



Designing an AI Solution for SOC Log Analysis Challenge

Building an AI system to make sense of SOC logs is a **multi-layered architectural problem**. Here's a step-by-step approach to designing a competitive solution:

Phase 1: Understanding the Problem Space

Before building, you need to understand what you're solving:

The Core Challenge: SOC teams receive thousands of log events per second. Most are noise (normal activity), but hidden within are critical threats. Your AI must:

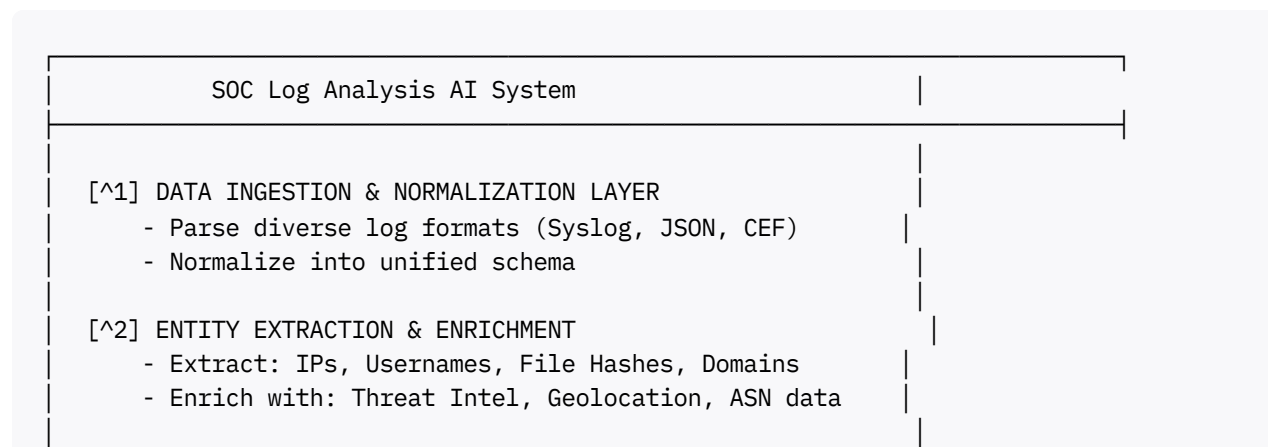
- Separate signal (real threats) from noise (normal operations)
- Explain *why* something is a threat in plain English
- Correlate logs across multiple sources (firewall, authentication, endpoint, database)
- Do this at scale without requiring manual rule creation

Key Constraints:

- False positives burn out analysts (alert fatigue)
- Missing attacks (false negatives) cause breaches
- Logs are messy: different formats, timestamps, field names
- Real-time processing is critical (not batch analysis)

Phase 2: High-Level Architecture Design

A robust solution typically has **5 main components**. This architecture balances accuracy, interpretability, and scalability:



```
[^3] ANOMALY DETECTION & CORRELATION ENGINE
- Detect deviations from baseline behavior
- Correlate events across logs
- Pattern matching (known attack signatures)
```

```
[^4] LLM-BASED INTERPRETATION & REASONING
- Summarize findings in natural language
- Generate incident narratives
- Provide evidence chains
```

```
[^5] ALERTING & ACTION LAYER
- Rank by risk/confidence
- Generate recommended actions
- Output to SIEM/Ticketing system
```

Phase 3: Detailed Component Design

Component 1: Data Ingestion & Normalization

Goal: Convert chaos into structure.

Approach:

- **Log Parsing:** Use regex, state machine parsers, or LLM-based parsing to extract fields
 - Input: Jan 15 10:23:45 firewall01 authpriv: Failed password for user admin from 192.168.1.100 port 22 ssh
 - Output: { timestamp: "2025-01-15T10:23:45Z", event_type: "auth_failure", source_ip: "192.168.1.100", user: "admin", service: "ssh" }
- **Schema Standardization:** Map all logs to a common schema (e.g., OASIS Common Event Format - CEF)
 - This allows downstream components to work uniformly

Technology choices:

- **Simple approach:** Regex + Pandas/PySpark
- **Moderate complexity:** Log parsing libraries (e.g., Logstash, Fluent Bit)
- **Advanced approach:** LLM-based parsing (GPT, Llama) for highly variable formats

Why this matters for your challenge: Judges will test with *real, messy* logs. If your normalization fails, everything downstream fails.

