# Windows Event Log Analysis: Implementation Details
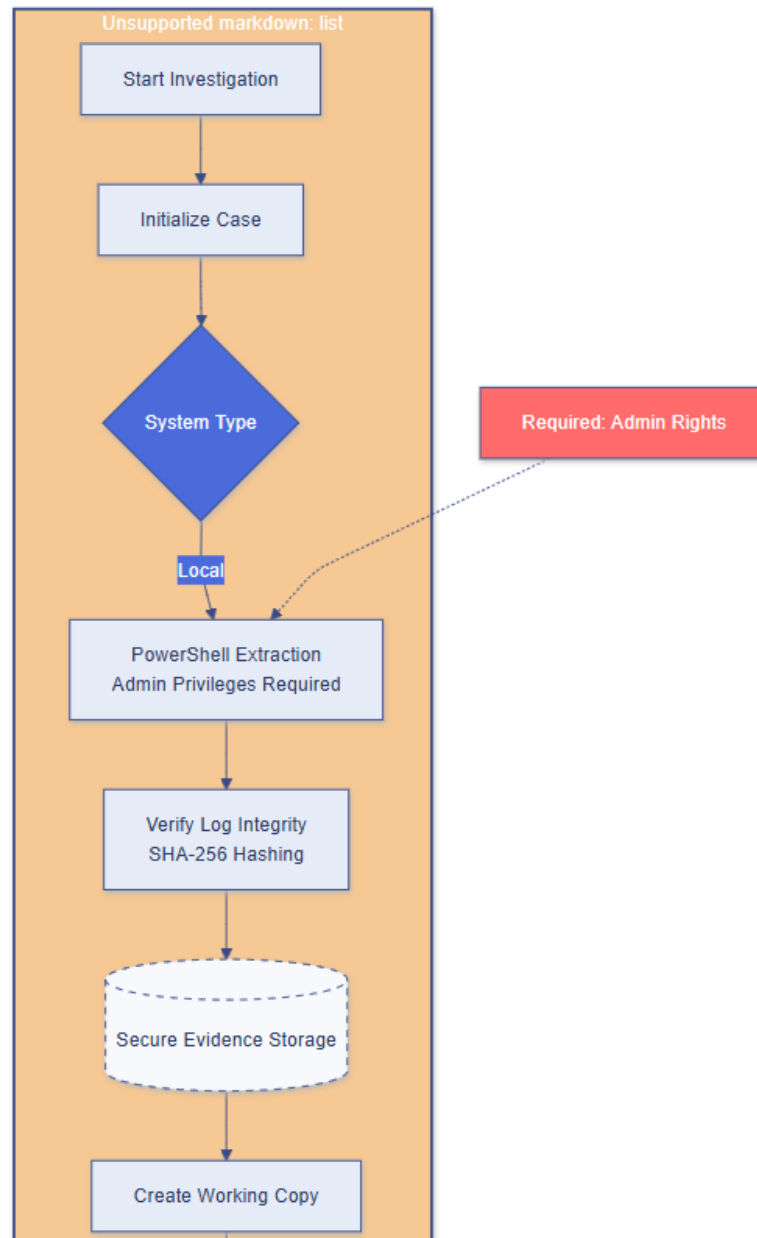
A Comprehensive Guide to Each Step

