# Data Engineering | Creation and configuration of an environment for analysis

**Preparing a Data Warehouse with BigQuery for the Autoforce Technical Challenge with *Google Cloud Platform* - GCP.**

## 1. Tools Used in Data Structure Implementation

1. **Kaggle**:

   - For dataset acquisition.

2. **Google Cloud Storage**:

   - Creating the data lake and storing CSV files.

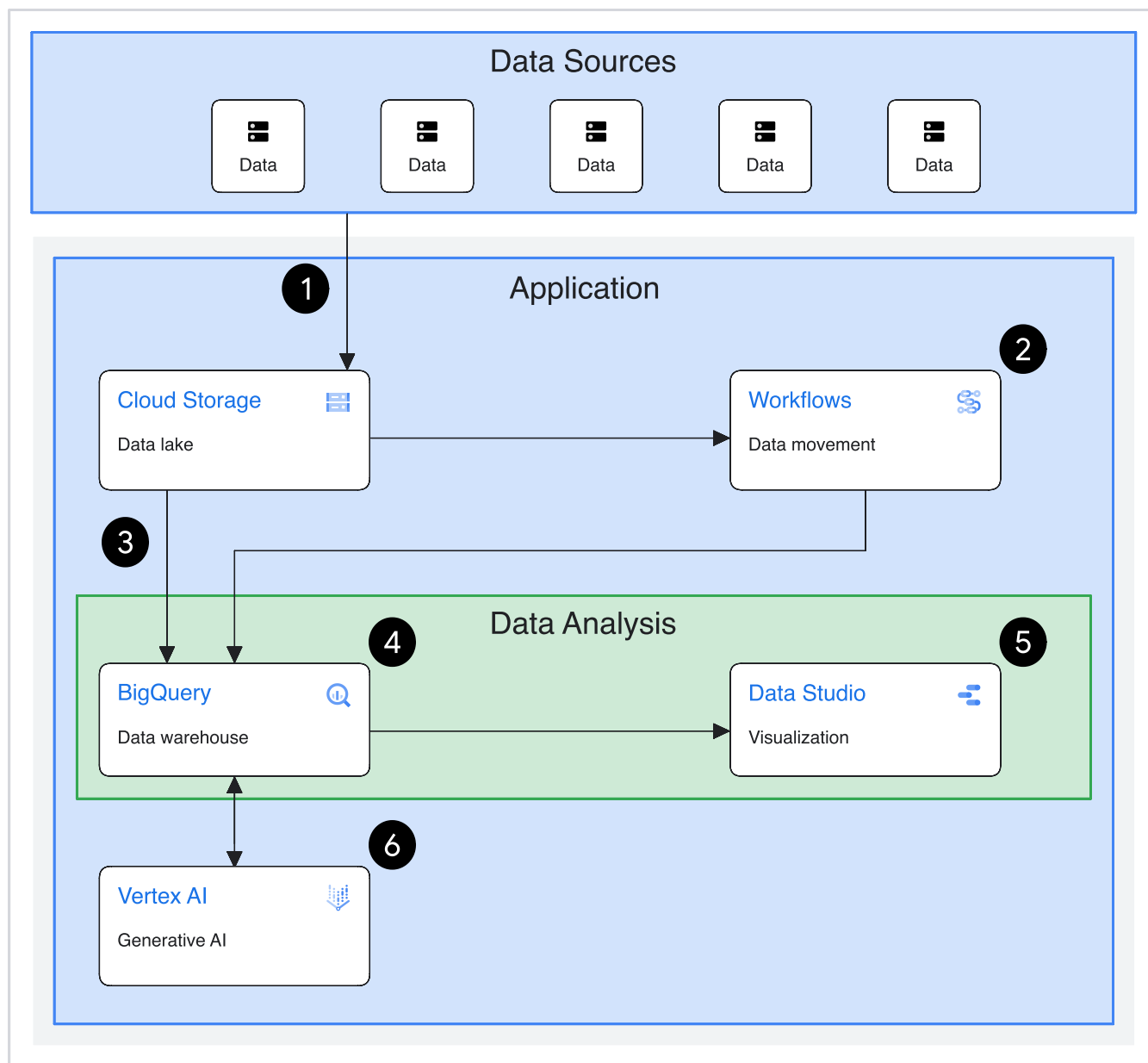3. **BigQuery Studio**:

   - Data warehousing and SQL querying.

4. **SQL**:

   - Language for data manipulation in BigQuery.

5. **Looker Studio**:

   - For data visualization and creating BI dashboards.

