

# IE 7374 ST: BI Data Integration

Accelerating Insights:

A Cloud-Based Data Pipeline for F1 Racing Analytics

Group 8

Anurag Palanki

Rohith Kanakagiri

## **INTRODUCTION**

Formula 1 (F1) is a world-renowned motor racing championship that has captivated millions of fans since its inception in 1950. Over the decades, the sport has evolved dramatically, with cutting-edge technology and innovation driving performance and shaping the industry's landscape. F1 racing teams and decision-makers rely on data-driven insights to optimize performance, increase efficiency, and sharpen their competitive edge.

## **PROBLEM DEFINITION**

While the F1 racing dataset contains a wealth of historical and current information, the data remains largely untapped, and the full potential of data-driven decision-making in the F1 world remains unrealized. The primary challenge is to design and implement a robust data pipeline in Google Cloud Platform (GCP) that efficiently loads data into a data warehouse, enabling users to access and analyze data for strategic and operational insights.

## **DATASET LINK:**

<https://www.kaggle.com/datasets/rohanrao/formula-1-world-championship-1950-2020>

## **Analysis Dimensions:**

The following dimensions can be explored in the F1 racing dataset:

- a. Driver Performance: Analyze individual drivers' performance over time, including wins, pole positions, and podium finishes. Identify key factors that contribute to a driver's success, such as team affiliation, car specifications, and driving styles.
- b. Team Performance: Evaluate the success of F1 teams in terms of constructors' championship points, race wins, and podium finishes. Identify trends in team performance and potential areas of improvement.
- c. Track Analysis: Investigate the impact of different tracks on race outcomes, examining factors such as track length, weather conditions, and average lap times.
- d. Historical Trends: Identify patterns in race outcomes, driver and team performance, and technological advancements over time.
- e. Race Strategy: Analyze pit stop strategies, tire choices, and fuel management to understand the impact on race results.

