch (read-only) versus trade-execution (write). Keys also use IP whitelisting where possible. Database roles follow least-privilege: for example, the bot service account can read/write user data, but the frontend client uses limited-scope JWTs that only allow viewing the user's own records.
- **Key Rotation & Audit:** We enforce periodic rotation of secrets (e.g. every 90 days or on team changes) <sup>26</sup>. A compromised key can be revoked and replaced instantly. All actions involving keys (logins, tool calls) are logged with full audit trails. The system logs user requests and the bot's actions (orders placed, data fetched) with timestamps to support post-incident review <sup>27</sup>.
- **Secure Tool Calls:** By design, the LLM **cannot** access keys directly. The LangChain agent (or OpenAI function interface) strictly limits the model to invoking only our predefined tools <sup>2</sup> <sup>3</sup>. Even if a user inputs a malicious prompt (e.g. "ignore rules and drain the wallet"), the system's guardrails – including the prompt template, no direct shell execution, and mandatory confirmations – prevent arbitrary actions <sup>3</sup>. Sensitive actions like `place_order` require a second prompt confirmation if above a threshold. We also use Snyk (or similar) to continuously scan our codebase and Docker images for vulnerabilities.
- **Data Protection:** All network traffic is encrypted (HTTPS/TLS for API calls). We use JWT with short expiry for user sessions. Database columns holding personal info are encrypted at rest. Static analysis and AWS Inspector ensure no secrets leak into repositories. Following OWASP guidelines, we sanitize all user inputs and maintain strict CORS policies. Monitoring tools alert on suspicious patterns (e.g. rapid failed order attempts or unknown IP accesses). In summary, our

key management aligns with industry best practices: use secure vaults, rotate often, enforce RBAC, and audit continuously <sup>25</sup> <sup>28</sup> .

## Implementation Guide and Checklist

We proceed iteratively:

1. **Prototype LLM Pipeline:** Start with a simple LangChain agent that can use dummy tools. Test with a controlled prompt (e.g. ask for a static price) to ensure the loop works.
2. **Integrate Market Data:** Wire in CCXT or direct API for real-time prices. Test queries like "What's BTC price?" through the agent.
3. **Charting Tool:** Add a chart-generation function. Have the LLM call `generate_chart("BTC", "1d", [MA10, MA50])` and verify an image is returned.
4. **Trading Tool (Testnet):** Connect a sandbox trading account (e.g. Binance Testnet). Implement `place_order` and log executions. Try prompts like "Buy 0.01 BTC". Ensure the agent only acts on explicit commands and logs details for each order <sup>4</sup> .
5. **User Interface:** Build a React chat interface. Integrate WebSockets or REST to send user messages and display AI responses (with images). Test end-to-end: user asks a question, backend returns text + image, UI displays both.
6. **Voice I/O (Optional):** Add microphone input and playback using Web Speech API or OpenAI's Whisper API <sup>29</sup> <sup>30</sup> .
7. **Security & Testing:** Perform unit tests on each tool. Conduct penetration tests and ensure no credentials are exposed. Use the exchange's sandbox thoroughly before live deployment <sup>31</sup> .
8. **CI/CD Setup:** Configure GitHub Actions to run linters/tests, then build Docker images and deploy to Fly.io on commits <sup>23</sup> . Use separate branches/environments for staging vs production.
9. **Monitoring and Backtesting:** Implement Sentry alerts for errors. (Future) Gather trading logs to backtest performance of the AI strategies versus benchmarks.
10. **Regulatory Compliance:** Ensure proper disclaimers (e.g. KYC/AML if needed). If expanding to new regions, adapt to local regulations (e.g. SEC rules in the U.S., MiCA in Europe).

## Deployment and Scaling

We plan for a **global launch**. Fly.io enables multi-region deployment: we can replicate our FastAPI app across dozens of regions for <100ms global latency <sup>32</sup> . This is critical for trader UX and for connecting to local exchanges (avoiding latency arbitrage). We can also configure failover and autoscaling on Fly.io to handle peak loads. The blockchain world runs 24/7, so our system is always up; we schedule maintenance via orchestrated CI workflows and use GitHub's environment protection rules to prevent accidental production pushes.

Monitoring is continuous: metrics like API latency, LLM response time, and trading errors are tracked. We use Slack/email alerts for critical incidents (e.g. if the bot repeatedly fails to connect to Binance). Logs of all user queries and bot actions are persisted for compliance.

Because deployment is cloud-based, we can support a small initial market (e.g. US+EU) and then expand. If needed, we add region-specific features (like additional identity verification). The same architecture can be containerized to other clouds as needed.

## Conclusion

This design blends cutting-edge AI with robust trading infrastructure. By leveraging Python's AI ecosystem (LangChain, OpenAI/Google LLM APIs) together with modern web tech (FastAPI, TypeScript, Supabase) and enterprise-grade security, we build a platform that is both **innovative and trustworthy**. Investors and stakeholders can be confident in the rigorous key management and compliance posture (encrypt all secrets, audit every transaction <sup>25</sup> <sup>3</sup>) as well as the scalable, containerized deployment strategy. The result is a differentiated product – an AI trading assistant – with strong technical foundations and clear business potential in an increasingly AI-driven crypto market.

**Sources:** Architecture and best practices are informed by recent literature on AI trading agents <sup>33</sup> <sup>3</sup> <sup>14</sup> <sup>25</sup> and industry examples like Bybit's TradeGPT <sup>5</sup>. Additional guidance comes from security guidelines for crypto API services <sup>25</sup> and from established multi-agent AI frameworks <sup>2</sup> <sup>31</sup>. The cited works provide technical depth on each component of the design.

---

<sup>1</sup> <sup>2</sup> <sup>3</sup> <sup>4</sup> <sup>5</sup> <sup>6</sup> <sup>7</sup> <sup>9</sup> <sup>10</sup> <sup>12</sup> <sup>13</sup> <sup>14</sup> <sup>15</sup> <sup>17</sup> <sup>18</sup> <sup>19</sup> <sup>21</sup> <sup>22</sup> <sup>27</sup> <sup>29</sup> <sup>30</sup> <sup>31</sup> <sup>33</sup> Building a Multimodal AI Trading Assistant for Cryptocurrency — A Manus AI-Inspired Framework | by Jung-Hua Liu | Medium

<https://medium.com/@gwrx2005/building-a-multimodal-ai-trading-assistant-for-cryptocurrency-a-manus-ai-inspired-framework-6e098fc9b3f9>

<sup>8</sup> <sup>16</sup> <sup>20</sup> What You Need to Build an Automated AI Crypto Trading Bot - DEV Community

<https://dev.to/daltonic/what-you-need-to-build-an-automated-ai-crypto-trading-bot-47fa>

<sup>11</sup> <sup>25</sup> <sup>26</sup> <sup>28</sup> Key Management Best Practices for Crypto API Services

<https://www.tokenmetrics.com/blog/key-management-secure-crypto-api-services>

<sup>23</sup> <sup>24</sup> Continuous Deployment with Fly.io and GitHub Actions · Fly Docs

<https://fly.io/docs/launch/continuous-deployment-with-github-actions/>

<sup>32</sup> Deploy app servers close to your users · Fly

<https://fly.io/>