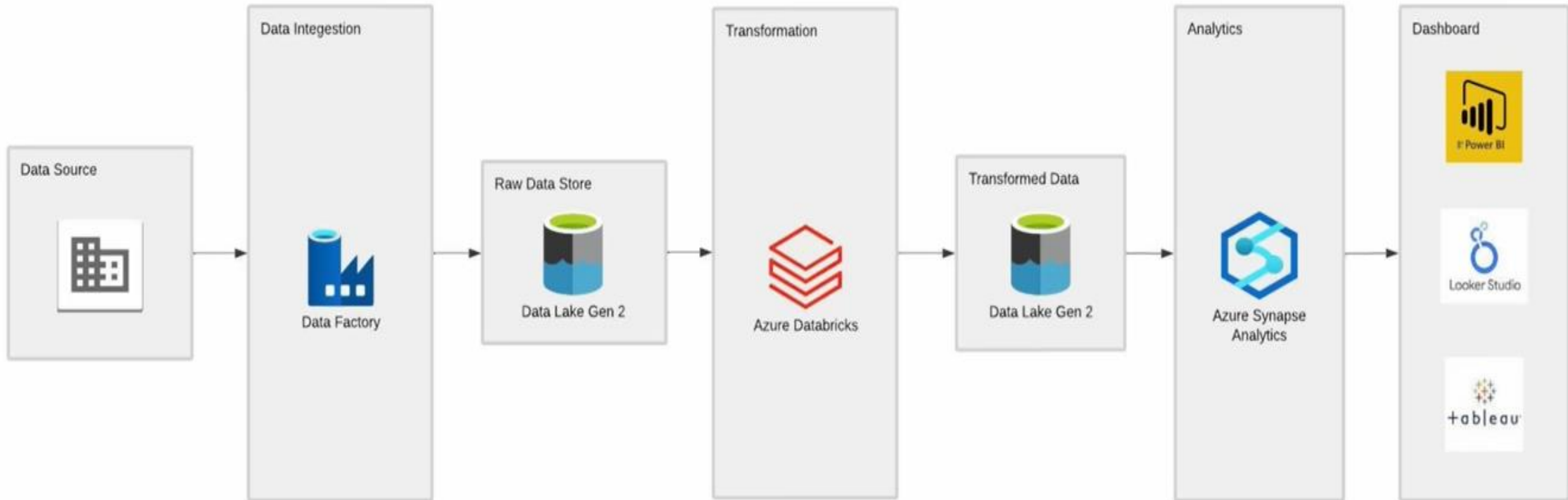# Index

1. **Project architecture**

2. **Step-wise execution details**
   - **Data Integration**
   - **Data Transformation**
   - **Data Modelling**
   - **Data Analytics**
   - **Visualization and Insights generation**

3. **Key Insights**
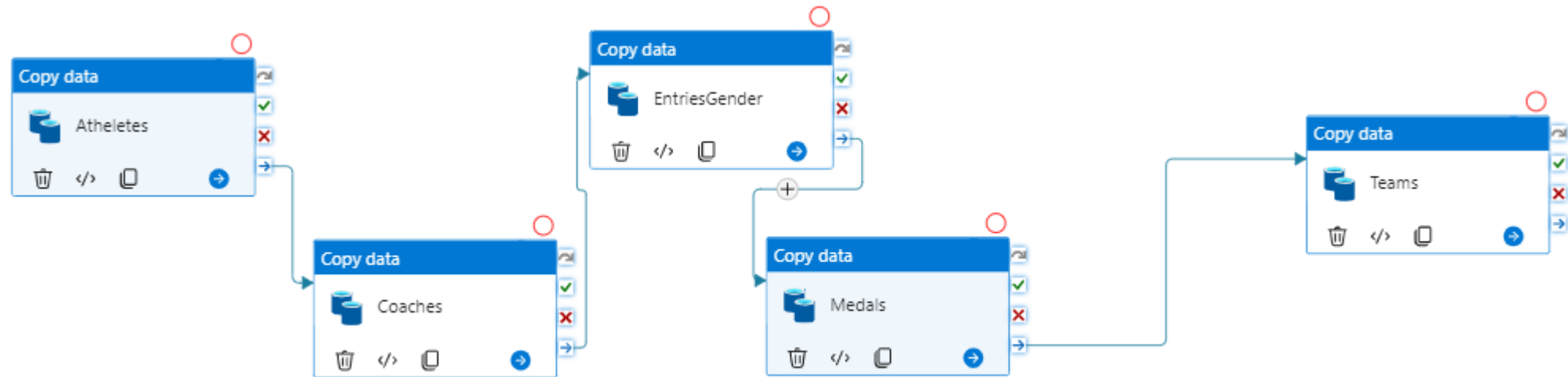
# Data Integration using Azure Data Factory