## Step 1: Event Log Collection and Initial Handling

Unsupported markdown: list

Start Investigation

Initialize Case

System Type

Required: Admin Rights

Local

PowerShell Extraction
Admin Privileges Required

Verify Log Integrity
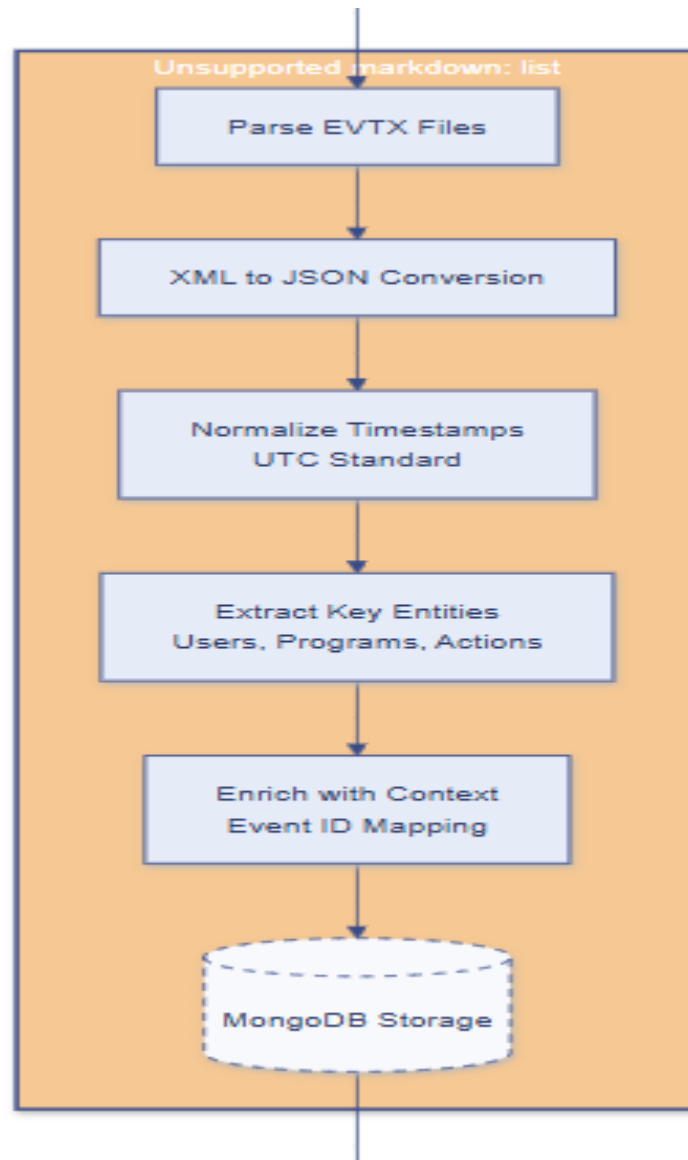SHA-256 Hashing

Secure Evidence Storage

Create Working Copy

Starting with the collection of event logs files .evtx from Windows, the system will first gather these logs. After collection, the system will create backup copies of these files, which will be maintained for the current session only.

For handling the .evtx files, we will develop a custom Python wrapper around the python-evtx library. This wrapper will help manage large log files efficiently by implementing streaming processing. Rather than loading entire logs into memory, the system will process them in chunks of 1000 events. Each chunk will go through initial validation to ensure data integrity and proper formatting. If the system encounters corrupted sections, it will mark these for special handling and continue processing the readable portions.

The secure storage system will implement a combination of encrypted file storage and a metadata database. This database, built on PostgreSQL, will maintain critical information about each log file, including source system details and collection timestamps. The system will also implement automated backup procedures that will create redundant copies of all received files.

