



# BUILD 2020 OLYMPIC GAMES DATA PROCESSING PIPELINE ON AZURE



**Mentor/Intern:** Nguyen Ngoc Thien / Tran Quoc Hai  
Le Trung Viet / Tran Thi Anh Thu  
Phan Thien Huu

# TABLE OF CONTENT



01

**Big Data  
Format**

02

**Pipeline & Demo**

03

**Visualization**

04

**Cost  
Management**



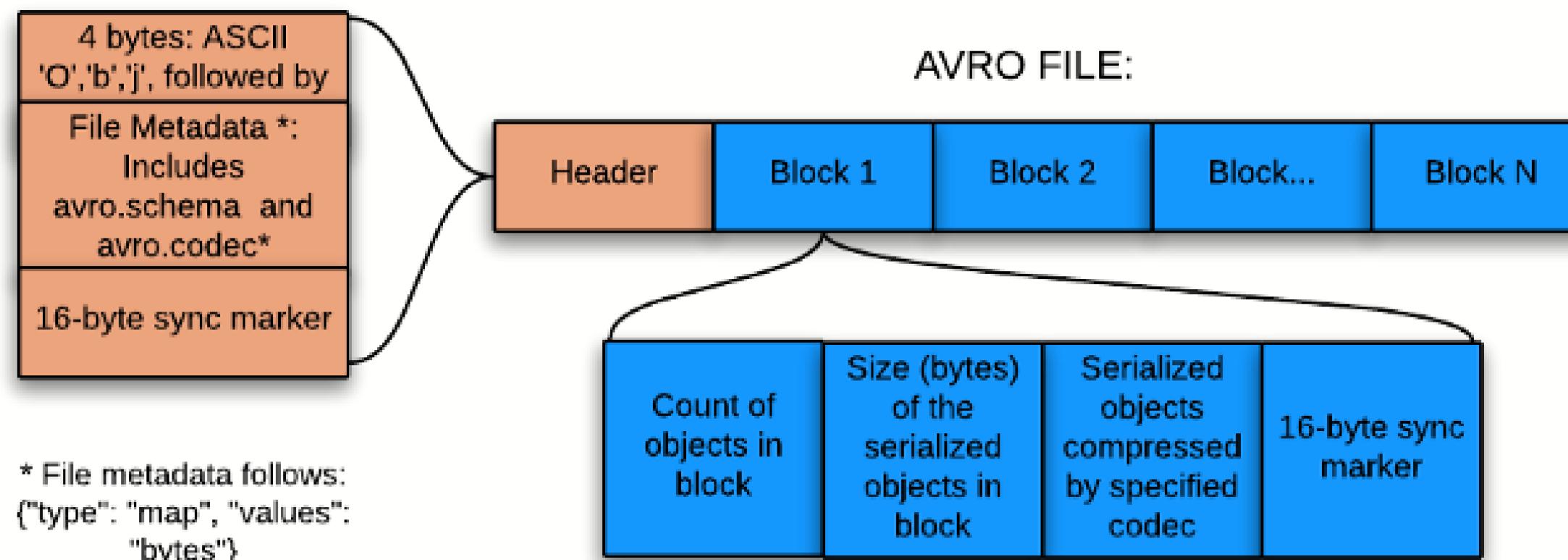
# BIG DATA FORMAT



# AVRO

Apache Avro is a data serialization format.

Apache Avro is a row-based storage.



# AVRO

## Flexibility

allows adding,  
editing, deleting

## Compact size

stored as  
compressed binary

## High speed access

provides efficient  
and fast write-to-  
read

## Multi-language support

Java, python, c++,  
and many more  
languages

## Easy integration

Hadoop, Spark,  
Kafka

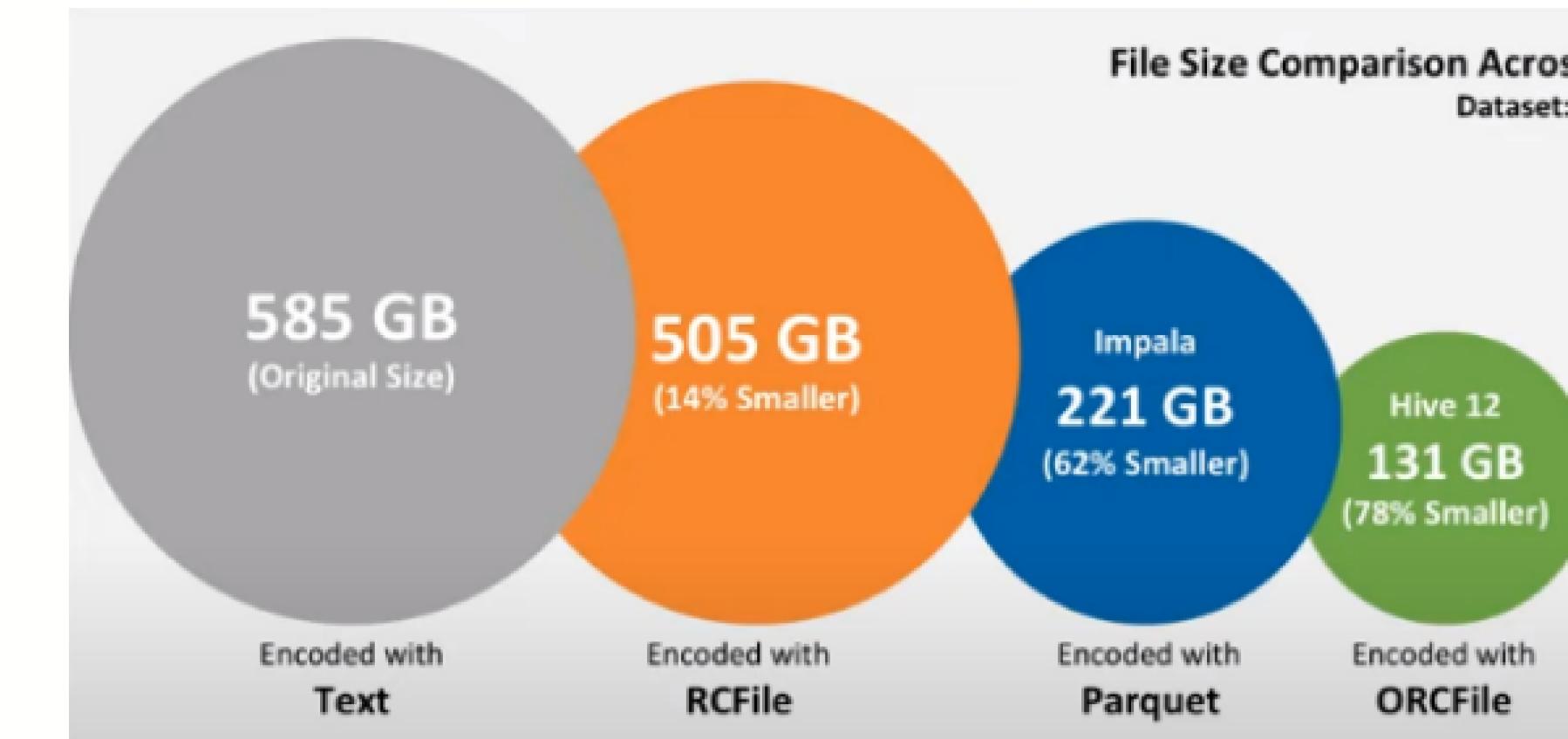
- • • • •
- • • • •
- • • • •
- • • • •

# ORC

- Column-based

Tên	Địa chỉ	SĐT	Email	Ngày sinh	Tổng tiền đã mua
John Smith	123 Main st	12345	john@...	1985-07-10	100000
Jane Doe	456 Oak St	6789	jane@...	1990-07-02	2000000