- In the first stage of the project, I integrated data from various sources using Azure Data Factory. This tool allowed me to create data-driven workflows for orchestrating and automating data movement and data transformation.

# Data Transformation using Azure Databricks

- The next step was to transform the data to make it suitable for analysis. I used Azure Databricks for this purpose, which is an Apache Spark-based analytics platform optimized for Azure. It allowed me to transform large volumes of data and make it ready for analysis.

# Data Storage using Azure Data Lake Gen 2

- After integrating the data, I stored it in a raw data store. For this, I used Azure Data Lake Storage Gen2 which provides secure, scalable, and cost-effective storage.

- I stored the transformed data back into Azure Data Lake Storage Gen2. Storing the transformed data separately ensures that the raw data remains untouched and can be used again if needed.

↑ Upload    🗗 Open in Explorer    🗑 Delete    → Move ⌄    ↻ Refresh    🔲 Open in mobile    🗐 CLI / PS    ⟑ Feedback

ⓘ Enabling SFTP on Azure Blobs has an hourly billing impact. Learn more about pricing.

⌃ Essentials                         JSON View

| | | | |
|---|---|---|---|
| Resource group (move) | : tokyo-olympic-data | Performance | : Standard |
| Location | : centralindia | Replication | : Read-access geo-redundant storage (RA-GRS) |
| Primary/Secondary Location | : Primary: Central India, Secondary: South India | Account kind | : StorageV2 (general purpose v2) |
| Subscription (move) | : Free Trial | Provisioning state | : Succeeded |
| Subscription ID | : d3010d30-a79a-4abc-a224-c7e68379d640 | Created | : 11/6/2023, 3:12:35 PM |
| Disk state | : Primary: Available, Secondary: Available | | |

Tags (edit)      : Add tags

# Data Analytics using Azure Synapse Analytics

- Finally, I analyzed the data and created a dashboard to visualize it. I used Azure Synapse Analytics, an integrated analytics service, for analyzing the data. For visualization, I used Power BI which helped me create comprehensive dashboards and reports.

# Data Visualization using Power BI

- For visualization, I used Power BI which helped me create comprehensive dashboards and reports.



## Tokyo Olympics 2020

| Total no. of players | Total no. of countries | Total no. of disciplines |
|---|---|---|
| 11.06K | 206 | 46 |

