

CASE STUDY

Enhancing Observability in ELT Pipelines for Data Engineers

Overview

- Company: Hevo
- Role: Senior Product Designer
- Project Timeline: 2 Months
- Target Users: Data Engineers
- Platform: ELT, B2B SaaS

Problem Statement

Data Engineers use our ELT SaaS platform to build and manage data pipelines, extracting data from sources like Salesforce and LinkedIn Ads, then loading it into their Data Warehouses. However, the platform's limited observability presented several challenges:

- Delayed Failure Detection – Engineers only realize pipeline failures after downstream processes break, leading to reactive troubleshooting.
- Manual Debugging Efforts – Without intuitive monitoring tools, engineers must manually sift through extensive logs to diagnose issues.
- Performance Bottlenecks – Pinpointing latency sources and optimizing pipelines is difficult due to a lack of real-time insights.
- Data Quality Assurance – Maintaining data integrity across pipeline stages is challenging without proactive monitoring.

These challenges led to increased downtime, high operational costs, and user frustration. Our objective was to design a comprehensive observability feature that empowers Data Engineers with actionable insights.

Secondary Research

Industry Challenges in ETL Pipeline Observability

1. Limited Visibility into Pipeline Health

- Many organizations lack real-time insights into data flows, making it difficult to detect failures or performance issues before they impact downstream reports.

Source: Monte Carlo Data

2. High Manual Debugging Effort

- 40% of an engineer's time is often spent troubleshooting failures due to a lack of structured monitoring.

In a large-scale data environment, manual log analysis is time-consuming and inefficient.

Source: IBM Research

3. Data Quality & Schema Drift Issues

- Unexpected schema changes (e.g., column name changes or missing data) can break transformations, leading to inaccurate reports.

In a survey, 60% of companies reported that schema drift led to incorrect business decisions.

Source: Gartner Research

4. Performance Bottlenecks & Scaling Challenges

- As data volume grows, pipelines slow down due to inefficient query execution or resource limitations.

Without throughput and latency monitoring, engineers struggle to optimize execution time.

Source: Data Engineering Weekly

Primary Research

To deeply understand the challenges faced by Data Engineers, we conducted:

1. User Interviews:

- Participants: 3 Data Engineers from mid-sized companies.
- Methodology: Semi-structured interviews to explore pain points related to pipeline observability.
- Key Insights:
 - Engineers were frustrated with late failure detection.
 - Debugging was time-consuming due to manual log analysis.
 - A strong need for a centralized dashboard with real-time pipeline visibility.

2. Competitive Analysis:

- Platforms Analyzed: Fivetran, Apache Airflow, StitchData.
- Findings:
 - Competitors provided basic logging but lacked proactive alerts & historical trend analysis.
 - Users relied on third-party tools for enhanced observability.

3. Data Analysis:

- Reviewed support tickets & feedback from the past year.
- Findings:
 - 40% of complaints were related to unexpected pipeline failures.
 - Users requested real-time alerts, detailed error reports, and performance monitoring tools.

User Journey

To visualize how Data Engineers interact with the platform and experience pain points, we mapped their journey:

1. Pipeline Setup & Monitoring

- Engineers create and configure pipelines using various data sources and destinations.
- They expect real-time visibility into pipeline status, but the current system provides only limited insights.
- They frequently check logs and dashboards to verify if the pipeline is running correctly.

2. Pipeline Execution & Failures

- Pipelines process large volumes of data, but failures occur due to schema mismatches, connectivity issues, or data transformations.
- Engineers detect failures only after data discrepancies emerge in the warehouse.
- Debugging involves manually combing through logs without clear root cause identification.

3. Failure Identification & Troubleshooting

- Engineers struggle with unclear error messages and a lack of structured logs.
- They need to compare logs from different runs to identify failure patterns.
- Support tickets are often raised due to insufficient self-service debugging capabilities.

4. Resolution & Optimization

- Engineers attempt fixes through trial and error, leading to additional run delays.
- Performance bottlenecks are hard to detect without real-time performance monitoring.
- A need for an Error Insights Screen becomes evident to proactively diagnose and resolve errors.

