# HPC Consult Ticket Analysis with SambaNova

**HIGH PERFORMANCE COMPUTING**

Author: Daisy Nsibu ; Mentors: Mike Mason, Tolulope Olatunbosun | HPC-SYS

## Overview

### Background

The High Performance Computing (HPC) Division uses RT, an older ticketing system that tracks user issues and interactions with the HPC consultants. With over 100,000 tickets and decades of interactions there is an amazing amount of useful information within those tickets however, the old ticket system makes it difficult to get that information out. The consultants want to know how this underutilized data can be used to better understand and serve their users.
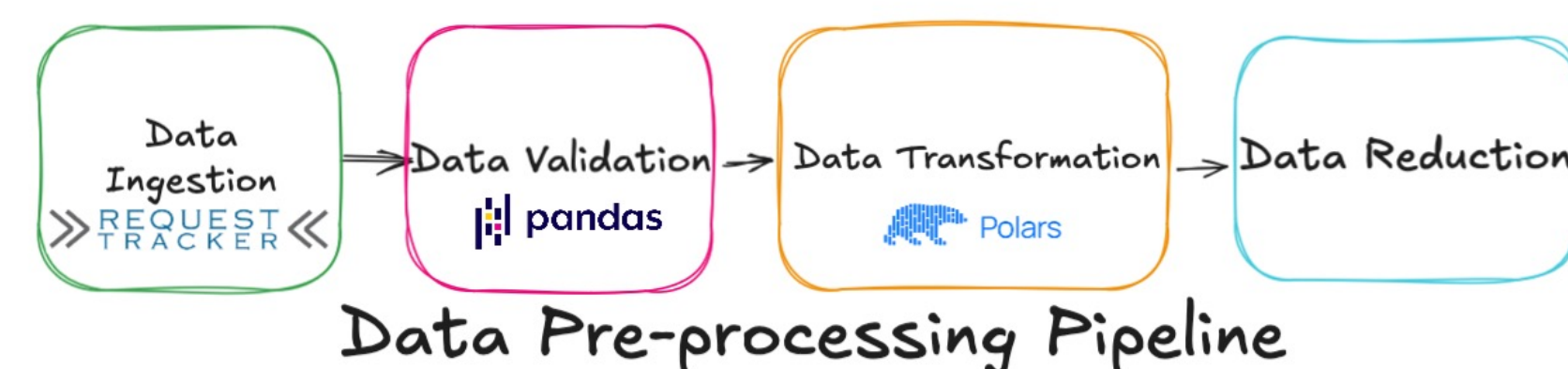
### Objectives

- Provide a better method to searching for tickets
- Develop a tool that uses AI/ML that provides comprehensive analysis of consult queue tickets.
- Use LLMs hosted by SambaNova to analyze the full ticket correspondence efficiently

## Methods

### Data Acquisition & Preprocessing

- Extracted over 100,00 tickets from 2008 to present day
- Source: Request Tracker (RT) ticketing system



Data Pre-processing Pipeline

### AI Implementation

- Utilized 🟣 **SambaNova** platform, powered by Reconfigurable Dataflow Units (RDUs)**,** for high-speed inference capabilities
- Implemented ∞ **Llama 3.3-Instruct-70B** model to:
  - Generate summaries
  - Predict and generate tickets categories
  - Perform sentiment and ticket analysis

### Web Application Development

- Built interactive dashboard using 👑 **Streamlit** framework
- Integrated 📊 **Plotly** for advanced interactive visualizations

## Conclusion

### Conclusion

With over 100,000 tickets and decades of interactions collected, we were able to develop the first of its kind, AI powered web application that consultants can use to easily search for tickets as well as explore the data.

### Future Directions

- Process ticket data using an advanced LLM:Llama 4 Maverick
- Live Ticket Integration: The web app connects to the RT API to collect real-time tickets
- Feedback ticket summaries, sentiment and predicted category to RT

## Results

These columns were generated using the LLM via SambaNova. This process involved developing and optimizing prompts for the model, then implementing batch processing to efficiently handle the data volume. Batch jobs were scheduled during off-peak hours (6:30pm-6:30am) to maximize computational resources.



