# San Francisco Crime Pattern Analysis with Snowflake, Airflow, dbt, and PowerBI

Pranav Reddy Gaddam[#1], Anvay Bhanap[*2]
Yashwanth Reddy Katipally[#3], Thejes Raj Gangadhar[*4]

#*Department of Applied Data Science, San Jose State University*
*1 Washington Sq, San Jose, CA 95192, USA*

[1]`pranavreddy.gaddam@sjsu.edu`

[2]`anvay.bhanap@sjsu.edu`

[3]`yashwanthreddy.katipally@sjsu.edu`

[4]`thejesraj.gangadhar@sjsu.edu`

*Abstract –* **This project aims to investigate the crime patterns present in San Francisco by analyzing incident reports from 2018 to 2024. Utilizing both historical and real-time data information, this study explores temporal trends, geographic crime distributions, and resolution effectiveness. The process integrates Snowflake for warehousing, Airflow for ETL workflows, dbt for ELT transformations, and PowerBI for dashboarding and visualizations. Findings reveal noticeable trends, including peak incidents during summer months, spatial clustering in areas with high population density, and an overall doubling of crime post-COVID. The proposed solution highlights how actionable insights can guide effective resource allocation and strategic policymaking, enabling the city to enhance its crime management efforts.**

*Keywords -* **ETL, ELT, COVID, dbt, DAG, Socrata**

## I. PROBLEM STATEMENT

San Francisco's crime rates are an everlasting concern for residents, tourists, and policymakers, with recent trends indicating further challenges. As of this year, the city has experienced a complex crime landscape characterized by persistent property crimes, increased drug-related offences, and significant incidents of theft and violence. Understanding the yearly and locational patterns of crime is essential for preventive measures and resource distribution to take place. Despite available data, deriving insights and conducting effective actions based on the data alone remains challenging. The objective of this report is to present a solution that can analyze the data over the past 6 years and glean patterns by time, area, and incident category to provide government officials and law enforcement with an evidence-based understanding of what measures need to be taken to make strides in crime prevention and policy development.

## II. INTRODUCTION

As crime trends evolve in urban areas, analyzing incident reports is crucial for proactive safety measures. San Francisco's social environment in 2024 reflects a complex outlook, from the staggering rise in property crimes and organized retail theft. Recent data also suggests a persistent rise in car break-ins, with over 6000 incidents reported in the first quarter of 2024 alone, and a notable surge of drug peddling in open air markets, especially in neighborhoods like the Tenderloin and South of Market. Violent crimes, though not increasing as drastically, remain a critical concern, with assault and robbery rates presenting a continuous challenge for law enforcement. San Francisco's unique socio-economic state, marked by its high-density localities and financial and cultural diversity, makes it a particular challenge to address issues in crime management. This project adopts an ELT approach to process crime incident data, utilizing Airflow for scheduling tasks and dbt for refining data transformations. By integrating these datasets into PowerBI dashboards, hidden trends can be brought to light, such as crime hotspots and seasonal spikes, supporting efforts to improve community safety.

## III. RESEARCH

San Francisco's crime landscape presents a complex challenge for residents, tourists, and policymakers alike. Persistent property crimes, increased drug-related offenses, and rising incidents of theft and violence call for a more detailed analysis of crime patterns across different times, locations, and categories. Several studies have contributed to understanding these dynamics, applying advanced methodologies such as machine learning, spatial analysis, and complex network analysis to uncover crime trends and provide actionable insights. Machine learning has been instrumental in crime prediction and analysis. For example, Rasool et al. [1] used machine learning models to forecast crime trends in San Francisco, highlighting how such techniques can aid law enforcement in anticipating crime occurrences and optimizing resource allocation. Similarly, Pradhan [2] conducted exploratory data analysis (EDA) on crime data, identifying key temporal and spatial patterns that could inform crime prevention strategies. The ability to predict crime occurrences using machine learning models helps policymakers address criminal activity proactively, rather than reactively. In addition to temporal analysis, spatial patterns play a significant role in understanding crime. Chen et al. [3] explored how changes in human activities, especially during the COVID-19 pandemic, influenced crime patterns in San Francisco. They found that stay-at-home mandates led to notable shifts in crime distribution, indicating the importance of integrating human behavior dynamics when analyzing crime data. This spatio-temporal understanding of crime is vital for resource allocation and planning effective prevention measures, as it helps target high-crime areas during specific times.