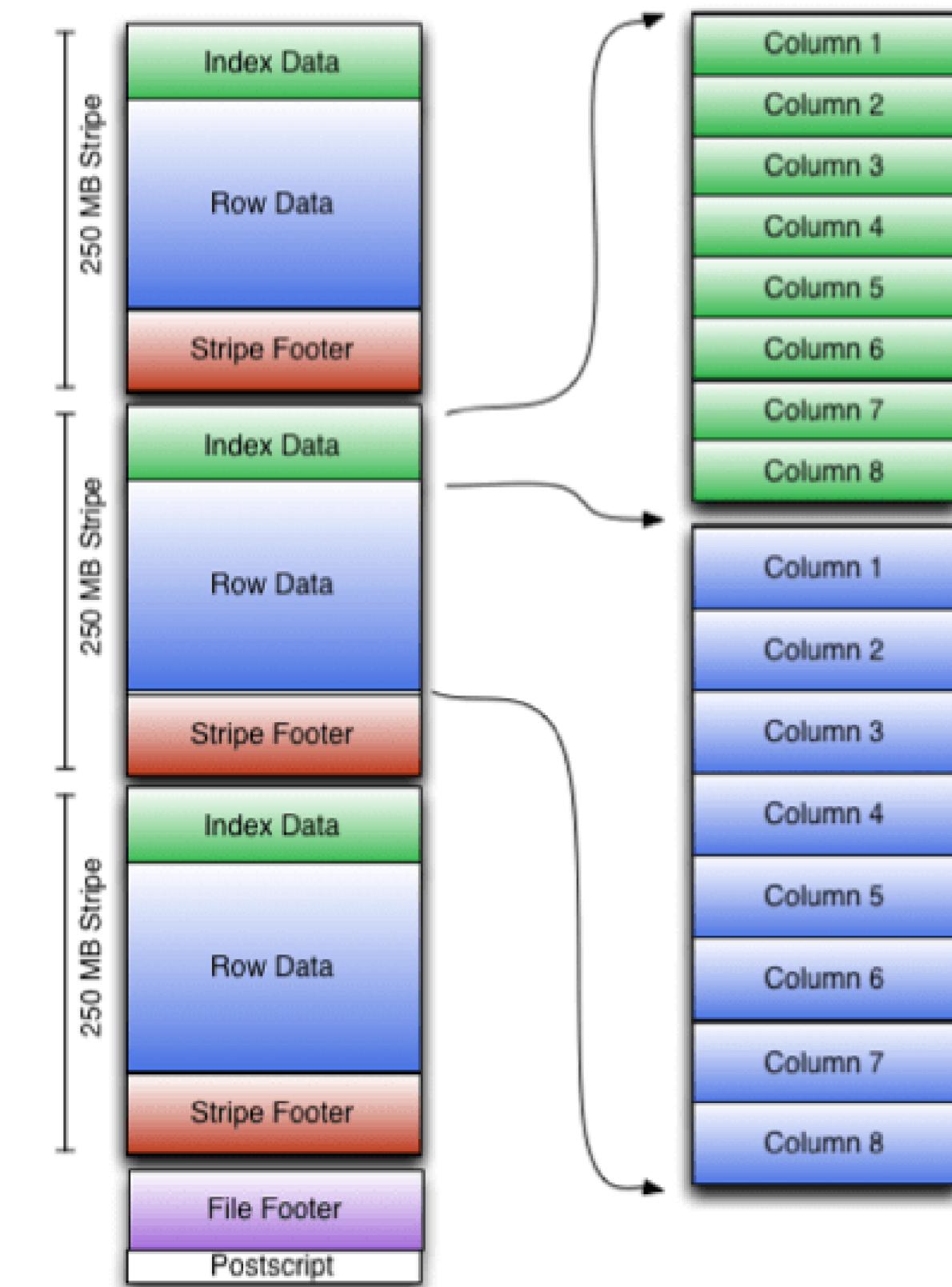
- Compression



# ORC

## Purpose:

- Speed Up Query
- Save Storage Space
- Efficiency in Big Data Processing
- Support Read



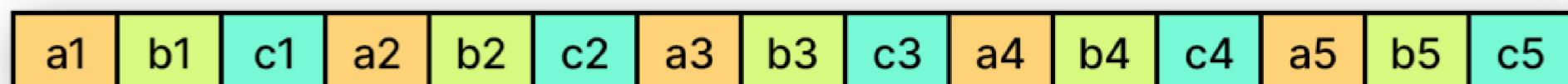
# PARQUET

Apache Parquet is a column-oriented data file format

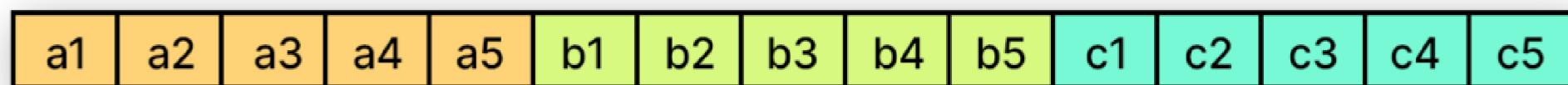
Logical table representation

a	b	c
a1	b1	c1
a2	b2	c2
a3	b3	c3
a4	b4	c4
a5	b5	c5

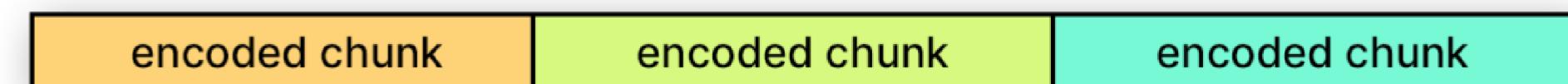
Row Layout



Column Layout



↓                    ↓                    ↓ encoding



# PARQUET

Free & Open source  
file

Language Agnostic

Column-based  
format

Highly efficient

Supports complex  
data types

# BENEFITS

Stores big data of any kind

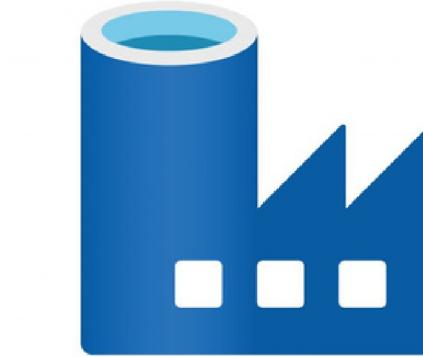
Saves on cloud storage space

Increases data throughput & performance

# COMPARISON



	AVRO	PARQUET	ORC
Schema Evolution	Best	Good	Better
Compression	Good	Better	Best
Row/Column Oriented	Row	Column	Column
Read/Write	Write	Read	Read
Splittability	Good	Good	Best