## Step 2: Data Transformation and Structuring

Unsupported markdown: list

```
Parse EVTX Files
        ↓
XML to JSON Conversion
        ↓
Normalize Timestamps
UTC Standard
        ↓
Extract Key Entities
Users, Programs, Actions
        ↓
Enrich with Context
Event ID Mapping
        ↓
MongoDB Storage
```
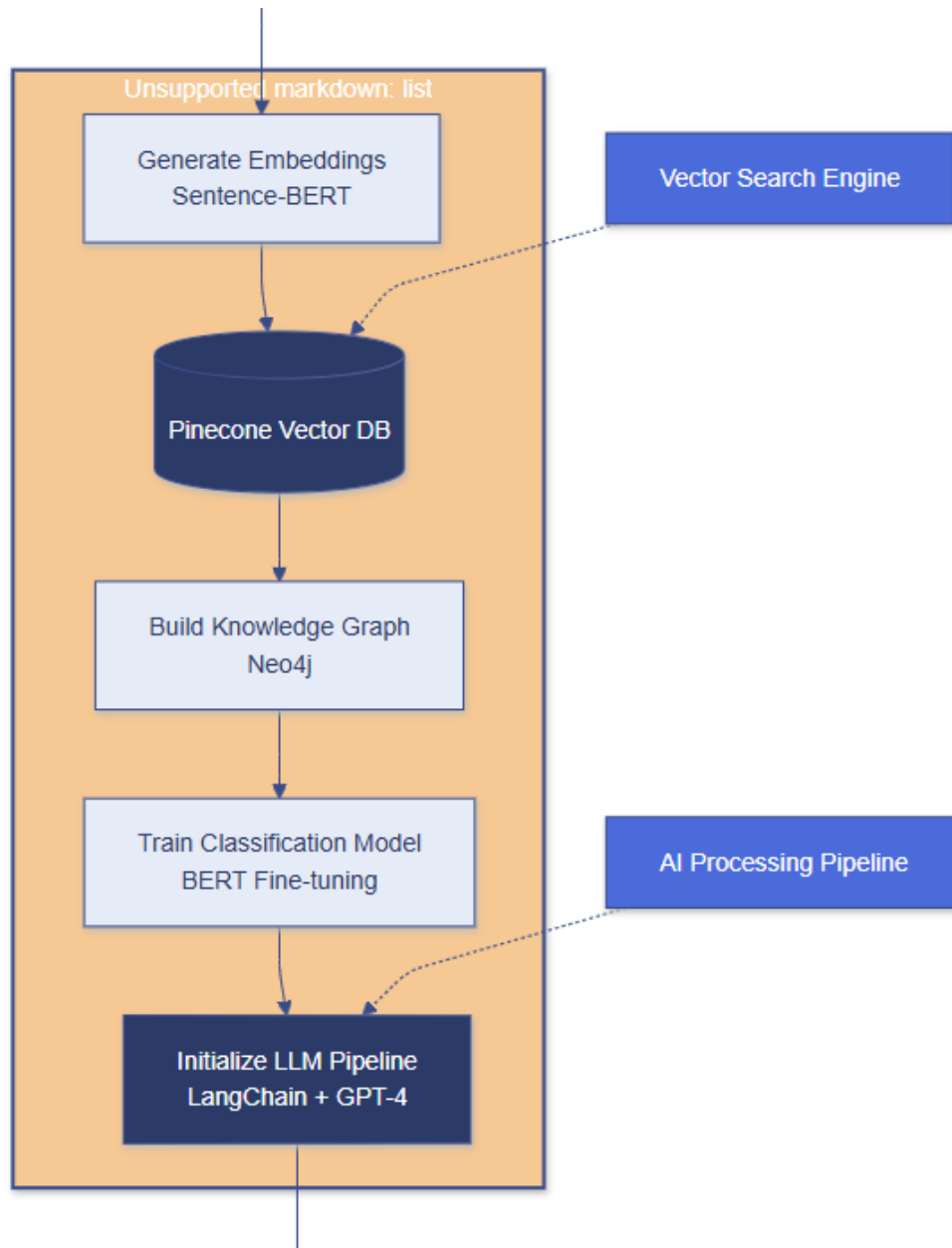
The transformation phase will convert raw Windows Event Logs into a structured format that our AI models can efficiently process. The system will utilize an ETL (Extract, Transform, Load) pipeline that will first extract the XML structure of each event and convert it into a normalized JSON format. During this conversion, special attention will be paid to timestamp normalization, ensuring all temporal data will be converted to UTC to maintain consistency across different time zones.

The pipeline will include specialized handlers for different event types. For instance, security events (Event ID 4600-4900 series) will receive extra processing to extract user accounts, access types, and affected resources. System events will undergo additional parsing to identify software installations, service changes, and system modifications. The system will implement parallel processing using Python's multiprocessing module to

handle multiple log files simultaneously, significantly reducing processing time for large datasets.

