# 📊 Web3 Data Pipeline for Model Ingestion

## 1. Project Overview

This project defines the architecture and implementation specifications for a robust, scalable data pipeline designed to extract, process, and transform raw **Web3 data** (primarily from the blockchain) into structured **feature sets** ready for **Machine Learning (ML)** model training and real-time inference.

The pipeline must handle the unique characteristics of blockchain data—namely, its immutability, high volume, and complex nested structure (logs and traces).

**Goal:** To establish a reliable, low-latency ETL (Extract, Transform, Load) workflow capable of feeding a continuous stream of up-to-date, high-quality Web3 features into an MLOps environment.

## 2. Architecture Overview

The system follows a modern **Medallion Architecture (Bronze, Silver, Gold)** pattern adapted for blockchain data, ensuring data quality increases at each stage.

### Core Architectural Layers

| Layer | Description | Data State | Destination |
|---|---|---|---|
| **Bronze (Raw Ingestion)** | Raw, schema-on-read data directly from the source (e.g., full blocks, raw transactions, un-decoded logs). | Immutable, raw, uncleaned. | Data Lake (S3, ADLS, GCS) |
| **Silver (Cleaned & Normalized)** | Data that has been cleaned, filtered, and decoded using smart contract **ABIs**. Normalized across different chains/protocols. | Structured, validated, normalized. | Data Warehouse/Lake Tables |
| **Gold (Feature Store Ready)** | Aggregated, time-series, and behavioral features ready for direct consumption by ML models. | Highly aggregated, time-series features. | Feature Store |

## 3. Pipeline Stages & Components Specification

### Stage 1: Extraction & Ingestion (E)

The primary challenge is moving data off-chain efficiently without resorting to slow, rate-limited RPC node polling.

| Component | Role | Data Target | Technology/Tooling |
|---|---|---|---|
| **Real-Time Streamers** | Captures **Smart Contract Event Logs** and pending transactions instantly. Crucial for real-time anomaly detection. | Events (`Transfer`, `Swap`, `Mint`), Pending Txns. | QuickNode Streams, Kafka/PubSub, custom WebSocket listeners |

| | | | |
|---|---|---|---|
| **Historical Batch ETL** | Extracts full historical block data, transaction details, and internal traces for long-term model training. | Blocks, Transactions, Traces. | Dedicated Blockchain ETL Provider, custom Python scripts (using high-performance RPC). |
| **Off-Chain APIs** | Gathers supplementary data for enrichment. | Price Data (OHLCV), Exchange Rates, Social/Sentiment Data. | CoinGecko API, proprietary exchange APIs, social media APIs. |

## Stage 2: Transformation & Feature Engineering (T)

This stage cleans the data and generates meaningful features.

### A. Decoding and Normalization (Silver Layer)

1. **ABI Integration:** Implement a lookup service that uses a contract's **ABI (Application Binary Interface)** to correctly decode the raw hexadecimal input data and event logs into structured fields (e.g., `uint256` to actual token amount).

2. **Data Cleaning:** Filter out failed transactions, remove duplicates, and handle potential malicious data patterns.

3. **Schema Enforcement:** Apply a strict schema to the parsed data to ensure consistency.

### B. Feature Engineering (Gold Layer)

Aggregate normalized Silver data into ML-ready features.

| Feature Type | Example Features | ML Use Case |
|---|---|---|
| **Liquidity & Finance** | Total Value Locked (TVL), 24h Trading Volume, Average Gas Price per hour, Liquidity Pool Depth. | Price Prediction, DEX Arbitrage modeling. |
| **Behavioral** | User transaction frequency (7-day window), token holding duration, "whale" vs. "retail" clustering. | User Segmentation, Market Manipulation detection. |
| **Temporal** | Features aggregated by time window (e.g., 5-minute, 1-hour). Calculate moving averages and volatility metrics. | Time Series Forecasting. |

## Stage 3: Loading & Consumption (L)

The final stage prepares data for immediate use by the ML platform.

| Component | Role | Requirements | Technology/Tooling |
|---|---|---|---|
| **Feature Store** | Centralized repository for feature serving. Must support low-latency reads for real-time inference. | Consistency between training and serving datasets. Point-in-time correctness. | Feast, dedicated Databricks Feature Store. |

| | | | |
|---|---|---|---|
| **MLOps Orchestration** | Manages the entire ML lifecycle: training, deployment, monitoring, and pipeline scheduling. | Scheduling pipeline runs (Airflow/Dagster), model versioning (MLflow), automated retraining triggers. | Airflow/Dagster, MLflow. |
| **Model Training Environment** | Trains models using the Feature Store data. | Scalable compute (Spark), GPU support (if necessary). | Databricks, AWS Sagemaker, Google Vertex AI. |

## 4. Key Technical Challenges & Solutions

| Challenge | Impact | Proposed Solution |
|---|---|---|
| **Data Volume & Throughput** | Blockchains generate massive amounts of data per day (millions of transactions/events). | Utilize **distributed streaming platforms** (Kafka, Spark Streaming) for processing and partitioned storage solutions in the Data Lake. |
| **Smart Contract Decoding** | Raw logs are unintelligible without the contract's ABI. | Build an **ABI Registry Service** that automatically fetches and manages ABIs. Use ABIs in the Silver layer transformation to decode logs dynamically. |
| **Data Consistency** | Ensuring all related data elements (e.g., block, transaction, receipt) arrive and are processed together. | Implement idempotent processing and robust schema validation at the Bronze/Silver boundary. Use unique identifiers (transaction hash, block number) for joining data. |
| **Real-Time Latency** | Models requiring near-instantaneous data (e.g., liquidations, front-running detection). | Prioritize event-driven architecture using dedicated streaming tools (QuickNode Streams) and ensure the Gold layer provides low-latency feature serving. |

## 5. Deliverables

1. Complete, executable data pipeline source code (Python/Scala).

2. Defined data schemas for the Bronze, Silver, and Gold layers.