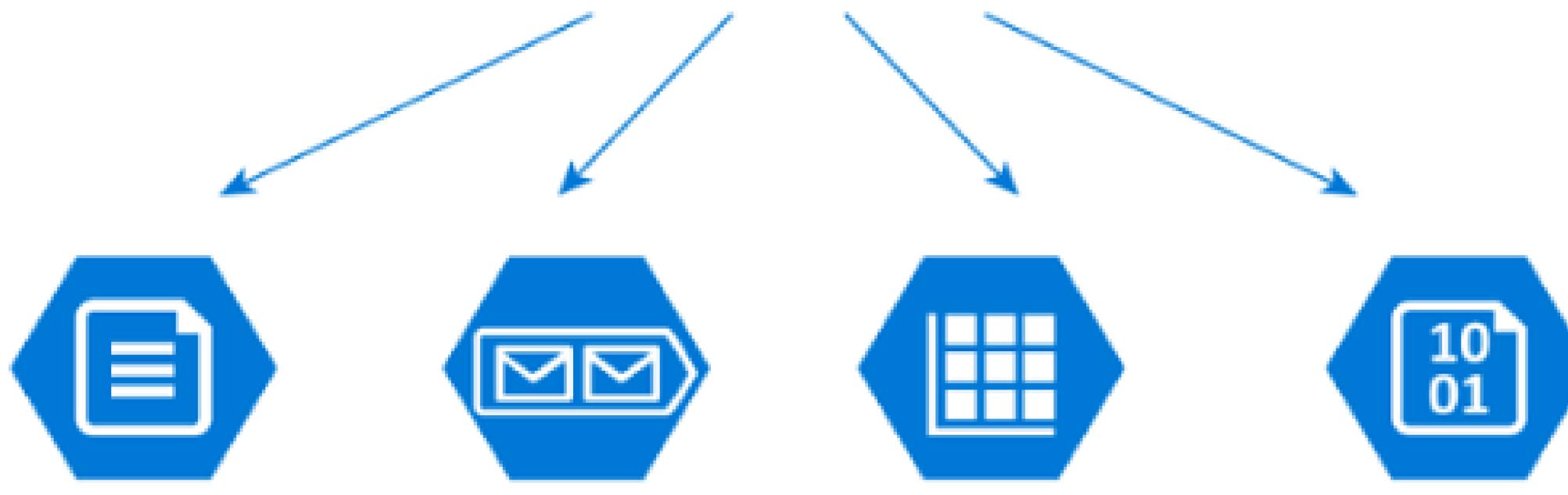


# PIPELINE & DEMO



# STORAGE ACCOUNT

## General-Purpose Storage Account



File

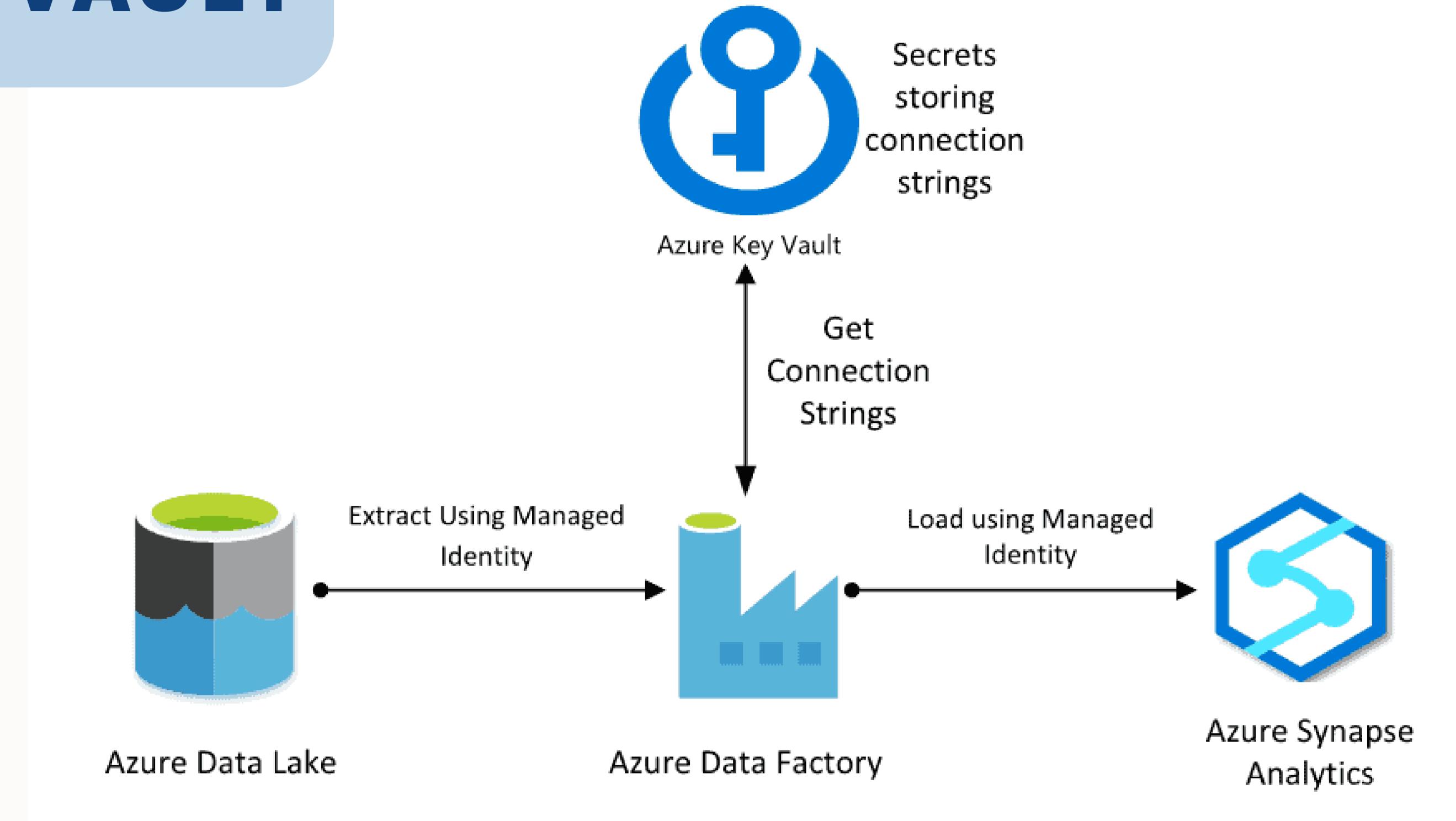
Queue

Table

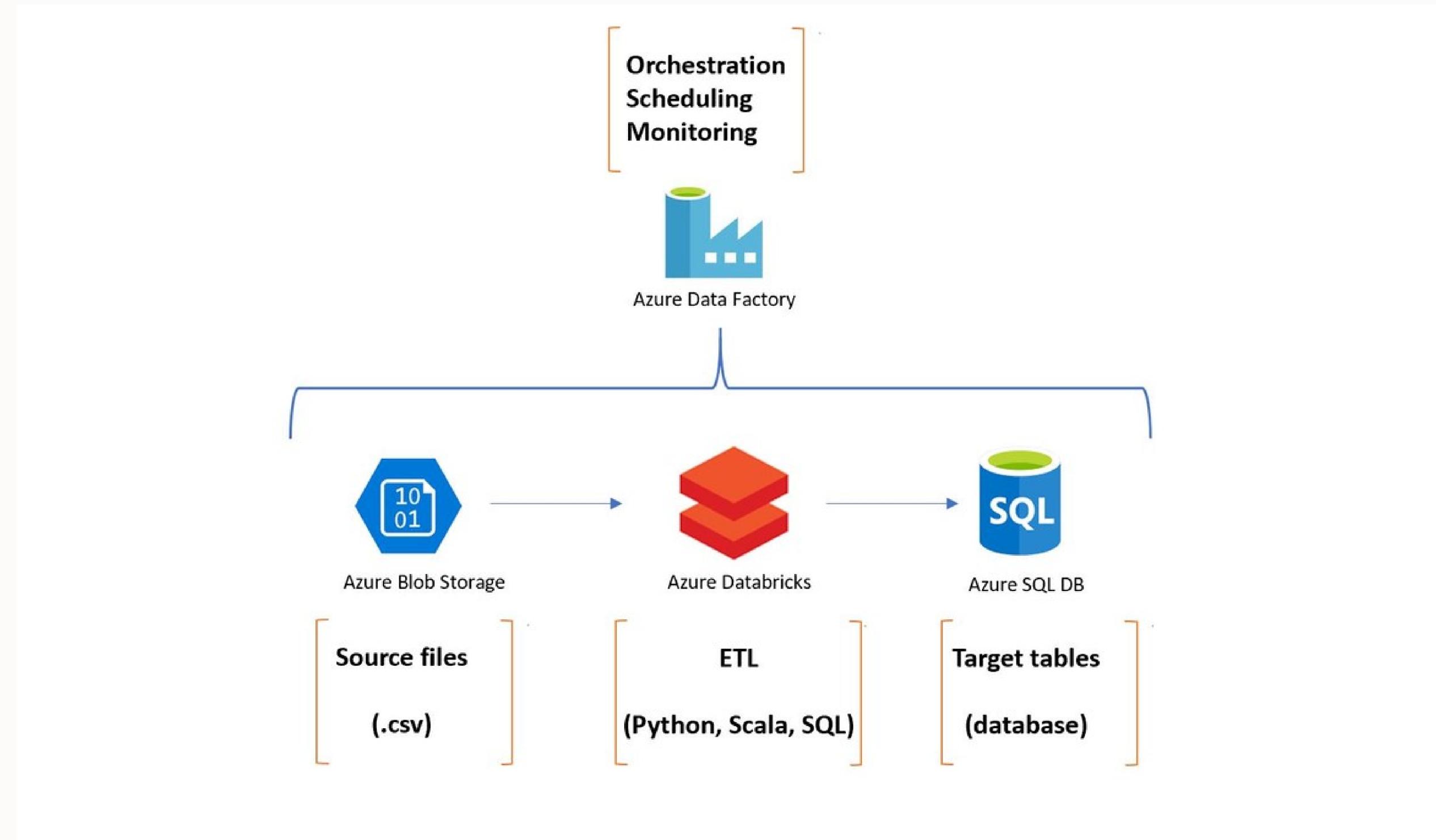
Blob

Storage Types

# KEY VAULT

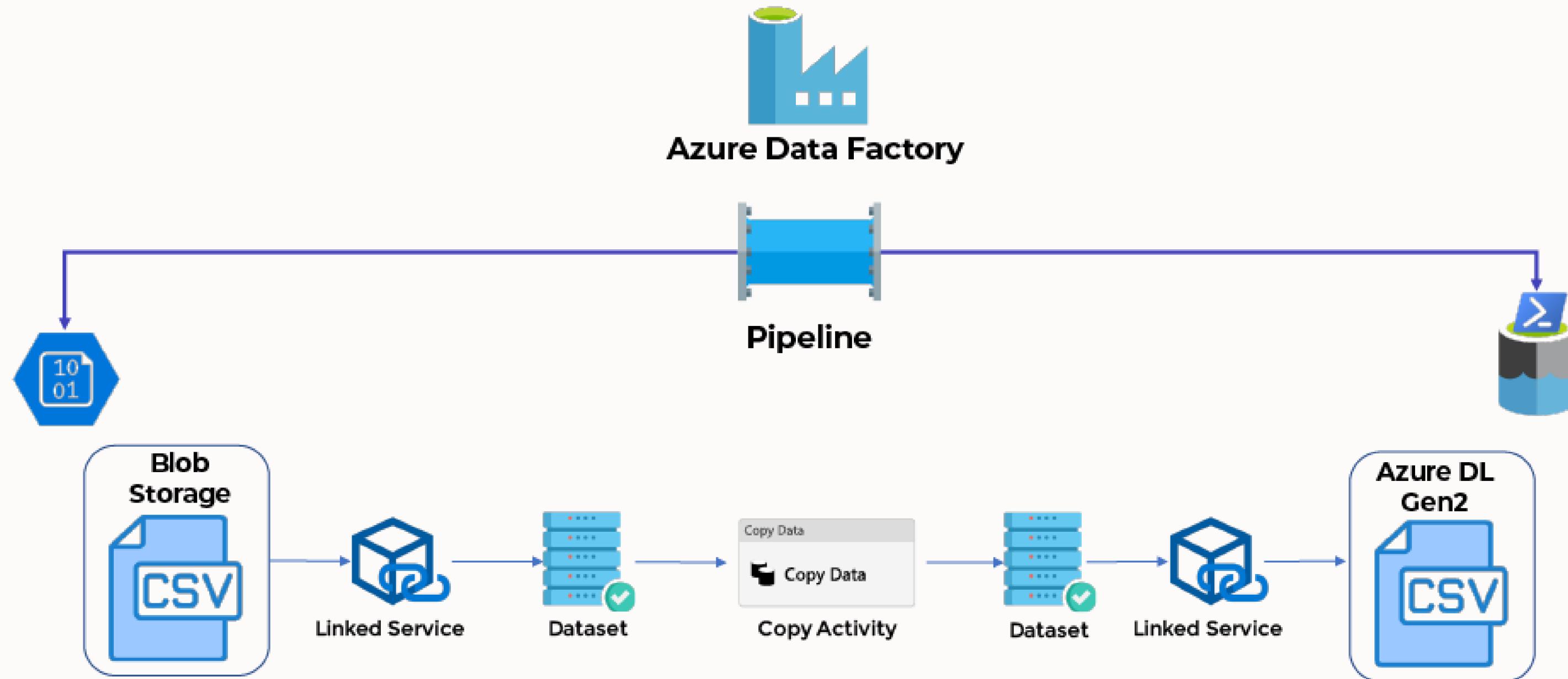


# DATABRICKS - SQL DATABASE



# DATA FACTORY

## ADF Architecture



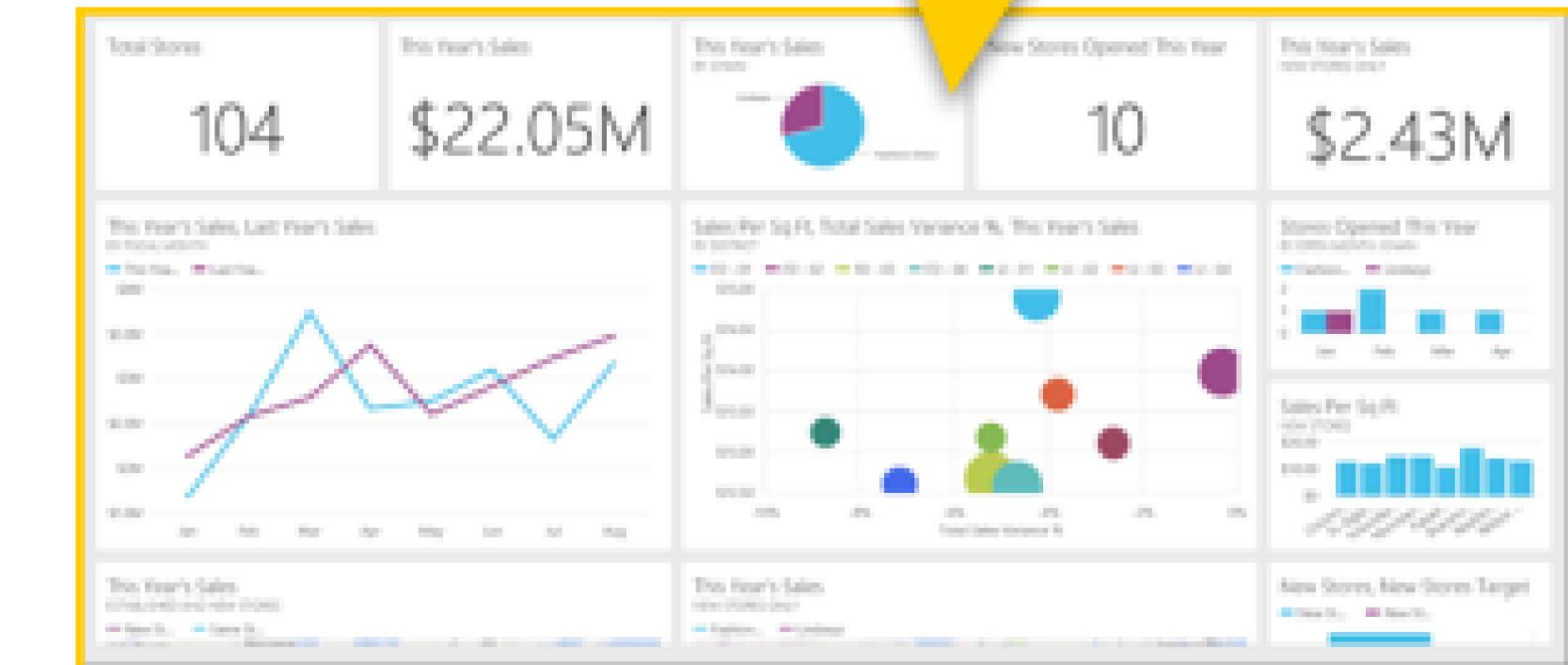
# POWER BI



**SQL**  
Database



**Power BI**



# DATA

## Tokyo 2020 Olympic Summer Games



Athletes



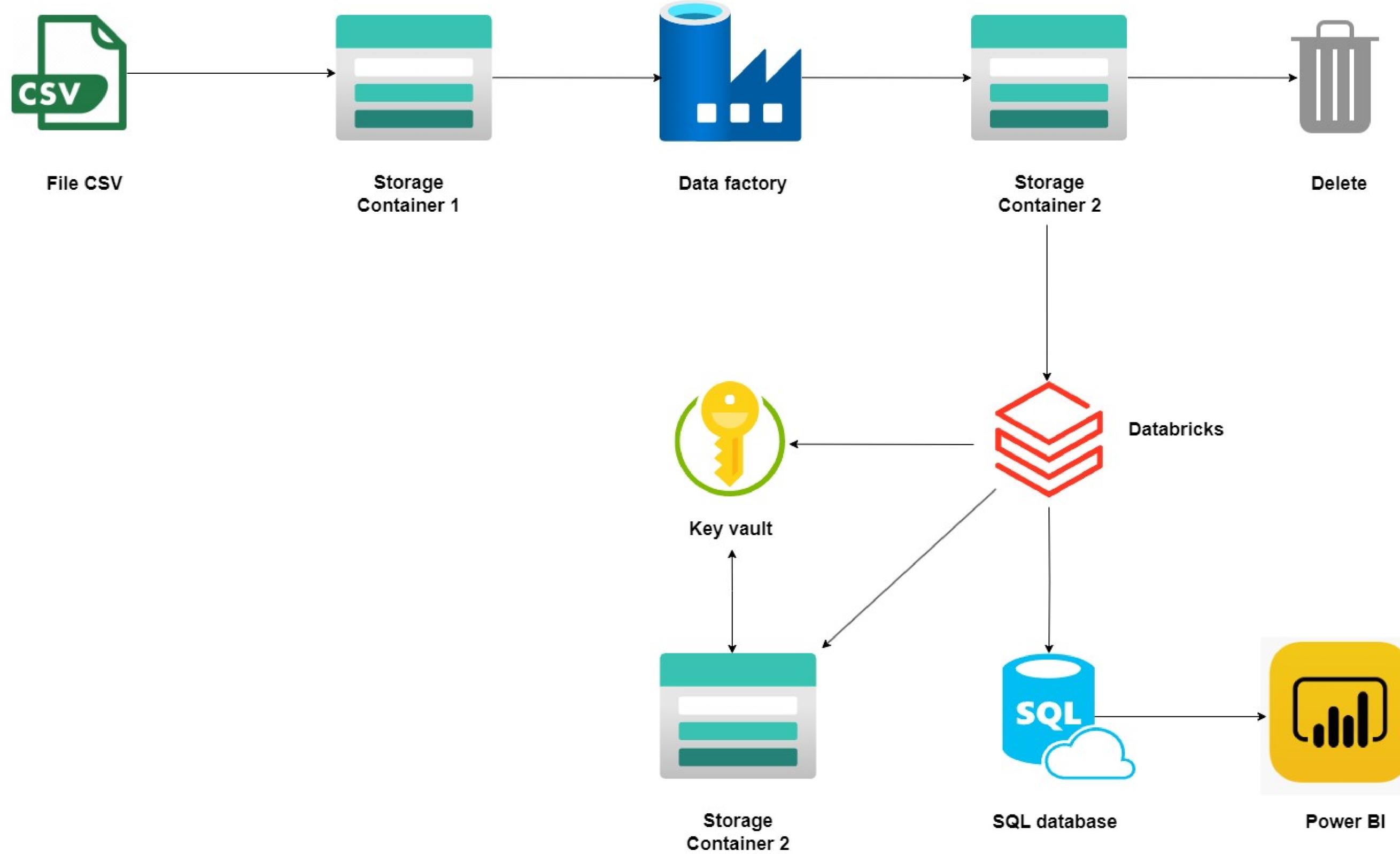
Coaches



Medals



# Azure Pipelines



# Read data by access key



Test try data Python

File Edit View Run Help Last edit was 23 minutes ago Provide feedback

Run all 22/8 Schedule Share

```
1 container = "container1"
2 storage_account_name = "huuphan"
3 spark.conf.set(
4     "fs.azure.account.key."+storage_account_name+".blob.core.windows.net",
5     dbutils.secrets.get("keyofme", "top"))
6 file_location = f"wasbs://{container}@{storage_account_name}.blob.core.windows.net/medals2.csv"
7 file_type = "csv"
8 df = spark.read.format(file_type).option("header", "true").load(file_location)
9 display(df)
```

▶ (2) Spark Jobs