Component 2: Entity Extraction & Enrichment

Goal: Turn raw strings into intelligence.

Approach:

- **NER (Named Entity Recognition):** Identify and classify entities
 - Use a pre-trained NER model (spaCy, transformers) trained on security data
 - Extract: IP addresses, domain names, file hashes (MD5, SHA256), usernames
- **Threat Intelligence Enrichment:** Look up entities against threat databases

```
IP: 192.168.1.100 → Is it on a known malware C2 list? → Flag it
Domain: example-login.com → Phishing domain? → Flag it
File Hash: abc123def → Known ransomware? → Flag it
```

Technology choices:

- **Libraries:** spaCy, Hugging Face transformers
- **External APIs:** VirusTotal, AlienVault OTX, AbuseIPDB
- **Self-hosted:** Run local threat intel feeds (YARA rules, Suricata)

Critical success factor: This is where many solutions fail. If you flag a benign IP as malicious without evidence, you lose points on false positives.

Component 3: Anomaly Detection & Correlation

Goal: Find the deviations that matter.

Approach:

- **Baseline Behavior Modeling:** Learn what's "normal"
 - For each user: typical login times, typical login locations, typical data access patterns
 - For the network: normal traffic volume, normal IPs, normal protocols
 - **Technique:** Use unsupervised learning (Isolation Forest, Local Outlier Factor) to learn the baseline
- **Deviation Detection:** Flag events that break the baseline
 - User logs in from a new country at 3 AM → Anomalous
 - Sudden spike in failed login attempts → Anomalous
 - Connection to a known C2 server → Anomalous
- **Cross-Log Correlation:** Link events across sources
 - Firewall log shows connection to 1.2.3.4
 - - Identity log shows user X authenticated at same time

- - Endpoint log shows process Y connecting to same IP
- = Correlated incident (not just one isolated event)

Technology choices:

- **Classical ML:** Isolation Forest, One-Class SVM
- **Deep Learning:** LSTM-based sequence anomaly detection (good for temporal patterns)
- **Graph-based:** Build a graph of entities and relationships; find abnormal patterns

Common pitfall: Detecting anomalies without explaining them. Your solution must preserve the "evidence trail" for later interpretation.

Component 4: LLM-Based Interpretation & Reasoning

Goal: Explain *why* this is a threat. This is where GenAI/LLMs shine—and where teams often struggle.

Approach:

Option A: Prompt Engineering + Few-Shot Learning (Fastest to implement, less control)

Prompt Template:

"Analyze these security events and explain if they constitute a threat:

- Event 1: [log line 1]
- Event 2: [log line 2]
- Event 3: [log line 3]

Threat Intelligence Context:

- IP 1.2.3.4 is known for: [malware C2]

Provide: (1) Threat assessment, (2) Confidence score, (3) Recommended action, (4) Evidence

Issues with this:

- LLMs can "hallucinate" (make up facts)
- No control over reasoning
- Can be expensive at scale

Option B: AI Agent Architecture (More robust, better control) ^[1] ^[2]

Build an **autonomous AI agent** that:

1. Receives anomalous events
2. **Autonomously decides what to investigate** (not told what to do)
3. Gathers correlated evidence from multiple logs
4. **Chains reasoning:** "First, I noticed a login failure. Then a success. Then a privilege escalation. This sequence suggests a brute-force attack followed by lateral movement."