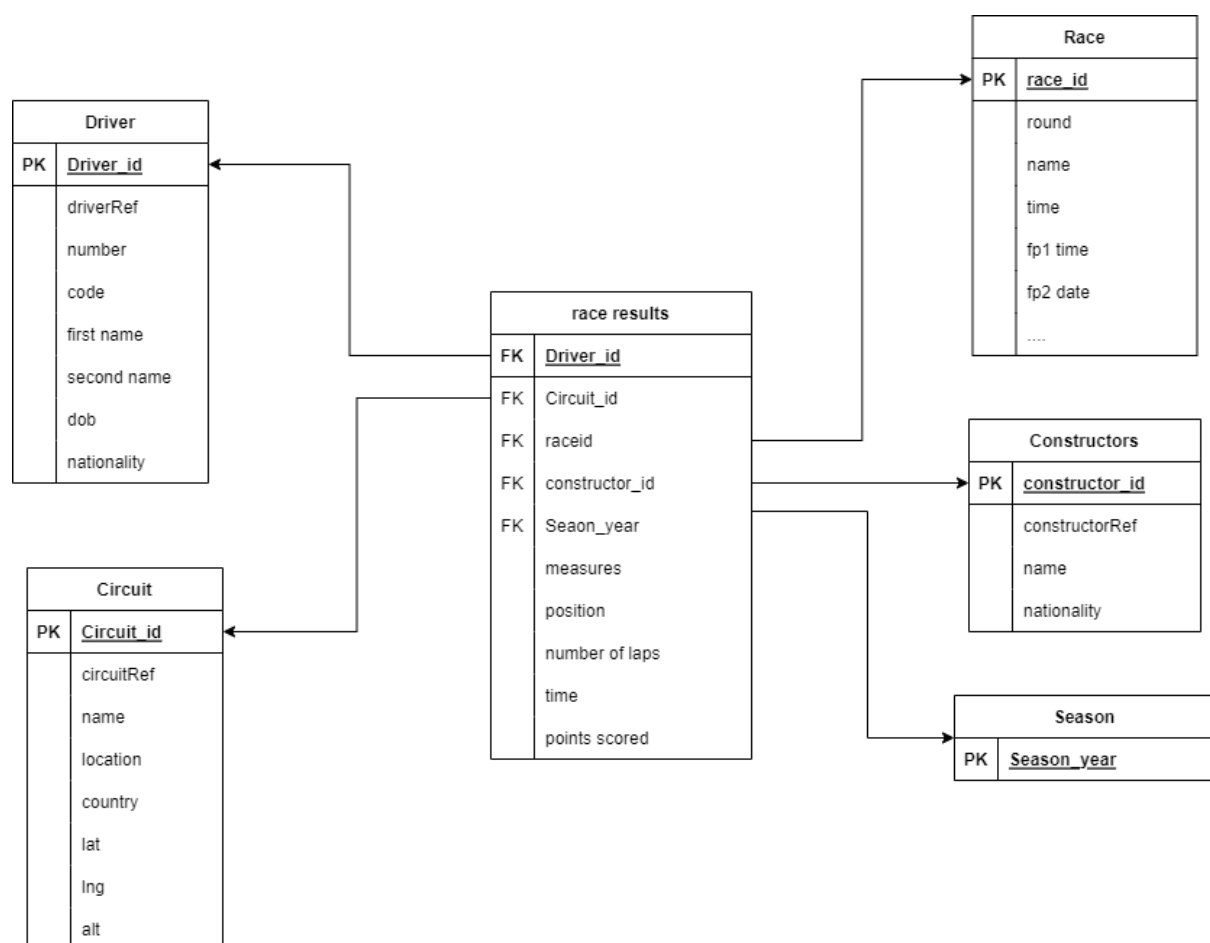
## Strategic and Operational Insights

By building a data pipeline and data warehouse for the F1 racing dataset, several strategic and operational insights can be achieved:

- Identify Success Factors:** Uncover the key factors that contribute to successful drivers and teams, enabling stakeholders to optimize their resources and strategies for future races.
- Optimize Race Strategies:** Analyze historical data to determine the most effective race strategies for specific tracks, weather conditions, and driver-team combinations.
- Benchmarking and Performance Tracking:** Monitor the performance of drivers, teams, and cars against historical data, enabling stakeholders to set ambitious yet realistic goals and track progress over time.
- Enhance Fan Engagement:** Leverage data-driven insights to create engaging and interactive content for fans, deepening their connection to the sport and fostering a more informed and passionate fanbase.

## Datawarehouse Schema

Here is the relational schema for our data warehouse with 5 dimensions and 1 fact table.



## Implementation

Google Cloud Platform (GCP) is a leading suite of cloud computing services offered by Google, providing a range of powerful tools and infrastructure that enable businesses to build, deploy, and scale applications, websites, and services on the same robust foundation as Google itself. GCP is known for its cutting-edge technology, ease of use, and seamless integration with other Google services. Given its robust suite of tools and scalability, GCP is an ideal choice for implementing a data pipeline to load data into a data warehouse.

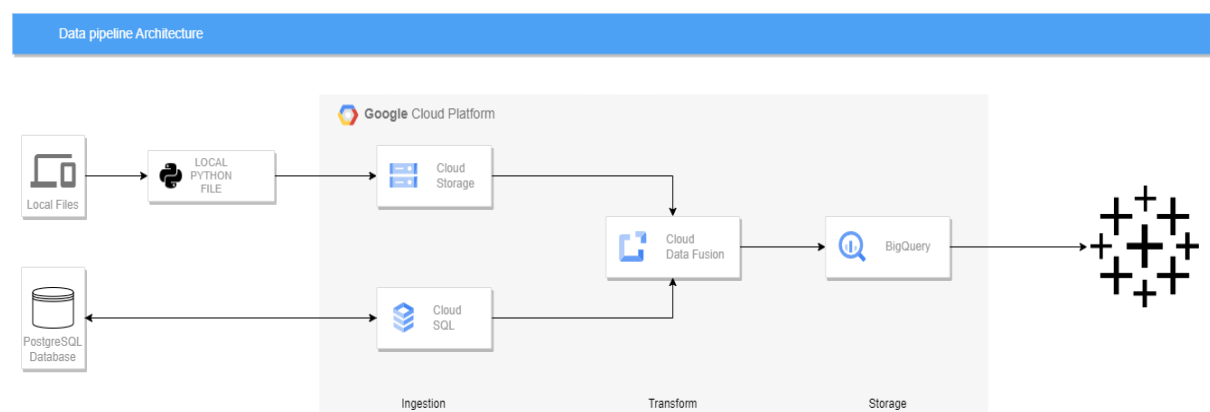
One of the primary reasons GCP is well-suited for this project is its ability to handle large-scale data processing tasks, which is essential given the extensive F1 racing dataset spanning over seven decades. GCP's suite of services, including Cloud Storage, Cloud SQL, Data Fusion, and BigQuery, enables seamless data ingestion, transformation, and storage, thereby streamlining the entire data pipeline process.