▶ df: pyspark.sql.dataframe.DataFrame = [c0: string, medal\_type: string ... 11 more fields]

Table +

	_c0	medal_type	medal_code	medal_date	athlete_short_name	athlete_name	athlete_sex
1	0	Gold Medal	1	2021-07-24 00:00:00.0	KIM JD	KIM Je Deok	null
2	1	Gold Medal	1	2021-07-24 00:00:00.0	AN S	AN San	X
3	2	Silver Medal	2	2021-07-24 00:00:00.0	SCHLOESSER G	SCHLOESSER Gabriela	X
4	3	Silver Medal	2	2021-07-24 00:00:00.0	WIJLER S	WIJLER Steve	null
5	4	Bronze Medal	3	2021-07-24 00:00:00.0	ALVAREZ L	ALVAREZ Luis	X
6	5	Bronze Medal	3	2021-07-24 00:00:00.0	VALENCIA A	VALENCIA Alejandra	null
7	6	Gold Medal	1	2021-07-24 00:00:00.0	CARAPAZ R	CARAPAZ Richard	M

## Delete unnecessary columns

```
1 df1 = df.drop('_c0', 'medal_date')
```

▶ df1: pyspark.sql.dataframe.DataFrame = [medal\_type: string, medal\_code: string ... 9 more fields]  
 Command took 0.08 seconds -- by ph.thienhuu@gmail.com at 22:29:54 22/8/2023 on 22/8

## handle null value and drop columns

```

1 from pyspark.sql.functions import *
2 from pyspark.sql.window import Window
3 # Xác định giá trị khác null trong cột theo từng nhóm
4 wd = Window.partitionBy('event')
5 df2 = df1.withColumn('NonNullValues', last('athlete_sex', ignorenulls=True).over(wd)) #lấy giá trị cuối cùng khi sắp xếp
tăng dần bỏ qua các giá trị null