## 2. Data Lake Preparation and Analysis Using Google Cloud

**Implementation scheme for the Technical Challenge on GCP.**

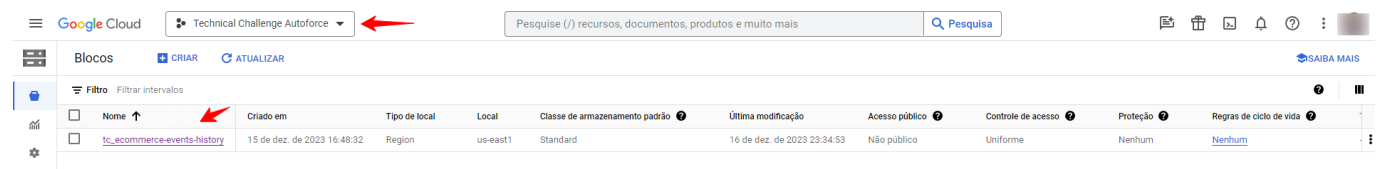**Data Lake Preparation on Google Cloud Storage**

After downloading all the data from Kaggle, a data engineering stage was set up for analysis. The data was then transferred to a **Google Cloud Storage** bucket in Lakehouse, which was created specifically for this challenge. Subsequently, this dataset underwent extensive manipulation and transformation within a **BigQuery Studio** project, aptly named *Technical Challenge Autoforce*.

1. **Data Acquisition**:

   o   Downloaded dataset from Kaggle.
   o   Dataset includes five CSV files with events data from Oct 2019 to Feb 2020.

2. **Google Cloud Storage Bucket Creation**:

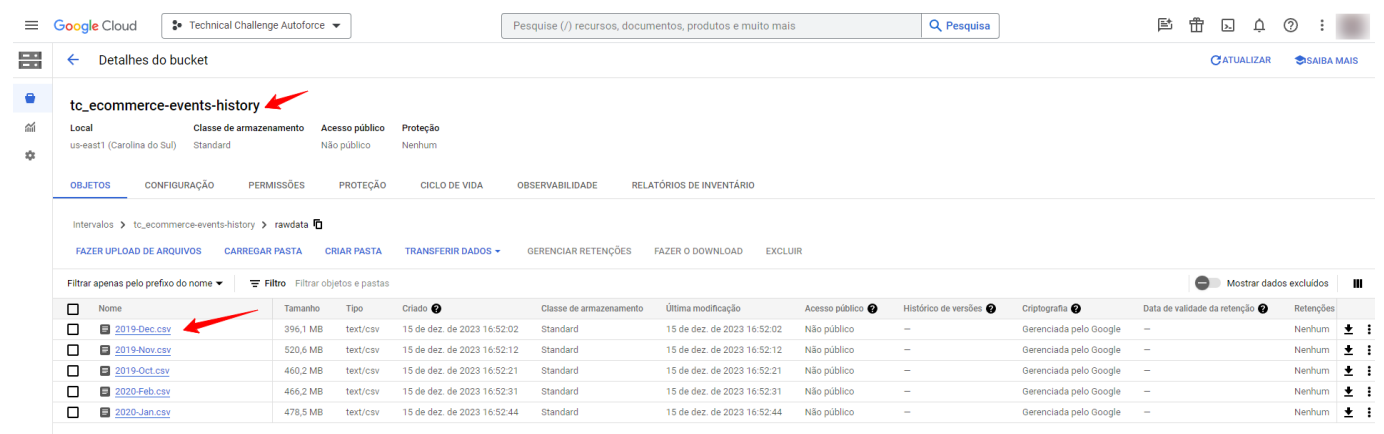   o   Created a bucket `tc_ecommerce-events-history` for data storage.

3. **Data Upload**:
    - Uploaded the CSV files to the `tc_ecommerce-events-history` bucket.