Additionally, GCP's serverless architecture ensures that resources are allocated efficiently and automatically scaled based on demand, allowing for cost-effective data processing and storage. This is particularly useful when dealing with large datasets, as it minimizes the need for manual intervention and ensures optimal performance.

Furthermore, GCP's security features and compliance certifications provide an additional layer of protection for the sensitive data being processed and stored, making it a reliable and secure choice for this project.

## Building the Data Pipeline

The data pipeline for the F1 racing dataset was built using several GCP services, following these steps:



### Loading the Dimensions and Fact Data

First, the dimension data was loaded into Google Cloud Storage buckets using a local Python file. This process involved reading the data from the source files, pre-processing it, and uploading it to the Cloud Storage buckets. Cloud Storage provides a scalable, durable, and cost-effective solution for storing and retrieving large amounts of data.

The fact data, which contains race results, was loaded into a Cloud SQL connected PostgreSQL server using pgAdmin. Cloud SQL is a fully-managed relational database service that simplifies database setup, maintenance, and administration. By using Cloud SQL, we can ensure the fact data is efficiently stored, managed, and retrieved for further processing.

### Transforming Data with Google Data Fusion

Google Data Fusion was used to transform the data from the Cloud Storage buckets and the Cloud SQL PostgreSQL server. Data Fusion is a cloud-native data integration service that helps users build and manage ETL (extract, transform, load) pipelines with a visual interface, making it easier to clean, prepare, and blend data from various sources.

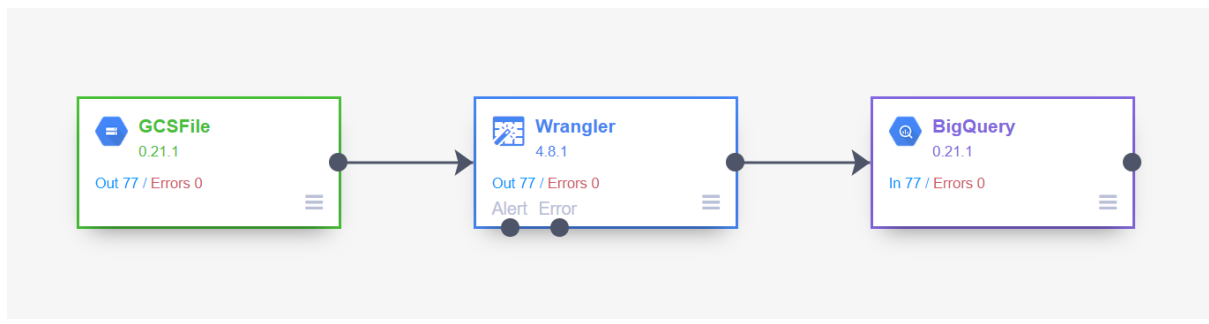
In this step, Data Fusion read the dimension data from Cloud Storage and fact data from Cloud SQL. Custom transformations were applied to clean and enrich the data, ensuring that it was in the appropriate format for further analysis. The transformed data was then prepared for loading into the data warehouse.

### Loading Data into BigQuery Dataset

After the transformation process, the data was loaded into a BigQuery dataset. BigQuery is Google's fully-managed, serverless data warehouse solution designed for large-scale data analytics. It enables super-fast SQL queries using the processing power of Google's infrastructure.

By loading the F1 racing dataset into BigQuery, we can perform complex queries and analysis on the data, uncovering insights into driver performance, team performance, track analysis, historical trends, and race strategies. This enables stakeholders to make data-driven decisions and drive improvements in the F1 racing world.

These is the data pipeline process image for loading the dimensions into the big query table.



The process is similar for the other dimensions as well.

The wrangler has transformation steps in it as show below.

Cloud Data Fusion | Pipeline

Wrangler Properties 4.8.1  
Wrangler - A interactive tool for data cleansing and transformation.