Complex network analysis offers another powerful tool for uncovering hidden patterns in crime data. Spadon et al. [4] applied this methodology to study the interconnections between crime events in urban areas. By treating crime data as a network, they were able to identify crime hotspots and areas with high criminal activity. This approach provides a more granular understanding of crime occurrences, which is essential for law enforcement agencies seeking to prioritize intervention in specific neighborhoods. Moreover, advanced statistical models like the Gaussian Cox process have been used to model crime events over both time and space. Shirota and Gelfand [5] applied this approach to analyze crime in San Francisco, focusing on how crime events unfold in a given location over time. Their research contributes to developing predictive models for identifying crime hotspots, helping law enforcement agencies adjust their strategies accordingly. These studies collectively demonstrate that analyzing crime through machine learning, spatial analysis, complex networks, and statistical models provides significant insights into criminal behavior in San Francisco. By applying these methodologies, policymakers and law enforcement agencies can craft evidence-based strategies for crime prevention, better allocate resources, and address the evolving challenges in urban crime dynamics.

## IV. METHODOLOGY

Analyzing crime trends in San Francisco needs a strong pipeline to integrate historical and live data, facilitating workable insights for city officials. With the city facing increasing crime rates, understanding the dynamics that cause such results is crucial for strategic financial allocation and law enforcement planning. The methodology detailed in the following sections provides a comprehensive overview of how these crime hotspots and yearly trends are found by integrating historical datasets with recent updates, securing an accurate and insightful crime management solution. Figure 1 shows how the ETL pipeline was constructed using Airflow, dbt, and PowerBI, implementing a straightforward process from data extraction to visualization. Initially, crime data was extracted from the San Francisco government data portal and loaded into Snowflake. Airflow automated the extraction and stages the processes through DAG to

maintain the scheduled updates. This ensured that historical data and live updates could coexist in a structured and accessible environment, prepared for the next phase. During transformation, dbt refined and aggregated the data to produce a structured dataset that highlighted crime patterns by district and category. This dbt framework optimized transformations directly into Snowflake, which maintained a high level of efficiency by avoiding unnecessary data transfers. Incremental updates were implemented using the datetime library to make sure that only the new records were processed, which significantly reduced resource consumption and improved system performance. Then, PowerBI dashboard visualized the collected data, and displayed interactive views of crime patterns, hotspots, and incident categories. Policymakers can explore the time-based trends to empower city officials to work towards solutions and preventative measures. These visualizations can aid strategic decision-making, confirming that resources are granted intentionally to areas with higher crime incidence, thereby optimizing public safety efforts and fostering stronger community engagement. Further implementation of ongoing updates was achieved through an ELT approach. Live data was integrated using Socrata API with raw entries loaded directly into Snowflake. Previously described aggregation methods were handled downstream in dbt. The model optimized data integration by using datetime based incremental updates, which minimized computational resources compared to ETL's traditional row-by-row merging. This dual implementation facilitated comprehensive data processing, which accommodated both historical records and real-time information while maintaining analytical scalability.



Fig. 1. System diagram of data pipeline for crime trend analysis using SFODP, Airflow, dbt, and Snowflake

A. *Extracting Historical Data from San Francisco Open Data Portal (see Python code)*

In order to establish the foundation for crime trend analysis, the project was initiated by sourcing historical data of the past 6 years from the San Francisco Open Data Portal (SFODP) as a CSV file. The file consisted of a variety of fields referencing the day and time of the event – Incident Datetime, Report Datetime, Incident Day of Week –, the type of crime that occurred –

Incident Category, Subcategory, Report Type, Resolution –, and the area in which it happened – latitude, longitude, neighborhood – as well as the nearest police district. All together, there were 862,000 rows of information gathered from all San Francisco neighborhoods. The preprocessing began with cleaning the data to remove missing values and standardizing the datetime format. Additionally, irrelevant columns, like Dispatch Number and Supervisor District, were removed from the working dataset as they were irrelevant to the overarching story being told through the data. By the end of the cleaning process, the 34 columns were reduced to a usable 13, as shown in Table 1. The resulting dataset was uploaded into a Snowflake data warehouse using its staging feature, ensuring secure storage. Though the data extraction process was manual, it was automated for consistency using a SnowSQL script integrated with Airflow. A Directed Acyclic Graph (DAG) was developed to regulate this upload procedure (Fig. 2), allowing for seamless integration into Snowflake. While this DAG was designed as a one-time execution, because of the static nature of a historical data CSV, it showcased the project's scalable approach to handling similar datasets in the future.