5. Provides a structured incident report with:

- Attack hypothesis
- Supporting evidence (with log line references)
- Risk score
- Recommended actions

Agent tools (things the agent can call autonomously):

- `query_logs(filter, time_range)` → Search logs
- `get_threat_intel(ip/domain)` → Look up threat data
- `correlate_events(events)` → Find relationships
- `generate_summary(events)` → Create narrative
- `calc_risk_score(evidence)` → Score severity

LLM orchestration:

- Use a small, efficient LLM (Llama-3-8B, Mistral) as the "brain"
- It decides which tools to call and interprets results
- Larger, more expensive models (GPT-4) only for complex reasoning (optional)

Technology stack:

- **Agent framework:** LangChain, Autogen, or custom orchestration
- **LLM:** Claude 3.5 Sonnet, GPT-4o, Llama-3-70B, Mistral 7B
- **Memory:** Vector database (Pinecone, Weaviate) to store log embeddings for fast retrieval

Why this wins competitions: Judges care about explainability. An agent that says "I found a brute-force attack because [specific log lines] show 50 failed logins in 2 minutes, followed by a successful login from the same IP" beats a black box that just says "Threat detected."

Component 5: Alerting & Prioritization

Goal: Don't overwhelm the analyst; surface the most critical threats.

Approach:

- **Risk Scoring:** Combine multiple signals into a single score

$$\text{Risk Score} = (\text{Confidence} \times \text{Weight}_1) + (\text{Impact} \times \text{Weight}_2) + (\text{Severity} \times \text{Weight}_3)$$

- Confidence: How sure are we this is a real threat? (0–1)
- Impact: How bad is it if true? (0–1)
 - Example: Privilege escalation to admin = HIGH impact
- Severity: How urgent? (0–1)
 - Example: Active lateral movement = HIGH urgency

- **Alert Thresholding:** Only alert if Risk Score > threshold
 - Prevents alert fatigue
- **Grouping & Deduplication:** Collapse similar alerts
 - Instead of "100 failed logins from same IP," create 1 alert: "Brute-force attack detected from 1.2.3.4"

Phase 4: Development Strategy (Step-by-Step)

Step 1: Start Simple (Week 1–2)

- Get data ingestion and normalization working
- Implement basic log parsing
- Test on small dataset
- **Goal:** Prove you can reliably extract fields from logs without losing information

Step 2: Add Anomaly Detection (Week 2–3)

- Implement a simple baseline model (e.g., Isolation Forest)
- Detect known attack patterns (brute force, data exfiltration)
- Measure false positives vs. recall
- **Goal:** Get baseline metrics working

Step 3: Integrate Enrichment & Intelligence (Week 3–4)

- Add entity extraction (NER)
- Integrate threat intel APIs
- Implement log correlation (linking related events)
- **Goal:** Improve detection accuracy without raising false positives

Step 4: Add LLM Interpretation (Week 4–5)

- Start with simple prompt engineering for summarization
- Gradually move to agent-based reasoning if time permits
- **Goal:** Make alerts understandable to analysts

Step 5: Polish & Optimize (Week 5–6)

- Reduce false positives aggressively
- Optimize latency
- Test edge cases
- **Goal:** Production-ready output

Phase 5: Technical Stack Recommendations

For parsing & data processing:

- Python + Pandas/Polars (data manipulation)
- spaCy or Hugging Face Transformers (NER)
- Pyspark (if handling massive log volumes)

For ML/anomaly detection:

- scikit-learn (Isolation Forest, One-Class SVM)
- PyTorch/TensorFlow (if using deep learning)
- OR use a managed service (e.g., AWS Lookout for Logs, Azure Sentinel)

For LLMs & agents:

- LangChain or LlamaIndex (agent orchestration)
- Ollama (run open-source LLMs locally, no API costs)
- OpenAI API (GPT-4o, if budget allows)
- Anthropic Claude (good balance of cost and performance)

For storage & retrieval:

- PostgreSQL + vector extension (pgvector) for embeddings
- Vector DB: Pinecone, Weaviate, or Qdrant (optional but helpful for retrieval)

For monitoring & visualization:

- Streamlit or Gradio (quick UI for judges to see results)
- Plotly (interactive dashboards)

Phase 6: How to Stand Out (Win Strategies)

1. Focus on False Positive Reduction First

- Any team can build something that flags everything as a threat
- Winners minimize false positives while maintaining high recall
- Show your FPR in the demo

2. Excellent Explainability

- Every alert must include: (1) What happened, (2) Why it's suspicious, (3) What to do next, (4) Evidence links
- Include a clickable link to the raw log lines
- Judges will test this ruthlessly

3. Handle Edge Cases

- What if a log field is missing?