*Figure 1: Interactive Dashboard Summary View with Dynamic Filtering Capabilities*



*Figure 2. Sentiment Over Time View of Filtered Data*



*Figure 3. Zoomed in View of the Distribution of Ticket Categories*

| RANK | CATEGORY | LLM-PREDICTED CATEGORY | LLM-GENERATED CATEGORY |
|---|---|---|---|
| 1 | ACCOUNTS | SOFTWARE RELATED | CLUSTER MAINTENANCE |
| 2 | MODULEFILES | AUTHENTICATION | PROJECT CREATION AND MANAGEMENT |
| 3 | LUSTRE | ACCOUNTS | STORAGE-OTHER |
| 4 | ENVIRONMENT- SYSTEM | ALLOCATIONS | SCHEDULED MAINTENANCE |
| 5 | WLM EDUCATION | LUSTRE | PROJECT CREATION |
| 6 | AUTHENTICATION | DST | USER ACCOUNT MANAGEMENT |

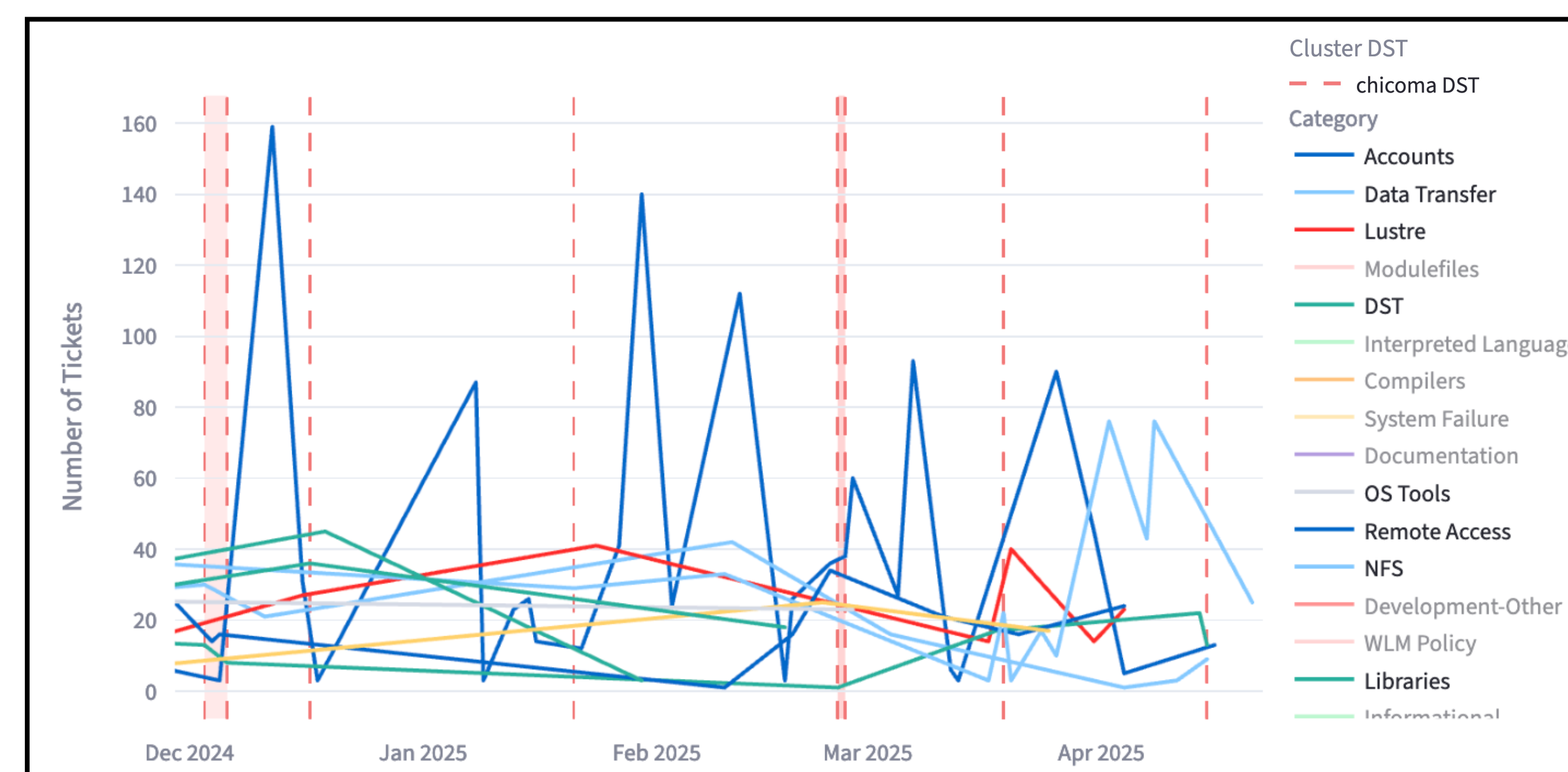*Table 1. Top 6 Consultant-Defined Categories with LLM-Predicted and LLM-Generated Categories*