TABLE I
DATA TYPES FOR INCIDENTS TABLE

| Column Name | Data Type | Constraints |
| --- | --- | --- |
| Incident ID | VARCHAR | PRIMARY KEY (datetime, ID) |
| Incident Datetime | TIMESTAMP | |
| Incident Category | VARCHAR | |
| Incident Day of Week | VARCHAR | |
| Analysis Neighborhood | VARCHAR | |
| Incident Subcategory | VARCHAR | |
| Intersection | VARCHAR | |
| Latitude | VARCHAR | |
| Longitude | VARCHAR | |
| Police District | VARCHAR | |
| Report Datetime | TIMESTAMP | |
| Row ID | VARCHAR | |
| Resolution | VARCHAR | |

Fig. 2. Airflow logs of loading historical data to Snowflake

*B. Implementing ETL Processes with dbt (see SQL code)*

Once the historical data was ingested, dbt (Data Build Tool) was implemented to perform the essential ETL processes to refine and transform the raw data. The transformations were centered on aggregating crimes by category and district (Fig. 3, Fig. 4) allowing the data to present associated patterns for any upcoming decision-making. The modular SQL-driven framework provided by dbt enabled all the creation of reusable models to dynamically adjust updates or new data sources. Furthermore, dbt facilitated data quality checks (Fig. 5, Fig, 6, Fig. 6) to ensure that accuracy and integrity were maintained prior to the analysis phase. By performing these transformations inside Snowflake, it was possible to keep the data transfer overhead to a minimum and focus on improving efficiency. The result generated a structured dataset ready for visualization and deeper exploration.



Fig. 3. dbt snapshot of incident categories in Snowflake

Fig. 4. dbt output table descriptions in Snowflake
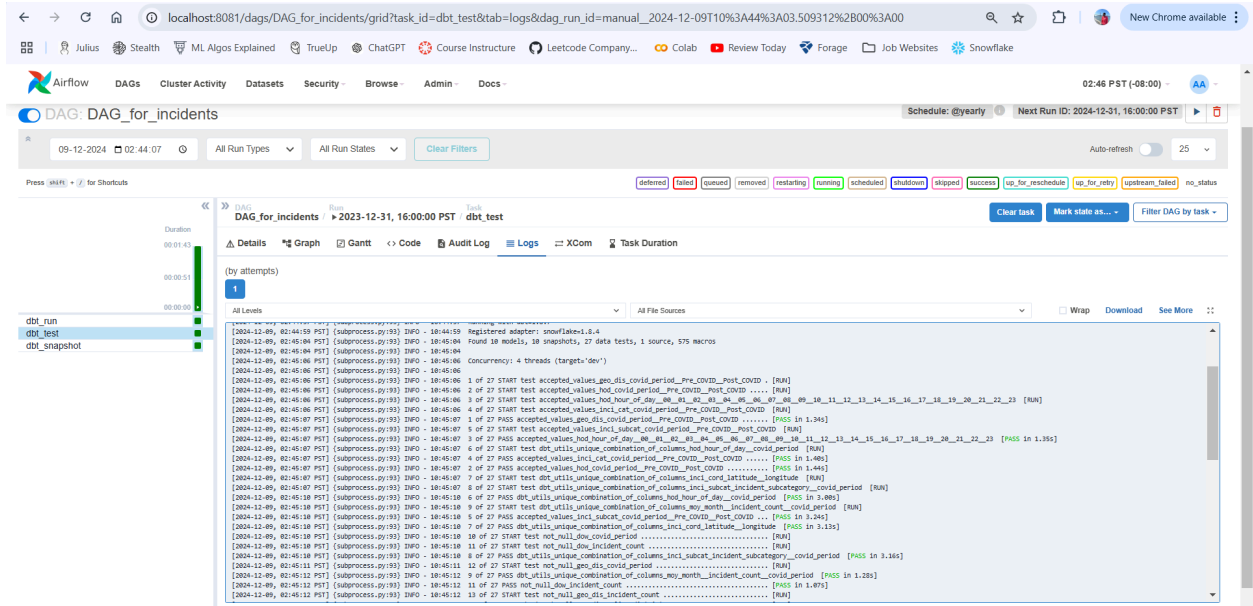