**Data Warehousing and Analysis in BigQuery**

Using **Google Cloud** tools and servers is the best option because it combines the resources of a *data lake* and a *data warehouse* into a single platform for storing, processing, and analyzing structured and unstructured data. This integration facilitates efficient storage, processing, and analysis of diverse datasets, encompassing both structured and unstructured data forms. Moreover, **Google Cloud's** seamless interoperability with **Looker Studio** significantly enhances its utility, enabling the creation of sophisticated, insightful dashboards essential for the visualization and interpretation of critical business intelligence.
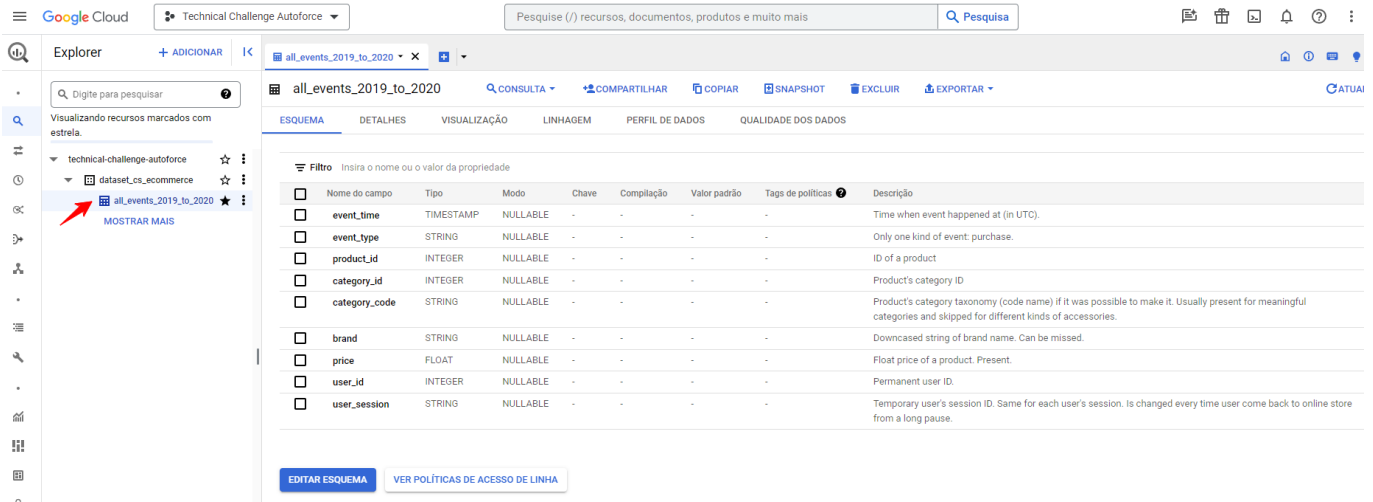
1. **BigQuery Studio Project Setup**:
    - Established the "Technical Challenge Autoforce" project in BigQuery.

2. **Data Import into BigQuery**:
    - Imported CSV files from Google Cloud Storage to BigQuery.

3. **Data Transformation and Unification**:
    - Merged individual tables into a single comprehensive table.

4. **SQL Queries and Analysis**:
   ○ Conducted SQL-based analysis for in-depth insights.



5. **Integration with Colab & Looker Studio**:
   ○ Leveraged Colab to analysis and Looker Studio for data visualization and BI insights.

---

## 3. Second part of the Technical Challenge Analysis - New dataset Makeup API

**Continuing the Case study with a New Challenge to Technical Analysis:**

> "Arbitrarily cross-reference product identifiers with the https://makeup-api.herokuapp.com/ API so that we have data about products and their respective categories."

So, to continue the analysis within **Google's BigQuery** environment, we needed to utilize **Google Cloud Function** to extract the bucket we created (`tc_ecommerce-events-history`). The dataset could then be exported to a table, similar to the process used with the other CSV files. And then later the data in the tables can be crossed.

To create an API with **Google Cloud** services that access the endpoint https://makeup-api.herokuapp.com/api/v1/products.json and downloads the complete dataset in CSV format, follow the steps below:

1. **Activation of required APIs on Google Cloud:**

- **Cloud Functions** to create a serverless function that will download the data.
- **Cloud Storage** if you want to store the generated CSV file.

2. **Creating a function in Cloud Functions:**

- Write a Node.js function supported by **Cloud Functions** that sends a GET request to the makeup API endpoint and converts the JSON response to CSV.

Which was to `main` file:

```javascript
const { Storage } = require('@google-cloud/storage');
const axios = require('axios');
const { Readable } = require('stream');

const storage = new Storage();

exports.downloadMakeupData = async (req, res) => {
  try {
    // Makes request to the makeup API
    const response = await axios.get('https://makeup-
api.herokuapp.com/api/v1/products.json');
    const products = response.data;

    // Convert the response to CSV
    const csvData = convertToCSV(products);

    // Bucket and file name in Cloud Storage
    const bucketName = 'tc_ecommerce-events-history';
    const destinationBlobName = 'source/makeup_data.csv';

    // Bucket and file reference in Cloud Storage
    const bucket = storage.bucket(bucketName);
    const blob = bucket.file(destinationBlobName);

    // Convert CSV string to readable stream
    const csvStream = new Readable();
    csvStream.push(csvData);
    csvStream.push(null);

    // Upload CSV to Cloud Storage
    csvStream.pipe(blob.createWriteStream({
      resumable: false,
      contentType: 'text/csv',
    }))
    .on('error', (error) => {
      res.status(500).send('Error uploading the CSV file: ' + error.message);
    })
    .on('finish', () => {
      res.status(200).send(`Successfully created ${destinationBlobName} in
${bucketName}`);
    });
  } catch (error) {
    res.status(500).send('Error downloading makeup data: ' + error.message);
  }
};
```

```javascript
function convertToCSV(objArray) {
    // Specify the columns you want to include
    const columns = ["id", "brand", "name", "price", "price_sign", "currency",
"rating", "category", "product_type", "tag_list", "created_at", "updated_at"];

    // Create the header with the selected columns
    let str = columns.join(",") + '\r\n';

    // Maps the data to the selected columns and joins with the separator
    return objArray.reduce((accumulator, nextObj) => {
      const row = columns.map(column => {
        let value = nextObj[column] === undefined || nextObj[column] === null ? "" :
nextObj[column].toString();
        if (value.includes(',') || value.includes('"') || value.includes('\n') ||
value.includes('\r')) {
          value = `"${value.replace(/"/g, '""')}"`;
        }
        return value;
      }).join(",");

      return accumulator + row + '\r\n';
    }, str);
}
```