Properties Documentation Metrics

Input Schema

Field	Type	Actions
driverId	int	[-] [0] [1]
driverRef	string	[-] [0] [1]
number	string	[-] [0] [1]
code	string	[-] [0] [1]
forename	string	[-] [0] [1]
surname	string	[-] [0] [1]
dob	string	[-] [0] [1]
nationality	string	[-] [0] [1]
url	string	[-] [0] [1]

Directives

Recipe

- 1 parse-as-simple-date :dob yyyy-MM-dd
- 2 lowercase :driverRef
- 3 lowercase :forename
- 4 lowercase :surname
- 5 lowercase :nationality
- 6 drop :url

WRANGLE

User Defined Directives(UDD)

Value

Output Schema

Field	Type	Actions
driverId	int	[-] [0] [1]
driverRef	string	[-] [0] [1]
number	string	[-] [0] [1]
code	string	[-] [0] [1]
forename	string	[-] [0] [1]
surname	string	[-] [0] [1]
dob	timestamp	[-] [0] [1]
nationality	string	[-] [0] [1]

For loading the race results fact we built the following pipeline which gets the circuit id and season year foreign keys from the race table. The race results in a Cloud SQL instance and the race dimension fact is in Google cloud store so we performed a Join before transforming the data.

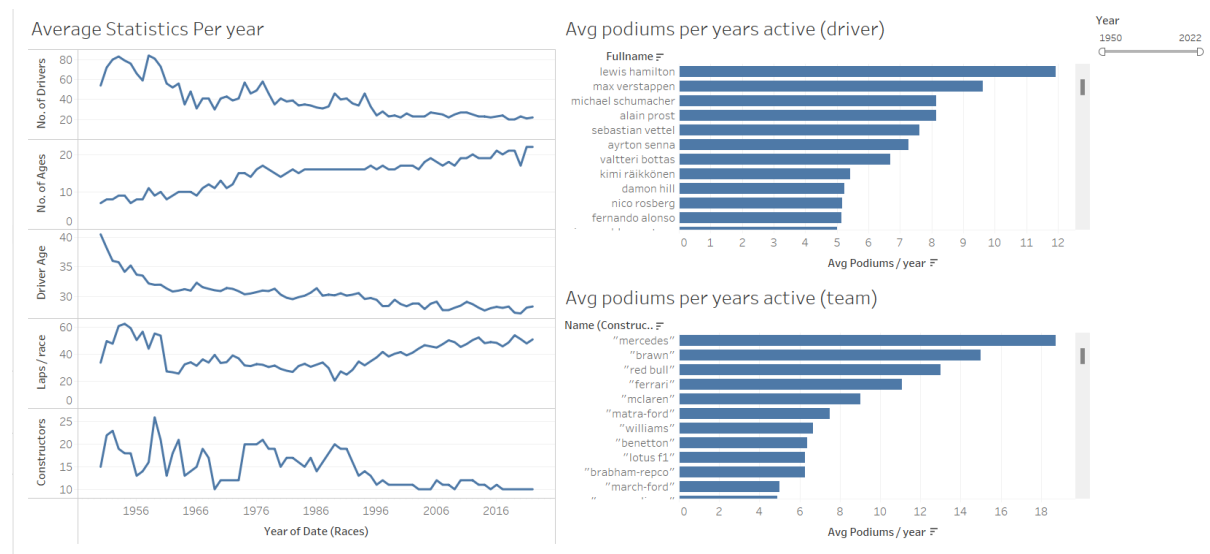


In the wrangler we dropped the redundant columns and pushed the data into big query

The data pipeline built using GCP services allows for efficient ingestion, transformation, and storage of the F1 racing dataset. By leveraging the power of GCP, the data pipeline ensures seamless data processing, scalability, and cost-effective performance, enabling stakeholders to unlock valuable insights and revolutionize the F1 racing world.