Fig. 5. Airflow logs of running all dbt models
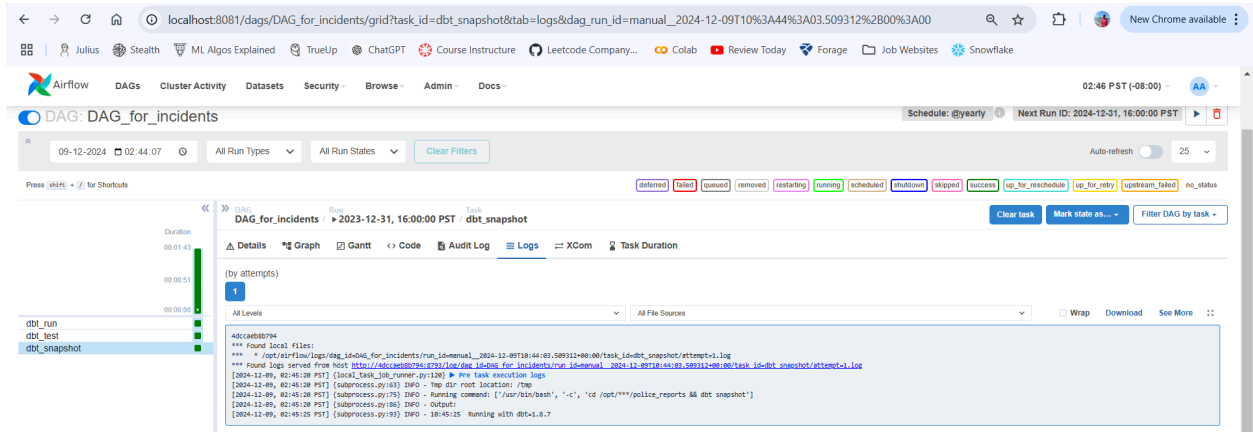
Fig. 6. Airflow logs of running all dbt tests


Fig. 7. Airflow logs of running all dbt snapshots

*C.  Creating Visualizations in PowerBI*

The transformed data was loaded into PowerBI to design an engaging and interactive dashboard (Fig. 8). The display visualized key crime statistics, like incidence amounts based on time in the year, localities, and category types. By implementing machine learning techniques to advance the visualization, geographic clustered heat maps were created, along with time series graphs that compared the difference in incidence pre- and post-COVID. Breaking down the dashboard, Figure 9 shows a bar chart presenting a trend of reported incidents based on years, though the scale can be adjusted to show daily and monthly trends, as well as those based on day of the week. The map presented in Figure 10 illustrates the spatial distribution of incidents, referencing areas with varying levels of activity with total incident counts in accordance with the temporal measurement in the previous bar chart. The two metrics in the bottom left corner of the dashboard present the total number of incidents since 2018 and the most dangerous San Francisco district, followed by the number of incidents that took place there. The grouped bar chart in Figure 11. shows the comparison of incident counts across the days of the week, along with how the amount of cases have changed after 2020. Finally, the line chart in the bottom right