6
7 # Thay các giá trị null bằng giá trị khác null đã tìm ra
8 df2 = df2.withColumn('athlete_sex', col('NonNullValues'))
9 df2 = df2.drop('NonNullValues')
10
11 # thay các giá trị null còn lại bằng ý nghĩa của dữ liệu
12 df3 = df2.withColumn('athlete_sex', when((col('athlete_sex').isNull()) & (col('event').like('%Women%')), 'W').otherwise(col('athlete_sex')))
13 df4 = df3.withColumn('athlete_sex', when((col('athlete_sex').isNull()) & (col('event').like('%Men%')), 'M').otherwise(col('athlete_sex')))
14
15 # từ ý dữ liệu của cột athlete_sex(thể hiện giới tính của nội dung thi đấu) chuyển qua event_sex
16 df4 = df4.withColumnRenamed('athlete_sex', 'event_sex')
17 medal_df = df4.withColumnRenamed('athlete_link', 'url')
18 medal_df = medal_df.drop('country_code', 'country', 'discipline_code', 'discipline')
19 display(medal_df)

```

# read and handle null value in file Athletes

Read Athletes CSV and transform

```

1 #Đọc địa chỉ file athletes
2 file_location1 = f"wasbs://{container}@{storage_account_name}.blob.core.windows.net/athletes.csv"
3 file_type = "csv"
4 athlete_df = spark.read.format(file_type).option("header", "true").load(file_location1)
5 display(athlete_df)

```

```

1
2     from pyspark.sql.functions import *