and as dependencies file:

```json
{
  "dependencies": {
    "@google-cloud/functions-framework": "^3.0.0",
    "@google-cloud/storage": "^7.7.0", // Substitua pelo número da versão mais
recente
    "axios": "^1.6.2" // Substitua pelo número da versão mais recente
  }
}
```

- For this code to work, you needed to install the @google-cloud and axios dependencies in your
  **Node.js** project.

In this version of the function in Node.js, the CSV file will be created within the tc_ecommerce-events-history bucket in the source directory with the name makeup_data.csv.

3. **Function deployment in Cloud Functions:**

- The function was implemented using the **Google Cloud** console or the CLI (*gcloud*).

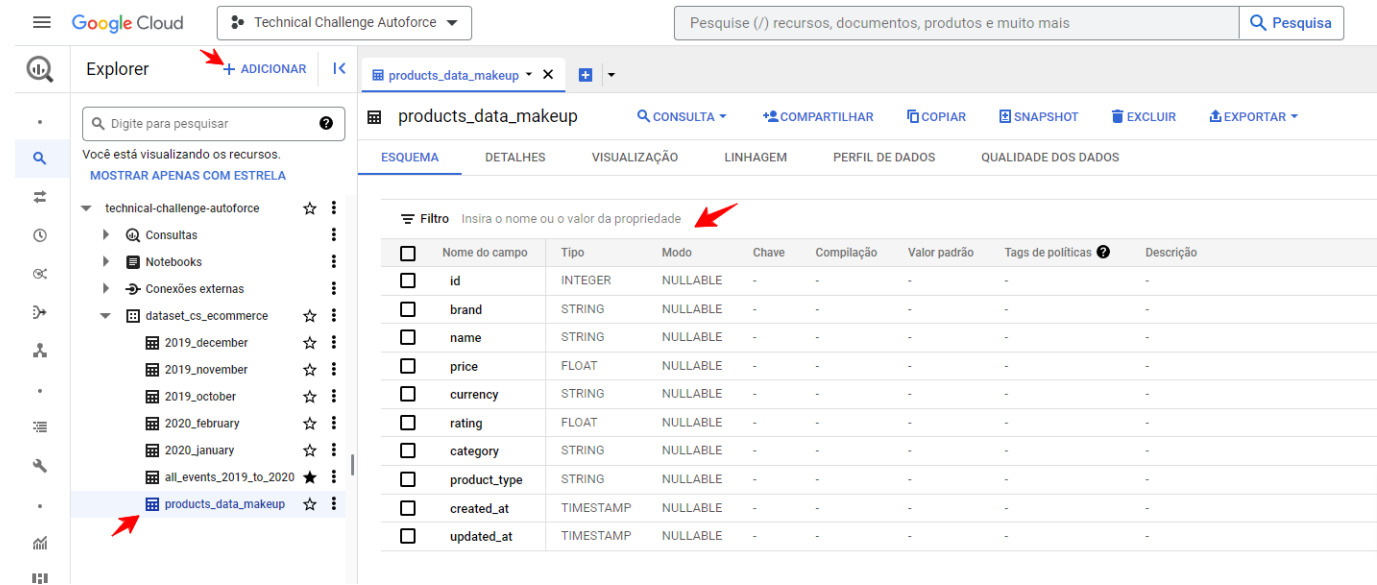4. **Call the created function:**

- A HTTP request was sent to the URL created by **Cloud Functions** in order to execute the function and
  generate the CSV file.

5. **CSV file access:**

- After executing the function, the CSV file was accessible in the designated **Cloud Storage** bucket.

6. **Data Import into BigQuery**

- Imported new makeup CSV file from Google Cloud Storage to BigQuery.



7. **SQL Queries and Analysis on Colab**:
   - Conducted new analysis on product categories as requested on **Google Colab** for in-depth insights.

---

## Benefits of Using Google Cloud Tools

- Robust solution for handling structured and unstructured data.
- Efficient data processing and analysis capabilities.
- Effective visualization and BI insight creation with Looker Studio.