corner (Fig. 12) also compares the difference in incident counts before and after COVID, but does so by hour of day. These insights serve to empower city officials in prioritizing resource allocation, identify areas in need of intervention, and formulate strategies for long-term crime reduction.
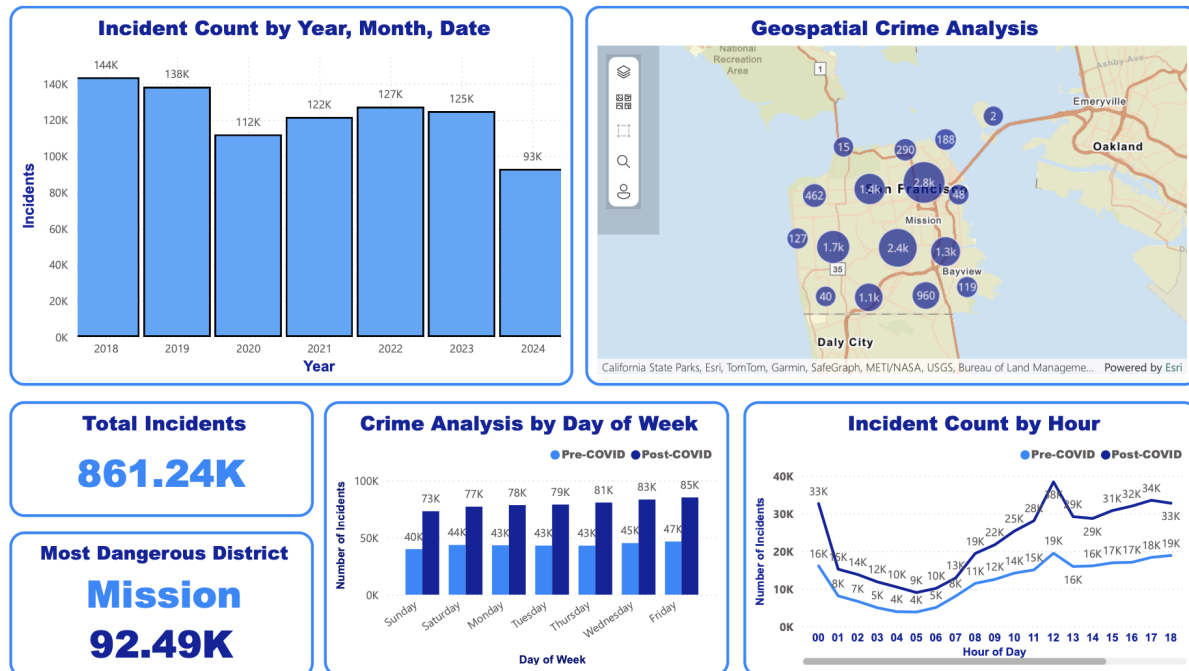


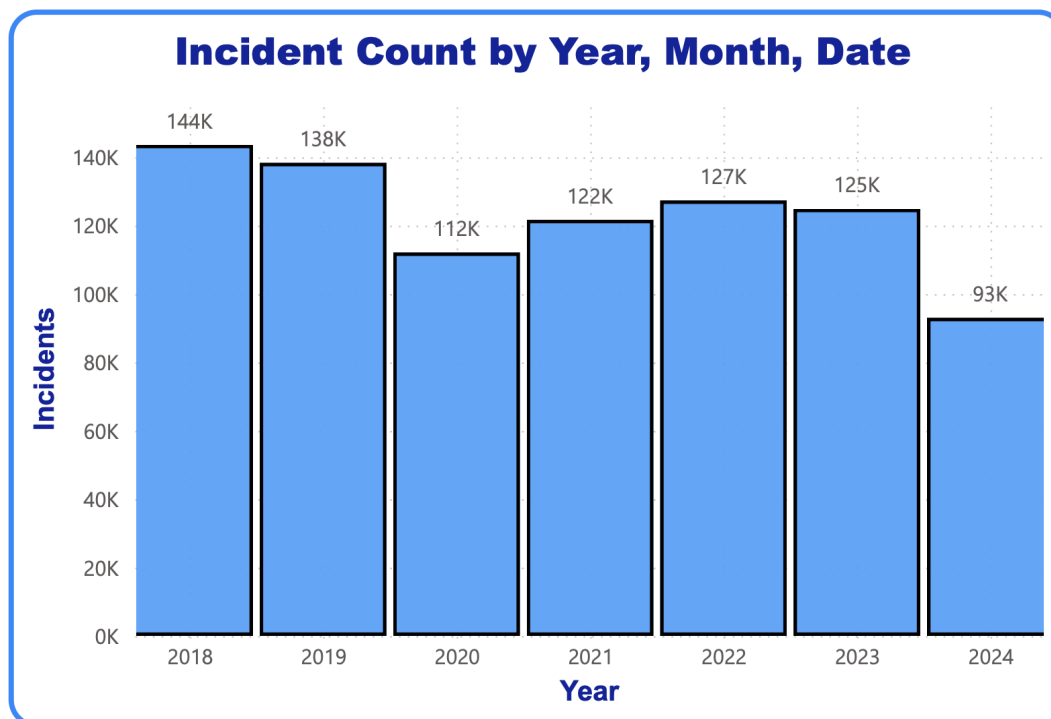Fig. 8. PowerBI dashboard depicting San Francisco Crime Analysis