3
4     # thay thế các giá trị null trong bảng
5     athlete_df_1 = athlete_df.withColumn('birth_place', when((col('birth_place').isNull()),(col('birth_country'))).otherwise(
6         (col('birth_place'))))
7     athlete_df_1 = athlete_df_1.withColumn('birth_place', when((col('birth_place').isNull()),(col('country'))).otherwise(col(
8         'birth_place')))
9     athlete_df_1 = athlete_df_1.withColumn('birth_country', when((col('birth_country').isNull()),(col('country'))).otherwise(
10        (col('birth_country'))))
11    athlete_df_1 = athlete_df_1.withColumn('residence_country', when((col('residence_country').isNull()),(col(
12        'residence_place'))).otherwise(col('residence_country')))
13    athlete_df_1 = athlete_df_1.withColumn('residence_country', when((col('residence_country').isNull()),(col('country'))).
14        otherwise(col('residence_country')))
15    athlete_df_1 = athlete_df_1.withColumn('residence_place', when((col('residence_place').isNull()),(col('country'))).
16        otherwise(col('residence_place')))
17    athlete_df_1 = athlete_df_1.withColumn('height_m/ft', when(col('height_m/ft').isNull(),'not collected yet').otherwise(col(
18        'height_m/ft'))))
19    display(athlete_df_1)