User Stories & Problem Resolution

User Story 1: Proactive Failure Detection

As a Data Engineer, I want to receive real-time alerts when my pipeline encounters an issue, so that I can resolve it before it affects downstream processes.

Solution: The observability feature introduced real-time alerts and notifications, enabling engineers to detect failures instantly and proactively take action.

User Story 2: Faster Debugging with Error Insights

As a Data Engineer, I want to view a detailed breakdown of pipeline errors so that I can quickly identify and resolve issues.

Solution: The Error Insights Screen provided categorized error logs, severity levels, and suggested fixes, allowing engineers to troubleshoot problems efficiently.

Success Metrics

- Adoption & Engagement Metrics
 - Feature Adoption Rate: % of users interacting with the new observability dashboard.
 - Pipeline Overview Page Views: Tracks how frequently engineers access detailed pipeline insights.
 - Alerts Interaction Rate: % of users who engage with alerts (click on notifications, adjust preferences).

Efficiency & Performance Metrics

- Mean Time to Detect (MTTD): Time taken to identify pipeline failures before vs. after the observability update.

We created mind maps linking common user challenges to potential solutions, ensuring every pain point had a corresponding feature enhancement.

Solutioning & Implementation. To address these challenges, we developed a comprehensive observability framework with three primary components:

Solutioning

1. Pipeline Listing Page (High-Level Overview)

- Aggregated Metrics:
 - Pipeline Health (Donut Chart) – A quick snapshot of running, failed, and pending pipelines.
 - Net Consumption – Number of events processed, converted to dollar value.

2. Pipeline Overview Page (Deep-Dive View)

- Top Section (Pipeline Metadata):
 - Pipeline Name, ID, Status, Last Run, Frequency of Runs.

3. Run History Page (Tracking Past Runs)

4. Error Insights Screen (Proactive Debugging)

- Summary Section:
 - Pipeline Name & ID
 - Run Date & Timestamp
 - Execution Status & Total Errors
 - Error Rate (%) & Throughput Impact

5. Error Breakdown:

- Visual Distribution of errors (Source, Destination, Transformation Errors)

6. Error Trends Over Time

7. Detailed Error Logs (Table View):

- Error Type, Message, Affected Records, Severity Level

8. Suggested Fixes:

9. Actionable Steps:

- Retry Run, Open Support Ticket, Compare Similar Errors, Download Logs

Problem Solved:

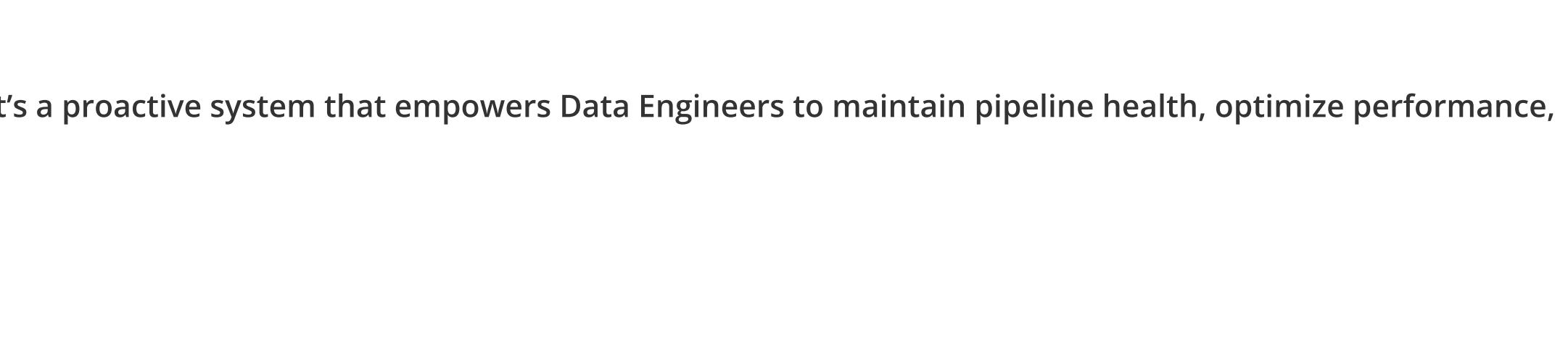
Engineers lacked a quick, aggregated view of all their pipelines, making it difficult to identify failing or under-performing pipelines at a glance.