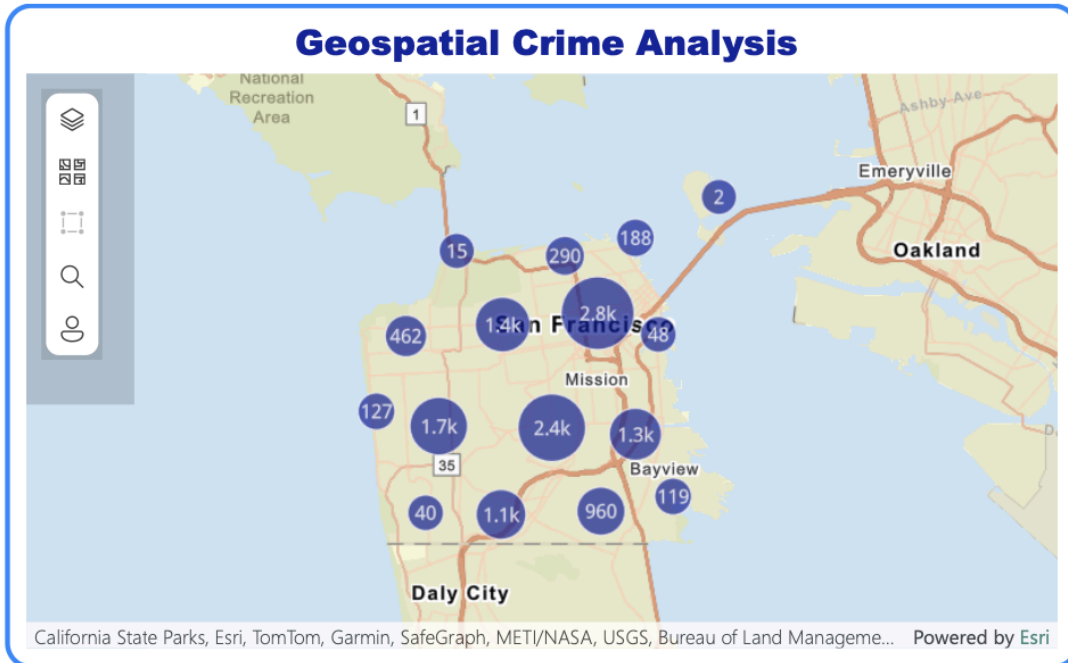Fig. 9. Incident Count by Year, Month, Date

Fig. 10. Geospatial Crime Analysis showing crime hotspots



Fig. 11. Crime Analysis by Day of Week (pre- and post-COVID)

Fig. 12. Incident Count by Hour (pre- and post-COVID)

### D. Integration of Live Governmental Data from Socrata (see *Python code*)

To complement the collection and representation of the historical dataset, new API integration was added from Socrata, an open data platform that allows users access to government data. This enhancement allowed the Snowflake table to maintain an updated repository of crime statistics. Figure 13 depicts how the Airflow DAG was modified to schedule and run weekly, ensuring that the system maintains current incident reports. Unlike the one-and-done execution for historical data, this DAG implements dynamic scheduling to automate the retrieval and ingestion process continuously.