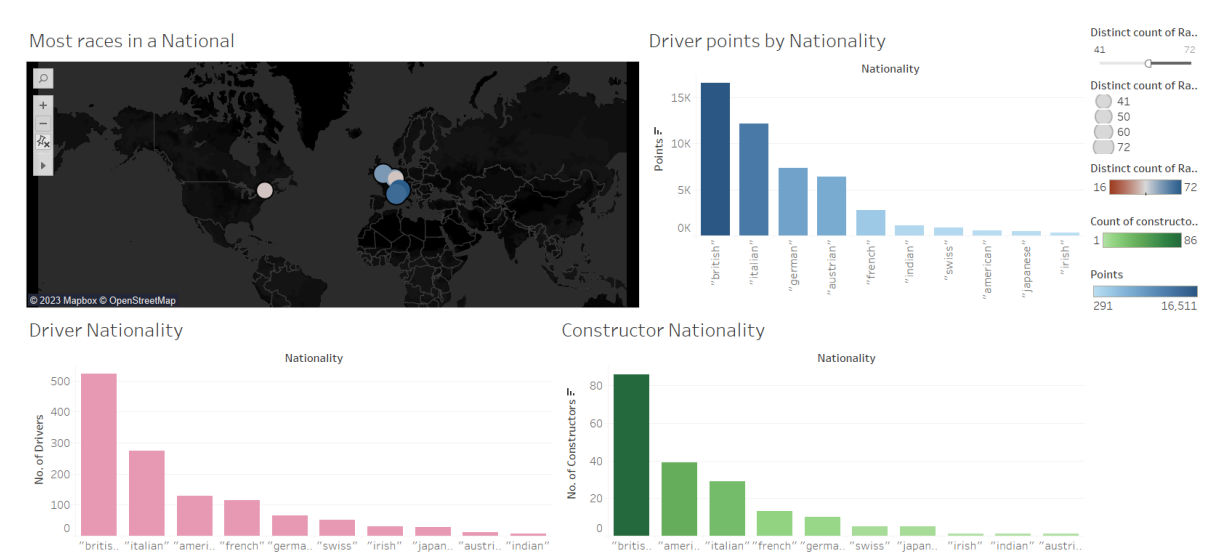
## Analysis

The big query data warehouse was connected to a tableau notebook to analyze the data and gain the insights that we wanted to gain. We first analyzed the basic statistics.



The number of races and the average age of the drivers increases as the years go on but the number of drivers, the number of issues the cars have decreases over time.

Sir Lewis Hamilton and Mercedes remain the most successful constructors and drivers by a huge margin over taking even Ferrari and red bull if we count the number of years they've been active in the sport.



### fastest circuits

Name (Circuits1)	Driver Full Name	Name (Constr..)	Min. Fastestlaptime
albert park grand prix circuit	charles leclerc	"ferrari"	1:20.260
autodromo enzo e dino ferrari	lewis hamilton	"mercedes"	1:15.454
autódromo internacional rodríguez	valtteri bottas	"mercedes"	1:17.774
autódromo internacional de alag.	lewis hamilton	"mercedes"	1:18.750
autódromo internazionale del m...	lewis hamilton	"mercedes"	1:18.833
autódromo josé carlos pace	valtteri bottas	"mercedes"	1:10.540
autodromo nazionale di monza	rubens barrichello	"ferrari"	1:21.046
bahrain international circuit	george russell	"mercedes"	0:55.404
baku city circuit	charles leclerc	"ferrari"	1:43.009
	lance stroll	"williams"	1:45.108
buddhi international circuit	sebastian vettel	"red bull"	1:27.249
circuit de barcelona-catalunya	giancarlo fisichella	"renault"	1:15.641
circuit de monaco	lewis hamilton	"mercedes"	1:12.909
circuit de nevers magny-cours	michael schumacher	"ferrari"	1:15.377
circuit de spa-francorchamps	kimi raikkonen	"mclaren"	1:45.108
circuit gilles villeneuve	valtteri bottas	"mercedes"	1:13.078
circuit of the americas	charles leclerc	"ferrari"	1:36.169
circuit park zandvoort	lewis hamilton	"mercedes"	1:11.097
circuit paul ricard	sebastian vettel	"ferrari"	1:32.740
fuji speedway	felipe massa	"ferrari"	1:18.426
hockenheimring	kimi raikkonen	"mclaren"	1:13.780
hungaroring	lewis hamilton	"mercedes"	1:16.627
indianapolis motor speedway	rubens barrichello	"ferrari"	1:10.399
istanbul park	juan pablo montoya	"mclaren"	1:24.770
jeddah corniche circuit	lewis hamilton	"mercedes"	1:30.734
korean international circuit	sebastian vettel	"red bull"	1:39.605
	jarno trulli	"lotus"	1:43.009
losail international circuit	max verstappen	"red bull"	1:23.196
marina bay street circuit	kevin magnussen	"haas f1 team"	1:41.905
miami international autodrome	max verstappen	"red bull"	1:31.361
nürburgring	max verstappen	"red bull"	1:28.139
red bull ring	carlos sainz	"mclaren"	1:05.619
seang por international circuit	sebastian vettel	"ferrari"	1:34.080