```

# create dim\_country and Id for athlete table

```

1 # sau khi thực hiện làm sạch
2 # tạo country_df
3 country_df = athlete_df_1.select('country_code', 'country')
4 country_df = country_df.distinct()
5 display(country_df)

```

```

1 #kiểm tra trước khi join
2 if athlete_df_1.select('url').distinct().count() == athlete_df.select('url').count() :
3     print('không có dữ liệu trùng lặp')
4 else :
5     print('Có dữ liệu trùng lặp')

```

► (5) Spark Jobs

không có dữ liệu trùng lặp

```

1 # tạo id cho bảng athlete
2 athlete_df_2 = athlete_df_1.withColumn('Athletes_Id', monotonically_increasing_id()+1)
3 athlete_df_2 = athlete_df_2.select('Athletes_Id', *athlete_df_2.columns[0:-1]).sort('Athletes_Id')
4 display(athlete_df_2)

```

## Join Athlete and Medals2

```

1 #Nối bảng athlete và bảng medal để tạo nên bảng fact
2 fact_df = athlete_df_2.join(medal_df, on='url', how='left_outer')
3 display(fact_df)

```

```

1 #Thay đổi những giá trị null ở medal_type thành no medal
2 fact_df_1 = fact_df.withColumn('medal_type', when(col('medal_type').isNull(),'No medal').otherwise(col('medal_type')))
3 fact_df_2 = fact_df_1.select('Athletes_Id', 'medal_type','country_code', 'discipline_code', 'discipline', 'event',
4 'event_sex')
4 display(fact_df_2)
5

```

```

1 # tạo discipline_df
2 wdl = Window.partitionBy('discipline_code')
3 fact_df_2 = fact_df_1.withColumn('NonNullValues', last('discipline', ignorenulls=True).over(wdl))
4 #lấy giá trị cuối cùng khi sắp xếp tăng dần bỏ qua các gtri null
5 # Thay các giá trị null bằng giá trị khác null đã tìm ra
6 fact_df_2 = fact_df_2.withColumn('discipline', col('NonNullValues'))
7 fact_df_2 = fact_df_2.drop('NonNullValues')
8
9 discipline_df = fact_df_2.select('discipline_code', 'discipline')
10 discipline_df = discipline_df.distinct()
11 display(discipline_df)

```

```

1 # Tạo bảng event
2 event_df = fact_df_2.select('discipline_code', 'event', 'event_sex')
3 event_df = event_df.distinct().dropna()
4 event_df = event_df.withColumn('Event_Id', monotonically_increasing_id() + 1)
5 event_df = event_df.select('Event_Id', *event_df.columns[0:-1]).sort('Event_Id')
6
7 # vì fact không có id event nên thực hiện tạo 1 khóa phụ để join bảng fact lấy id event
8 event_df = event_df.withColumn('event_sub', concat(col('discipline_code'), col('event'), col('event_sex'))))
9 event_df = event_df.withColumnRenamed('discipline_code', 'E_discipline_code')
10 display(event_df)

```

```

1 # tạo sub_key
2 fact_df_3 = fact_df_2.withColumn('discipline', concat(col('country_code'), col('discipline')))
3 fact_df_3 = fact_df_3.withColumnRenamed('discipline', 'sub_key')
4 # tạo khóa phụ để join với event
5 fact_df_3 = fact_df_3.withColumn('event_sub', concat(col('discipline_code'), col('event'), col('event_sex'))))
6 display(fact_df_3)

```

```

1 # join vs bảng event
2 fact_table = fact_df_3.join(event_df, on='event_sub', how='left_outer')
3 fact_table = fact_table.select('Athletes_Id', 'medal_type', 'country_code', 'discipline_code', 'Event_Id')
4 display(fact_table)

```

```

1 # xóa khóa phụ đã tạo để dc lại bảng event ban đầu
2 event_df = event_df.drop('event_sub')
3 event_df = event_df.withColumnRenamed('E_discipline_code', 'discipline_code')
4 display(event_df)

```

```
1 # import bảng coaches, làm sạch và tạo ID
2 file_location2 = f"wasbs://{{container}}@{{storage_account_name}}.blob.core.windows.net/coaches.csv"
3 file_type = "csv"
4 coaches_df = spark.read.format(file_type).option("header", "true").load(file_location2)
5 coaches_df_1 = coaches_df.withColumn('sub_key', concat(col('country_code'), col('discipline')))
6 display(coaches_df_1)
```

```
1 #Tạo id
2 coaches_df_2 = coaches_df_1.withColumn('Coaches_Id', monotonically_increasing_id()+1)
3 coaches_df_2 = coaches_df_2.select('Coaches_Id', *coaches_df_2.columns[0:-1]).sort('Coaches_Id')
4
5 display(coaches_df_2)
```

```
1 # tạo bảng total huy chương theo từng Quốc gia
2 total = fact_table.filter(col('medal_type') != 'No medal')
3 total.groupBy('country_code').count().display()
```

# Save dataframes to Azure Storage

```
1 # Đường dẫn đến Azure Blob Storage
2 blob_storage_path = f"wasbs://container2@huuphan.blob.core.windows.net/transformed/event"
3 blob_storage_path1 = f"wasbs://container2@huuphan.blob.core.windows.net/transformed/country"
4 blob_storage_path2 = f"wasbs://container2@huuphan.blob.core.windows.net/transformed/discipline"
5 blob_storage_path3 = f"wasbs://container2@huuphan.blob.core.windows.net/transformed/coaches"
6 blob_storage_path4 = f"wasbs://container2@huuphan.blob.core.windows.net/transformed/athlete"
7 blob_storage_path5 = f"wasbs://container2@huuphan.blob.core.windows.net/transformed/fact_table"
8
9 # Lưu DataFrame vào Azure Blob Storage dưới dạng file CSV
10 event_df.coalesce(1).write.parquet(blob_storage_path, mode="overwrite")
11 country_df.coalesce(1).write.parquet(blob_storage_path1, mode="overwrite")
12 discipline_df.coalesce(1).write.parquet(blob_storage_path2, mode="overwrite")
13 coaches_df_2.coalesce(1).write.parquet(blob_storage_path3, mode="overwrite")
14 athlete_df_2.coalesce(1).write.parquet(blob_storage_path4, mode="overwrite")
15 fact_table.coalesce(1).write.parquet(blob_storage_path5, mode="overwrite")
```