Fig. 13. Airflow logs of loading API data to Snowflake

Instead of using the traditional MERGE statement for handling duplicates, a more effective strategy was used for incremental updates. The MERGE function is reliable in most cases but here it consumes a significan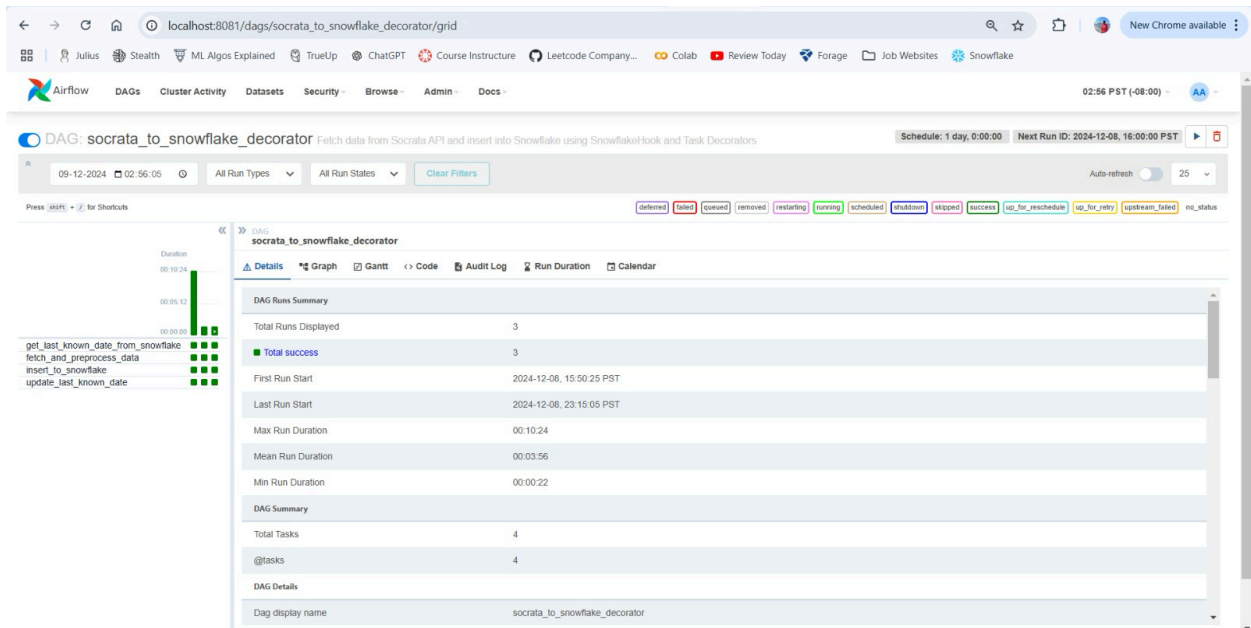t amount of computing resources as it individually processes each row of the dataset. To avoid this limitation, the datetime library was employed to track the maximum timestamp in the current table. Python's timedelta() function allowed the DAG to calculate the next starting point for ingestion by adding one minute to the maximum timestamp. Each DAG execution updated the last_known_time variable, ensuring that only the new records were fetched and integrated. This approach significantly reduced resource usage and improved the pipeline's performance. To ensure that the DAG checks entries properly, error handling has been implemented using the try catch method where if no new records are found, task implementation will automatically complete and this has been ensured for all the following tasks as well.

*F. Automating Data Consistency and Accuracy*

The updated DAG ensured consistency and data accuracy through thorough testing before deployment. Incremental updates were validated against the raw data to confirm the accuracy of the last_known_time mechanism. The utility of Airfow's task dependency management further confirmed that updates occurred in the correct sequence, restricting any duplication or data loss. This automated solution enabled the pipeline to maintain a balance between efficient and reliable transformation while seamlessly handling large and dynamic datasets.

## V. RESULTS

The dashboard offers a multifaceted presentation of crime data through a variety of useful visualizations. The Incident Count by Year, Month, Date bar chart illustrates a yearly trend in reported incidents, peaking in 2019 at 144K and declining to an expected average of 125K in 2024 (recognizing that the 93K does not include the total for December). The Geospatial Crime Analysis clustered heatmap highlights the geographic distribution of crime, with hotspots appearing in Financial District, Mission, and the Tenderloin., visualized with proportional circles and yearly crime amounts of 2.8K and 2.4K. The overall metrics in the bottom left corner summarize the number of incidents over the past six years, totalling at 861.24K, while the most dangerous district identified Mission with the highest number of incidents: 92.49K. The Crime Analysis by Day of Week bar chart compares the incident counts across different days, distinguishing between pre-COVID and post-COVID patterns, with the highest activity observed towards the end of the work week (Wednesday through Friday). Lastly, the Incident Counts by Hour line chart examines hourly trends, showing the bulk of incidents during the afternoon, peaking at 12 PM, following a gradual rise in the evening, from 4 – 6 PM, with another peak around midnight.

