

# AI-Orchestrated Remediation with Human-in-the-Loop (HITL)

## 1. The Strategic "Why": Why AI is Mandatory

When managing **100+ runbooks**, a traditional system is "blind." It executes a script because a keyword matched, without understanding *why*.

We need AI to bridge the gap between "Raw Data" and "Human Understanding":

- **The Scaling Problem:** Manually mapping 100+ unique alerts to 100+ GitHub pipelines creates a "Maintenance Trap."
- **The Data Problem:** Splunk logs are technical and messy. AI converts them into **Plain English** so an on-call engineer can make a decision in seconds, not minutes.
- **The Safety Problem:** We cannot trust 100% autonomous scripts for critical infrastructure. AI provides the **Reasoning**, but the Human provides the **Permission**.

---

## 2. How the "Smarter" Workflow Works

Instead of a direct trigger, we insert an **Interactive Layer** (Slack, Microsoft Teams, or ServiceNow).

### Step 1: Intelligent Interpretation (Splunk → AI)

Splunk triggers an alert. The AI reads the log and searches the **Vector Database** of 100+ runbooks. It finds the best match based on the *meaning* of the error.

### Step 2: The Human-Readable Proposal (The "Smart" Notification)

The AI does not run the GitHub pipeline immediately. Instead, it sends a message to the engineer that looks like this:

**AI Support Bot: Incident Identified**

**What happened?** > I detected a "Connection Timeout" on the **Payment-Service** in **Region: US-East-1**.

**Analysis:** > This matches **Runbook #84 (Database Proxy Restart)** with 95% confidence. The logs suggest the proxy is hung due to an idle connection spike.

**Proposed Action:** > Trigger GitHub Action: `db-proxy-remediation.yml`

- **Target:** `proxy-01-prod`
- **Action:** Soft Restart

### 3. Why This is "Smarter" Than Traditional Automation

Feature	Traditional Automation	AI + Human-in-the-Loop
<b>Logic Type</b>	"Keyword" matching (fragile).	"Semantic" matching (understands intent).
<b>User Interface</b>	Raw logs or technical alerts.	<b>Human-Readable Summary</b> (Plain English).
<b>Risk Control</b>	All or Nothing (Auto-run or Manual).	<b>Permission-Based:</b> AI prepares, Human authorizes.
<b>Context</b>	Doesn't know if it worked before.	<b>Self-Correction:</b> AI can say: " <i>I tried this 10 mins ago and it failed; I recommend a different runbook now.</i> "

---

### 4. Technical Architecture: The Permission Flow

1. **Detection:** Splunk sends a raw alert to the **AI Middleware** (Python/FastAPI).
  2. **Reasoning:** \* AI summarizes the logs.
    - o AI queries a **Vector DB** containing 100+ Runbook descriptions.
    - o AI extracts parameters (Cluster ID, Namespace, etc.).
  3. **Interaction:** AI formats a "Proposal" and sends it to a Chatbot (Slack/Teams) with **Action Buttons**.
  4. **Permission:** The Human clicks "**Approve.**"
  5. **Execution:** The Middleware receives the "Approve" signal and calls the **GitHub repository\_dispatch API**.
  6. **Confirmation:** Once the GitHub Action finishes, the AI updates the chat thread: "*Success! Service is healthy again. View Logs [Link].*"
-

## 5. Value to the Organization

### A. Massive Reduction in MTTR (Mean Time to Resolution)

Engineers no longer need to log into Splunk, find the runbook, and then manually trigger a GitHub Action. The AI does the "homework" and puts the **"Fix Button"** right in front of them.

### B. "Zero-Code" Maintenance

When the team creates a **101st runbook**, they simply save a Markdown file in the repository. The AI automatically "reads" it and begins suggesting it when relevant alerts fire. **No new code is required to support new fixes.**

### C. Reduced Cognitive Load

On-call engineers are often tired or stressed. By providing a **Human-Readable Summary** and a **suggested fix**, the AI reduces the chance of human error during an outage.

## 6. Expected ROI

- **MTTR Reduction:** From ~20 minutes (triage) to **<2 minutes** (assisted remediation).
- **Engineering Efficiency:** Saves **20-25 hours/month** in manual mapping and script maintenance.
- **Cost Savings:** Reduces GitHub Action waste by **30%** via intelligent suppression of redundant triggers.