- What if timestamps are wrong?
- What if an IP is spoofed?
- Show you thought about this

4. **Correlation Across Logs**

- Bonus points if you link network logs + identity logs + endpoint logs
- Show a "kill chain" narrative (e.g., "Attack progressed as: recon → exploitation → lateral movement")

5. **Real-Time Performance**

- If processing takes 10 minutes for 10,000 logs, that's too slow
- Aim for < 1 second per alert generated
- Show latency metrics in your demo

6. **Beautiful Presentation**

- Create a simple web dashboard showing:
 - Top threats (ranked by risk)
 - Attack narrative (with evidence)
 - Statistics (FPR, recall, processing time)
- Judges appreciate good UI

Phase 7: Testing & Validation Strategy

Golden Test Dataset:

1. **Clean logs (50%):** Normal activity. Your system should be quiet here.
2. **Attack logs (40%):** Known attacks injected. Measure recall.
3. **Edge cases (10%):** Benign-but-suspicious logs. Test reasoning.

Metrics to track:

- Precision (of all alerts, how many are real threats?)
- Recall (of all real threats, how many did we catch?)
- F1-score (balance of precision and recall)
- False positive rate
- Processing latency
- Cost (API calls, compute time)

Validation process:

- Split logs: 70% training, 30% held-out test
- Tune parameters on training set
- Final evaluation on test set (judges won't see this)

- Expect judges to introduce *new* attacks you haven't seen before

Quick Example: Building the First MVP

Here's a minimal viable product to get started:

```
# Pseudocode for a basic SOC AI system

import pandas as pd
from sklearn.ensemble import IsolationForest
import requests

def analyze_soc_logs(raw_logs):
    # Step 1: Parse & normalize
    normalized = parse_logs(raw_logs) # dict with fields: timestamp, user, ip, action, e

    # Step 2: Extract entities & enrich
    enriched = enrich_with_threat_intel(normalized)

    # Step 3: Detect anomalies
    baseline = train_baseline_model(training_logs)
    anomalies = detect_anomalies(enriched, baseline)

    # Step 4: Generate explanations
    for anomaly in anomalies:
        summary = generate_summary_with_llm(anomaly)
        risk_score = calculate_risk(anomaly)

        if risk_score > THRESHOLD:
            alert = {
                "risk": risk_score,
                "summary": summary,
                "evidence": anomaly["log_lines"],
                "action": get_recommended_action(anomaly)
            }
            yield alert
```

Key Takeaway

Winning solutions balance **three competing forces**:

1. **Accuracy**: Find the real threats (high recall)
2. **Specificity**: Avoid false alarms (high precision)
3. **Explainability**: Make analysts trust the output

Most teams optimize for (1), sacrificing (2) and (3). The winners optimize for all three. Start with #2 (reduce false positives), then layer on #1 and #3.

1. <https://arxiv.org/abs/2509.05936>
2. <https://radiantsecurity.ai/learn/ai-agents/>
3. <https://www.paloaltonetworks.in/cyberpedia/ai-in-threat-detection>
4. <https://radiantsecurity.ai/learn/evaluate-ai-soc-analysts/>
5. <https://www.semanticscholar.org/paper/LLMeLog:-An-Approach-for-Anomaly-Detection-based-on-He-Jia/8729c2ffaa24a855acd60b1fbe59f93c1553584f>
6. <https://corelight.com/resources/glossary/ai-threat-detection>
7. <https://www.ijirmps.org/papers/2025/1/232046.pdf>
8. https://www.youtube.com/watch?v=xdUR8-_P3DU