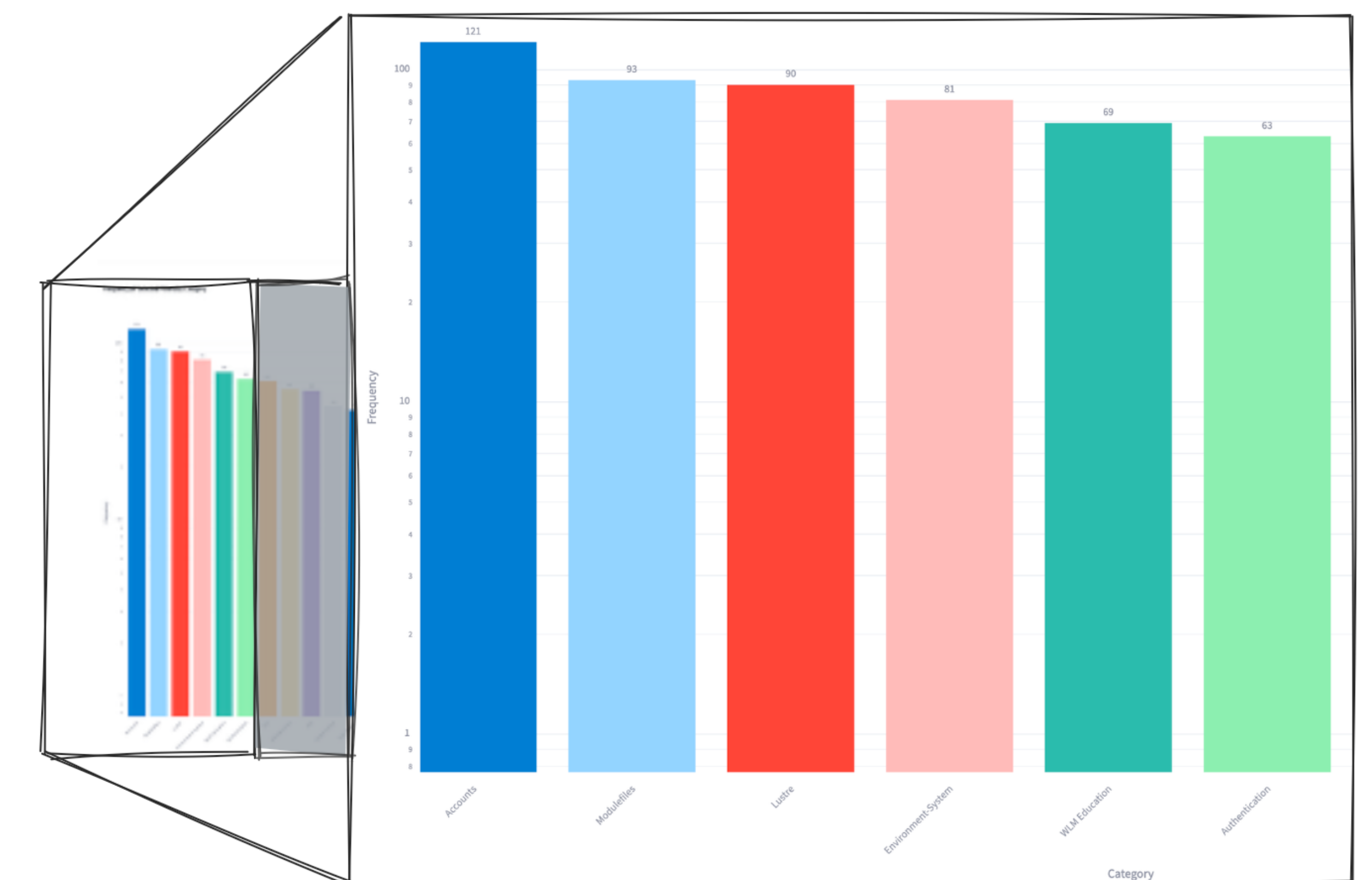*Figure 4. Number of Tickets Created by Date and Category with Chicoma DST Overlay*

1. Account Access Issues (533 tickets, 33.8%)

**User Question:**
Users frequently report being unable to log into their accounts on specific clusters. Common symptoms include authentication failures, expired credentials, or permission issues.

**Resolution:**
The standard resolution involves verifying account status in the user database, resetting credentials if necessary, and ensuring proper group permissions are set. For expired accounts, the renewal process includes contacting the PI or account manager to request an extension.

2. Job Submission and Scheduling Issues (400 tickets, 25.4%)

**User Question:**
Users experience issues with submitting jobs, including errors with job scripts, scheduling conflicts, and node failures.

**Resolution:**
The resolution involves troubleshooting the job script, checking for scheduling conflicts, and investigating node failures. In some cases, users may need to adjust their job submission scripts or seek assistance from the HPC consulting team.

3. Cluster Maintenance and Downtime (130 tickets, 8.2%)

**User Question:**
Users frequently report issues related to cluster maintenance and downtime. Common symptoms include unexpected downtime, delayed job execution, or resource unavailability.

**Resolution:**
The standard resolution involves providing updates on cluster maintenance schedules, notifying users of planned downtime, and ensuring that clusters are returned to service as soon as possible after maintenance.

*Figure 5. Top 3 User Issues and Resolution for Filtered Data*