## VI. ANALYSIS

The development of the San Francisco Crime Analysis project presented several challenges that required innovative solutions to ensure efficiency and scalability. One significant hurdle was the API limitations, as the Socrata API only allowed retrieval of 50,000 records per query. To overcome this, the complete historical dataset was downloaded in CSV format for bulk processing, while the API was reserved for live incremental updates, adhering to its record limits. Another challenge involved the inefficiency of row-by-row `INSERT INTO` statements for

large datasets, which resulted in slow processing times. This was resolved by adopting a staging approach, where historical data was temporarily staged in Snowflake before being ingested in bulk using `COPY INTO` commands, significantly reducing runtime. For live data updates, the initial use of `MERGE` statements proved overly complex and computationally expensive. Instead, a more streamlined last-known-date approach was implemented, fetching only new records by querying the API based on the latest `incident_datetime` in Snowflake. Additionally, data preprocessing revealed critical fields with null or missing values, which were systematically handled by dropping invalid rows and validating the dataset before ingestion. Another challenge was ensuring historical data integrity while enabling time-series analysis. This was addressed through dbt snapshots, which captured and preserved changes in historical data. These strategies collectively ensured the pipeline's robustness, enabling seamless integration of both historical and live data while maintaining accuracy and efficiency. The project highlights the importance of balancing performance optimization with effective data handling for real-world analytical pipelines.

## VII. CONCLUSION

In conclusion, this project was able to demonstrate the potential of integrated data warehousing and analytics in addressing urban safety challenges. Data-driven approaches have already shown promising results in the city's crime management. For example, the SFPD used predictive policing technologies and targeted analytics to direct a 15% reduction in property crimes between 2021 and 2023, especially in the high-risk neighborhoods of the Tenderloin and Mission District. Through the implementation of data-informed strategies, law enforcement have been able to deploy resources more purposefully, learning to proactive crime prevention. By combining historical and real-time datasets, the project uncovered genuine insights into the common incidental patterns. Future work can focus on predictive modeling to anticipate high-risk periods and integrating socio-economic data for a hologist understanding of crime drivers. The results will inevitably inform policymakers in crafting helpful and effective solutions based on defined data to build a safer San Francisco.

## VIII. REFERENCES

[1]     AzharGhafoor, "Forecasting the Trends and Patterns of Crime in San Francisco using Machine Learning Model," ResearchGate, https://github.com/AzharGhafoor/PDF_XSS_PAYLOADS (accessed Dec. 9, 2024).
[2]     I. Pradhan, "Exploratory Data Analysis and crime prediction in San ...," Exploratory Data Analysis And Crime Prediction In San Francisco, https://scholarworks.sjsu.edu/cgi/viewcontent.cgi?referer=&httpsredir=1&article=1633&context=etd_projects (accessed Dec. 9, 2024).
[3]     T. Chen, "Spatio-temporal stratified associations between urban human activities and crime patterns: a case study in San Francisco around the COVID-19 stay-at-home mandate," ResearchGate, https://www.researchgate.net/publication/361117742_Spatio-temporal_stratified_associations_between_urban_human_activities_and_crime_patterns_a_case_study_in_San_Francisco_around_the_COVID-19_stay-at-home_mandate (accessed Dec. 9, 2024)..
[4]     G. Spadon et al., "Complex network tools to understand the behavior of criminality in urban areas," arXiv.org, https://arxiv.org/abs/1612.06115 (accessed Dec. 9, 2024).
[5]     S. Shirota and A. E. Gelfand, "Space and circular time log Gaussian Cox processes with application to crime event data," arXiv.org, https://arxiv.org/abs/1611.08719 (accessed Dec. 9, 2024).