The transformed data will then be enriched with additional context. The system will maintain a mapping database of common Event IDs to human-readable descriptions, which will be used to annotate each event. Entity extraction will also be implemented at this stage, identifying and tagging key elements like usernames, program names, and system components. This enriched data will be stored in the primary analysis database, which will be built using MongoDB to handle the semi-structured nature of event log data efficiently.
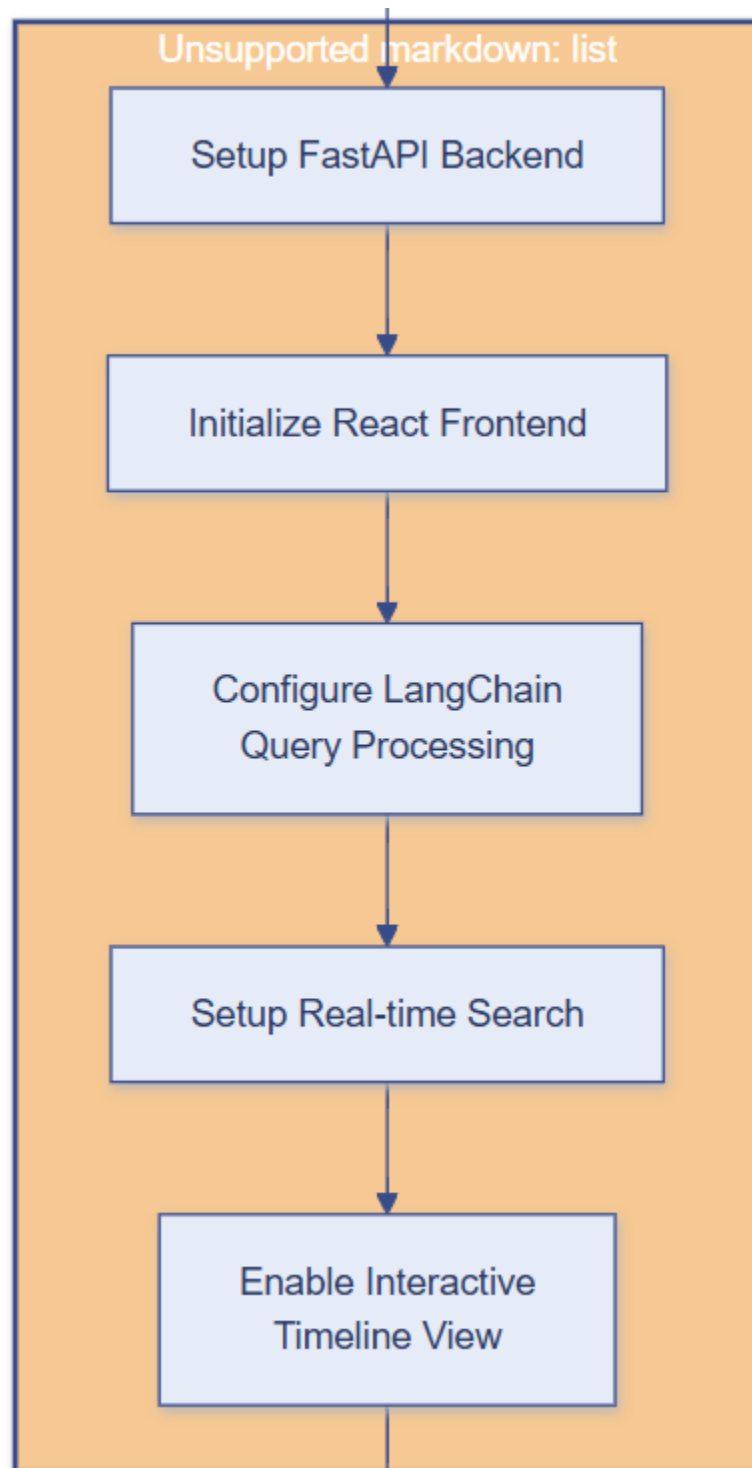
## Step 3: AI Model Application



The AI processing layer will combine several sophisticated models working in concert. At the core, the system will use a fine-tuned BERT model specifically trained on Windows Event Log data to classify events by their significance and type. This model will be trained on a large dataset of labeled event logs, allowing it to identify patterns that might indicate suspicious activity.

The AI pipeline will be implemented using LangChain as the orchestration framework. This will allow the system to combine multiple LLMs effectively. The system will utilize GPT-4 through the OpenAI API for natural language understanding and Claude for reasoning about event sequences. The pipeline will first create embeddings of event descriptions

using Sentence-BERT, storing these in a Pinecone vector database for efficient similarity searching.

