

[Open in app ↗](#)[≡ Medium](#) [Search](#) [Write](#) [3](#)

◆ Get unlimited access to the best of Medium for less than \$1/week. [Become a member](#) X



# Keeping Cool Under Pressure: How We Streamlined Spark Production Job Failure Handling for a Major Logistics Client



Anurag Sharma 6 min read · Aug 25, 2025



...

In the fast-paced world of logistics, where data and time are the lifeblood of operations, spark production job failures can send ripples of panic through teams, clients, and managers. As a team working for one of India's largest logistics companies, we faced this challenge head-on. Our client's operations relied heavily on timely and accurate data processing, and any hiccup in our production jobs could disrupt their supply chain, impacting everything from delivery schedules to customer satisfaction.

When production jobs failed, the ripple effect was immediate: the client's support team would escalate the issue to our manager, who would then scramble to get a developer to investigate. The catch? Accessing the on-premise Spark cluster required being on the company's network, meaning developers often had to rush to the office — sometimes late at night or on weekends. This wasn't just inconvenient; it created stress, disrupted work-life balance, and occasionally led to delays in resolving issues.

In this blog, we'll walk you through how we transformed this chaotic process into a streamlined, automated solution that kept clients, managers, and developers calm and in control, even when production jobs went awry.

## The Problem: A Reactive Fire Drill

Our logistics client processed massive datasets daily to manage inventory, optimize routes, and ensure timely deliveries. These datasets orchestrated through Apache Airflow, running on an on-premise spark cluster. When a job failed, the process to address it was far from ideal:

- 1. Client Escalation:** The client's support team would notice a job failure (e.g., missing data or delayed reports) and immediately contact our manager.
- 2. Manager's Response:** The manager would then call the developer responsible for the job, often outside regular hours.
- 3. Network Dependency:** Since the Spark cluster was only accessible within the company's network, developers had to physically come to the office or use a VPN (which wasn't always reliable) to investigate logs and diagnose the issue.
- 4. Panic and Delays:** This reactive approach created a sense of urgency and frustration. Sometimes, developers could resolve the issue quickly, but other times, they'd refuse to come to the office immediately, leading to delays until the next working day.

This wasn't sustainable. The constant firefighting was burning out the team, straining client relationships, and putting unnecessary pressure on everyone involved. We needed a proactive solution that would minimize disruptions, provide clarity on job failures, and allow developers to maintain a healthy work-life balance.



## The Solution: Automating Failure Detection and Communication

To tackle this challenge, we designed a robust, automated process that reduced manual intervention, provided immediate insights into job failures, and kept all stakeholders informed without requiring developers to rush to the office.

Here's how we did it:

### Step 1: Automated Log Scanning with Python

The heart of our solution was a Python script integrated with Apache Airflow, which monitored and analyzed Spark job failures. Since our cluster was on-premise, we could directly access the driver logs stored in a designated directory.

The script was designed to:

- **Trigger on Failure:** Using Airflow's failure callback mechanism, we configured the script to run automatically whenever a Spark job failed. The script received critical metadata, such as the job name, execution timestamp, and DAG ID, from Airflow.
- **Parse Driver Logs:** The script parsed the Spark driver's stdout logs to identify error messages, stack traces, or other relevant details. We used regular expressions to extract specific error patterns (e.g., "NullPointerException," "OutOfMemoryError," or custom business logic errors).
- **Run Diagnostic Queries:** To provide deeper context, the script executed a set of predefined diagnostic SQL queries against the source and target

datasets. These queries checked for common issues, such as:

*Is the source data count greater than zero?*

*Are there duplicate records on the business key in the source data?*

*Is the target table accessible and writable?*

*Are there any schema mismatches between source and target? The results of these queries were formatted as a simple true/false or numerical summary, making it easy to pinpoint the root cause.*

## Step 2: Automated Notifications via Email

Once the script gathered the necessary information, it compiled a concise yet detailed report and emailed it to the entire team, including developers, leads, the manager, and the client's support team.



The email included:

- **Job Details:** Job name, execution time, and DAG ID.
- **Error Summary:** A snippet of the error message or stack trace from the driver logs.

- **Diagnostic Results:** True/false or numerical outputs from the diagnostic queries.
- **Severity and SLA:** Based on predefined rules, the script classified the issue as “minor” (SLA: 1 day) or “major” (SLA: 2–3 days). For example, a missing source file might be flagged as minor, while a schema mismatch affecting downstream processes would be major.

The email subject was formatted for clarity, e.g., [PROD JOB FAILURE] Job XYZ Failed — Minor Issue (SLA: 1 Day).

### Step 3: Empowering the Team to Respond

With the email in hand, the developer or lead could quickly assess the situation without needing immediate cluster access. If the error was straightforward (e.g., a missing file or a simple configuration issue), they could reply to the email thread with their analysis and proposed fix, along with the severity level. For more complex issues, the team could plan their response within the defined SLA, avoiding the need for late-night office visits.



To ensure the solution remained effective, we:

- **Refined Error Patterns:** Regularly updated the script's regex patterns to catch new types of errors based on past failures.
- **Expanded Diagnostics:** Added more diagnostic queries as we identified recurring issues, such as network timeouts or resource contention.
- **Monitored SLA Adherence:** Tracked resolution times to ensure we met client expectations and used this data to fine-tune our process.

## The Impact: A Calmer, More Efficient Workflow

The automated failure detection and notification system transformed how we handled production job failures.

Here's how it made a difference:

- **Reduced Panic:** By providing immediate, actionable insights via email, we eliminated the need for frantic calls between the client, manager, and developers. Everyone had the same information at the same time, fostering transparency and trust.
- **Improved Work-Life Balance:** Developers no longer had to rush to the office to diagnose issues. They could review the email, assess the severity, and plan their response, often from home or during regular hours.
- **Faster Resolution Times:** The diagnostic queries and error summaries gave developers a head start on troubleshooting, reducing the time to identify and fix issues.
- **Stronger Client Relationships:** The client appreciated the proactive communication and clear SLAs, which reassured them that issues were being handled efficiently. This strengthened their confidence in our team.
- **Scalable Process:** The solution was flexible enough to handle new jobs and error types, making it a long-term fix for our growing workload.

## Lessons Learned

Building this process taught us a few key lessons:

1. **Automation is a Stress-Reliever:** Automating repetitive tasks like log scanning and notifications frees up mental bandwidth for problem-solving.

2. **Clear Communication is Key:** Providing all stakeholders with the same information at the same time reduces miscommunication and builds trust.
3. **SLAs Set Expectations:** Defining clear timelines for minor and major issues helps manage client expectations and prioritizes team efforts.
4. **Iterate and Improve:** Continuously refining the script and diagnostics ensured the solution stayed relevant as our systems evolved.

## Conclusion

In the high-stakes world of logistics, where every minute counts, production job failures can feel like a crisis. By implementing an automated failure detection and notification system, we turned chaos into calm. Our Python script, integrated with Airflow, empowered our team to diagnose issues quickly, communicate effectively, and resolve problems without sacrificing work-life balance. This solution not only kept our clients and manager happy but also made our team more efficient and resilient.

If you're dealing with similar challenges in your production environment, consider automating error detection and communication. It's a small investment that can yield big returns in team morale, client satisfaction, and operational efficiency.

Spark

Big Data

Automation

Data Analytics

Data Analysis