- Displays a high-level overview of all active pipelines.

Pipeline Health Donut Chart provides an instant snapshot of healthy, degraded, and failed pipelines.

Net Consumption & Performance Metrics (Avg Execution Time, Peak Throughput, Avg Throughput, Avg Lag, Total Downtime, Avg Error Rate) help engineers assess performance trends.

Pipeline Cards display source, destination, and key throughput data, allowing engineers to quickly drill into problem areas.



Problem Solved:

Engineers needed detailed insights into individual pipelines but had to rely on raw logs to understand performance and issues.

- Pipeline Metadata Section displays essential details like Pipeline ID, Status, Last Run Time, and Frequency.

Charts for Event Consumption & Latency Trends help engineers visualize data flow patterns.

Performance Metrics (Avg Lag, Avg Throughput, Peak Throughput, Downtime) allow engineers to analyze pipeline health over time.

Key Metrics:

- Avg Lag, Avg Throughput, Peak Throughput, Total Downtime

Detailed Error Logs:

- Date & Timestamp

Execution Status (Success/Failed)

Execution Time

Errors Detected

Records Processed

Avg Throughput & Peak Throughput

Run History Page (Tracking Past Runs)

Run Details:

Event Consumption:

Latency:

Average Throughput:

Peak Throughput:

Average Lag:

Downtime:

Problem Solved:

Engineers struggled to diagnose failures due to lack of structured error categorization and root cause analysis.

- Error Breakdown Graphs (Pie chart or Bar Graph) categorize errors into Source Errors, Transformation Errors, Destination Errors.

Trend Graph highlights spikes in errors over time.

Performance Metrics (Avg Lag, Avg Throughput, Peak Throughput, Downtime) help engineers analyze pipeline health over time.

Key Metrics:

- Avg Lag, Avg Throughput, Peak Throughput, Total Downtime

Detailed Error Logs:

Event Type:

Affected Records:

Severity Level:

Error Message:

Suggested Fixes:

Actionable Steps:

Retry Run, Open Support Ticket, Compare Similar Errors, Download Logs

Problem Solved:

Engineers lacked a quick, aggregated view of all their pipelines, making it difficult to identify failing or under-performing pipelines at a glance.

- Displays a high-level overview of all active pipelines.

Pipeline Health Donut Chart provides an instant snapshot of healthy, degraded, and failed pipelines.

Net Consumption & Performance Metrics (Avg Execution Time, Peak Throughput, Avg Throughput, Avg Lag, Total Downtime, Avg Error Rate) help engineers assess performance trends.

Pipeline Cards display source, destination, and key throughput data, allowing engineers to quickly drill into problem areas.

Problem Solved:

Engineers needed detailed insights into individual pipelines but had to rely on raw logs to understand performance and issues.

- Pipeline Metadata Section displays essential details like Pipeline ID, Status, Last Run Time, and Frequency.

Charts for Event Consumption & Latency Trends help engineers visualize data flow patterns.

Performance Metrics (Avg Lag, Avg Throughput, Peak Throughput, Downtime) allow engineers to analyze pipeline health over time.

Key Metrics:

- Avg Lag, Avg Throughput, Peak Throughput, Total Downtime

Detailed Error Logs:

Event Type:

Affected Records:

Severity Level:

Error Message:

Suggested Fixes:

Actionable Steps:

Retry Run, Open Support Ticket, Compare Similar Errors, Download Logs

Problem Solved:

Engineers struggled to diagnose failures due to lack of structured error categorization and root cause analysis.

- Error Breakdown Graphs (Pie chart or Bar Graph) categorize errors into Source Errors, Transformation Errors, Destination Errors.

Trend Graph highlights spikes in errors over time.

Performance Metrics (Avg Lag, Avg Throughput, Peak Throughput, Downtime) help engineers analyze pipeline health over time.

Key Metrics:

- Avg Lag, Avg Throughput, Peak Throughput, Total Downtime

Detailed Error Logs:

Event Type:

Affected Records:

Severity Level:

Error Message:

<h3