The knowledge graph construction will utilize Neo4j to maintain relationships between events. The system will implement custom algorithms that will identify causal relationships between events, tracking system state changes and user activities over time. This graph structure will be crucial for answering complex queries about system behavior and user actions

## Step 4: Query Interface Development

Unsupported markdown: list

```
Setup FastAPI Backend
        |
        v
Initialize React Frontend
        |
        v
Configure LangChain
Query Processing
        |
        v
Setup Real-time Search
        |
        v
Enable Interactive
Timeline View
```
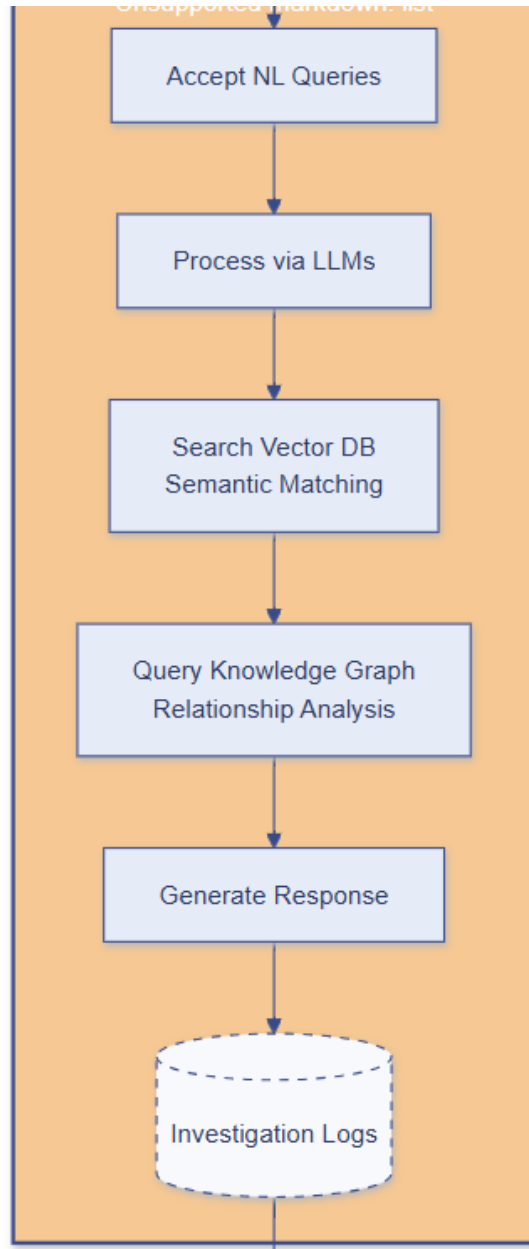
The query interface will be built as a modern web application using FastAPI for the backend and React for the frontend. The interface will implement a chat-like interaction model where investigators can ask natural language questions about the system's history. The system will include a custom query processing pipeline that will convert natural

language questions into structured queries against the event database and knowledge graph.

The interface will include real-time suggestion features that will help investigators formulate effective queries. As they type, the system will analyze the partial query and suggest relevant event types, time ranges, or entities that might be of interest. A context-aware response system will be implemented that will maintain conversation history and can reference previous queries and findings.

Behind the scenes, the system will use a combination of LLM-powered processing and traditional search algorithms. When a query comes in, it will first go through an intent classification model to determine the type of information being requested. Based on the intent, the system will either perform a direct database query, a graph traversal operation, or engage the LLM pipeline for more complex analysis.

## Step 5: Investigation Phase Management

Accept NL Queries

Process via LLMs

Search Vector DB
Semantic Matching

Query Knowledge Graph
Relationship Analysis

Generate Response

Investigation Logs

During the investigation phase, the system will maintain an active session that will track all queries and findings. The system will implement a structured investigation methodology where each query and its results will be automatically documented. LlamaIndex will be used to maintain an efficient index of all discovered information, allowing for quick reference and cross-referencing during the investigation.

The investigation management system will include automatic anomaly detection running in the background. As investigators interact with the system, the models will continuously analyze the data for patterns that might indicate suspicious activity. These potential

findings will be queued and presented to the investigator at appropriate moments, helping guide the investigation without interrupting the current line of inquiry

Reporting (Optional): Provide a session report.