### Successful Team at Track

This treemap visualizes the success of different teams across various circuits. The largest categories are 'mexico city grand prix' and 'indian grand prix', both dominated by the 'ferrari' team (orange). Other significant categories include 'portuguese grand prix' (dominated by 'mclaren', blue), 'são paulo grand prix' (dominated by 'united states', dark blue), and 'bahrain grand prix' (dominated by 'ferrari', orange). The '70th' category shows a mix of teams including 'sakhir grand prix' (orange), 'williams' (red), and 'vanwall' (teal).

### dangerous track

The bar chart displays the number of issues for various circuits. The 'circuit' (brown) has the highest number of issues, exceeding 200. Other circuits with significant issue counts include 'circuit' (orange, ~100), 'autod.' (pink, ~90), 'circuit' (yellow, ~85), 'nurbur.' (dark blue, ~80), 'silverst.' (teal, ~75), 'autodr.' (red, ~70), 'red bull' (green, ~60), 'hocken.' (purple, ~55), and 'hungar.' (dark purple, ~50).

The figure consists of four charts:

- Average issues faced per race:** A line chart showing the average probability of issues from 1951 to 2020. The y-axis is 'Avg prob race' (0 to 4) and the x-axis is 'Date (Races)'. The line fluctuates between approximately 1.5 and 4.0.
- Top reasons to not finish a race:** A horizontal bar chart showing the count of status reasons. The y-axis is 'Count of status' (0K to 2K). The x-axis lists reasons: Engine, Gearbox, Suspension, Transmission, Electrical, Brakes, Withdraw, Clutch, Fuel system, Turbo, Hydraulics, Overheating, Ignition, Oil leak, Throttle, Out of fuel, Harsh start, Wheel, Winch, Oil pressure, Brake pump, Diff, Tyre, Fuel leak, Steering, Collision damage, Radiator, Puncture, Power Unit, Wheel bearing, and Injection. The 'Engine' reason has the highest count, exceeding 2K.
- Drivers with most car issues:** A bubble chart showing the number of issues for various drivers. The size of each bubble represents the number of issues. The bubbles are colored by nationality: american (dark blue), austrian (light blue), brazilian (orange), british (green), finnish (dark green), french (teal), and italian (pink). 'irola' is the largest bubble, indicating the most issues.
- Teams with the most issues with cars (per race):** A bubble chart showing the number of issues per race for various teams. The size of each bubble represents the number of issues per race. The bubbles are colored by constructor: alfa romeo (dark blue), arrows (light blue), benetton (green), brabham (orange), ensign (pink), ferrari (dark blue), force india (green), and jordan (brown). 'ensign' is the largest bubble, indicating the most issues per race.

This is a dashboard is to analyse the average issues per race over the years. Then the reason for most issues being the engine. Then a bubble plot showing the drivers and constructors having the most issues.



## **Conclusion**

In summary, our team has utilized a data warehouse to gain valuable insights into the world of Formula 1 racing. By connecting a Big Query data warehouse to a Tableau notebook, we were able to analyse data from the five dimensions (circuit, driver, constructor, race, status and season) and one fact table (race results) to gain a deeper understanding of the sport.

Through our analysis, we were able to identify trends and patterns in driver and team performance, historical data, race locations, and car issues. We found that the number of races and average driver age increases over time, but the number of drivers and car issues decreases. We also discovered that Sir Lewis Hamilton and Mercedes remain the most successful drivers and constructors in the sport by a significant margin.

By analysing the nationality and nation hierarchies of teams, drivers, and circuit dimensions, we were able to identify the dominance of European natives in the sport. We also found that Italy, Monaco, and Britain are the top locations for races and the nationality from where most constructors are from.

Our analysis of the fastest lap times and drivers and teams with the highest average points at each circuit, as well as the dashboard that analyses the average issues per race over the years, will help teams to optimize their race strategy and identify potential areas for improvement in their car design and performance.

Overall, utilizing a data warehouse and Google Cloud Platform has been crucial in gaining these valuable insights, as it allowed us to easily store, manage, and analyse large amounts of data from various sources. These insights can help teams to make informed decisions and improve their performance in future races.