### Country-wise medals

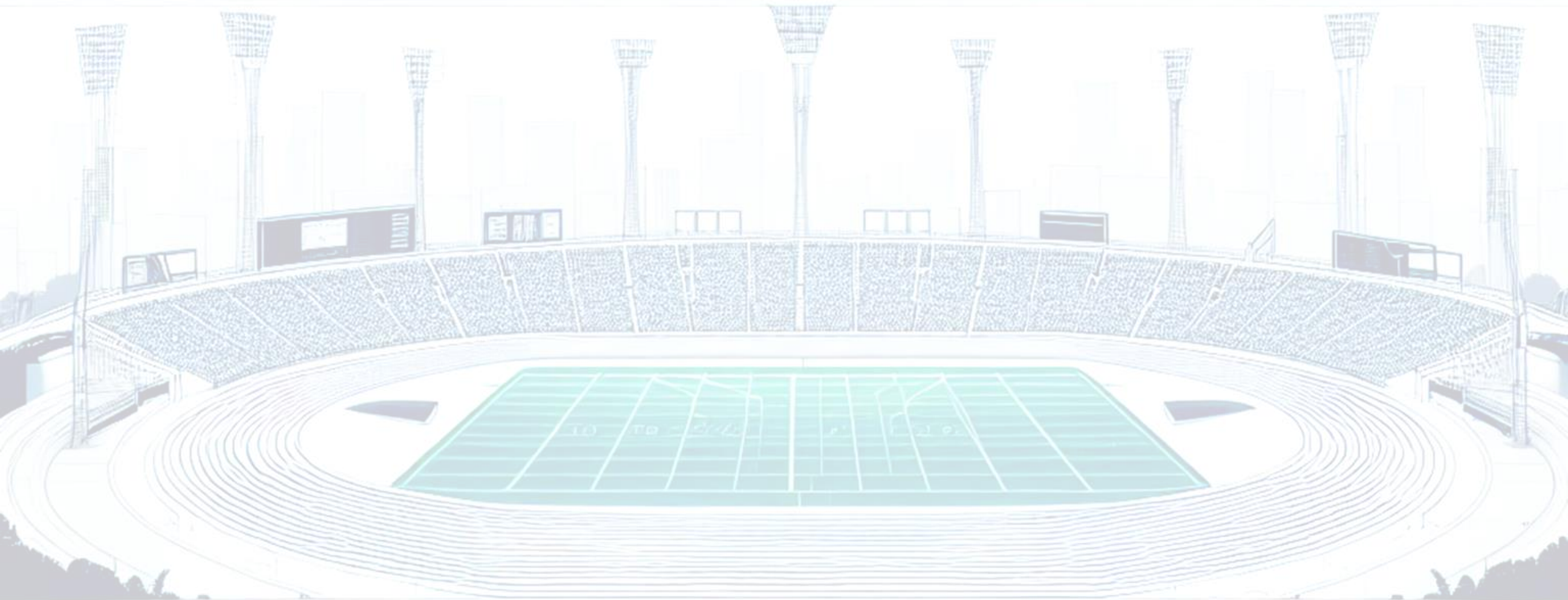| Country | Gold medals | Silver medals | Bronze medals | Total medals | Rank by Total medals |
|---|---|---|---|---|---|
| United States of America | 39 | 41 | 33 | 113 | 1 |
| People's Republic of China | 38 | 32 | 18 | 88 | 2 |
| ROC | 20 | 28 | 23 | 71 | 3 |
| Great Britain | 22 | 21 | 22 | 65 | 4 |
| Japan | 27 | 14 | 17 | 58 | 5 |
| Australia | 17 | 7 | 22 | 46 | 6 |
| Italy | 10 | 10 | 20 | 40 | 7 |
| Germany | 10 | 11 | 16 | 37 | 8 |
| Netherlands | 10 | 12 | 14 | 36 | 9 |

# Key Insights generated from Power BI report

- More than 10 thousand players participated in Tokyo Olympics 2020, the exact number of participants were 11,060 plyers participated from various countries.

- More than 200 countries participated in Tokyo Olympics 2020.

- The total number of disciplines in Tokyo Olympics were 46, some of the important disciplines were:
    - Athletics
    - Swimming
    - Badminton
    - Table Tennis.

- USA topped the medal tally with 113 total medals followed by China with 88 medals followed by Great Britain with 71 total medals.

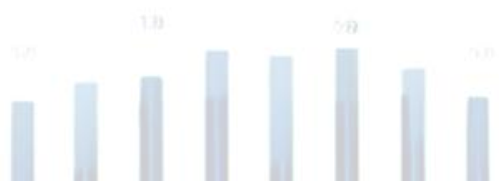- India won total 7 medals and ranked 33 among all the countries

# Thank You