▶ (25) Spark Jobs

Command took 8.62 seconds -- by ph.thienhuu@gmail.com at 22:34:52 22/8/2023 on 22/8

# Recording data to SQL database

```
1 server_name = 'huu-server'
2 database_name = 'olympic2020'
3 #Tạo secrets cho username và password sql
4 user_name = dbutils.secrets.get(scope ='keyofme', key ='usernamesql')
5 pass_word = dbutils.secrets.get(scope ='keyofme', key ='passwordsql')
6
7 list_table = [
8     (event_df, 'Dim_Events'),
9     (country_df, 'Dim_Countries'),
10    (discipline_df, 'Dim_Disiplines'),
11    (athlete_df_2, 'Dim_Athletes'),
12    (coaches_df_2, 'Dim_Coaches'),
13    (fact_table, 'Fact_Table')
14 ]
15 for i, j in list_table :
16     i.write\
17         .mode("overwrite")\
18         .format("jdbc")\
19         .option("url", f"jdbc:sqlserver://{{server_name}}.database.windows.net;databaseName={{database_name}};" )\
20         .option("dbtable", j)\ 
21         .option("user", user_name)\ 
22         .option("password", pass_word)\ 
23         .option("driver", "com.microsoft.sqlserver.jdbc.SQLServerDriver")\ 
24         .save()
```

```

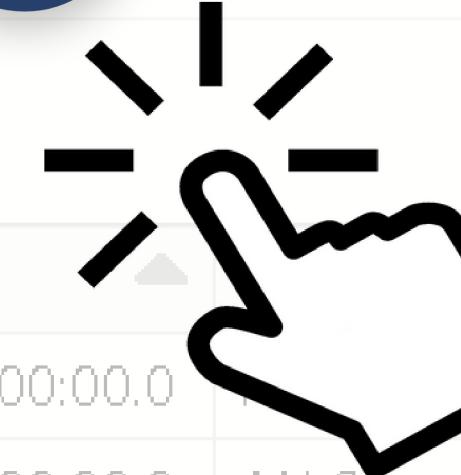
1 container = "container1"
2 storage_account_name = "huuphan"
3 spark.conf.set(
4     "fs.azure.account.key."+storage_account_name+".blob.core.windows.net",
5     dbutils.secrets.get("keyofme", "top"))
6 file_location = f"wasbs://{container}@{storage_account_name}.blob.core.windows.net/medals2.csv"
7 file_type = "csv"
8 df = spark.read.format(file_type).option("header", "true").load(file_location)
9 display(df)

```

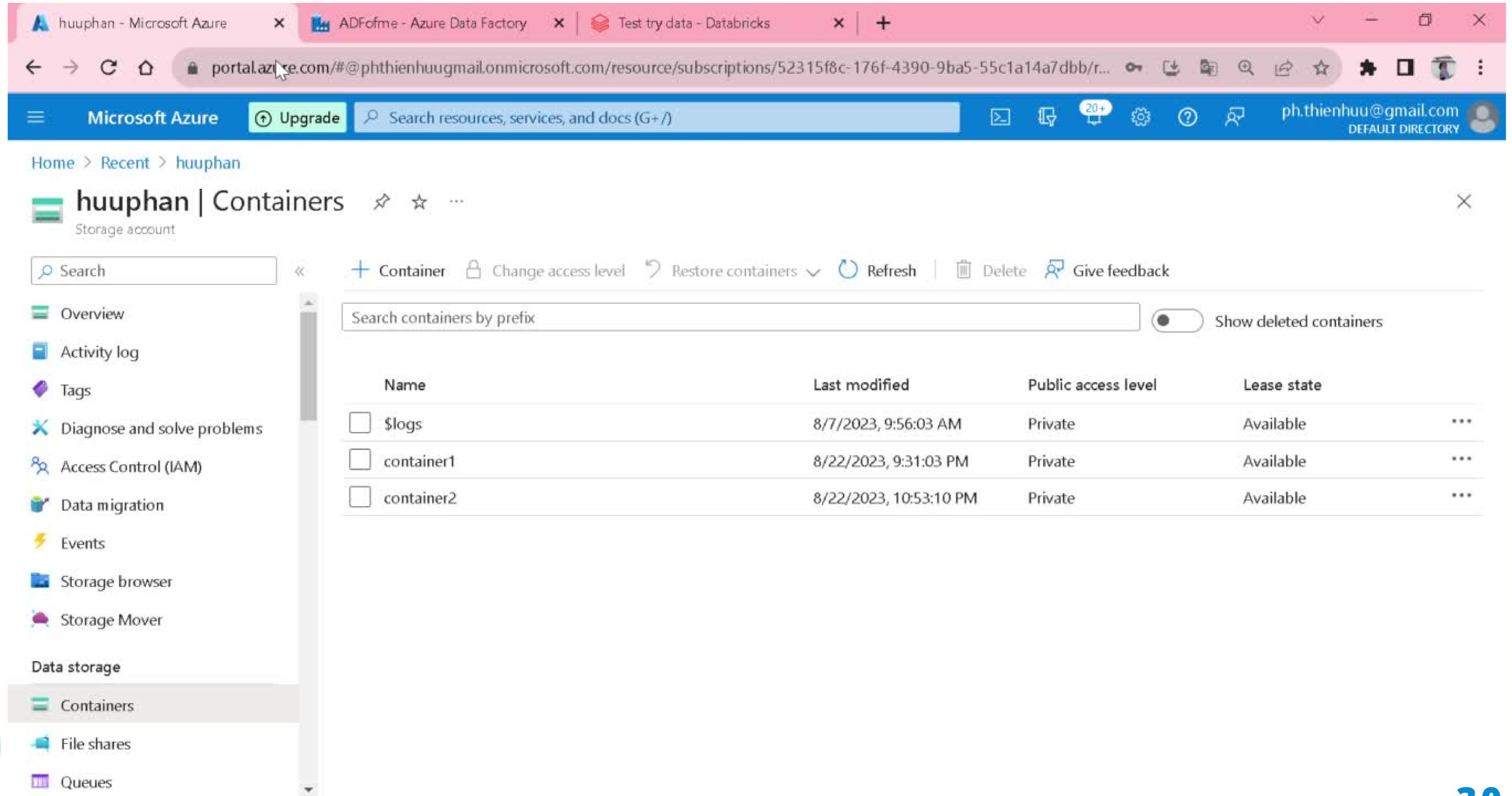
▶ (2) Spark Jobs

▶ 📈 df: pyspark.sql.dataframe.DataFrame = <DataFrame [c0, medal\_type, medal\_code, medal\_date, sport\_name, athlete\_name]> [1 more fields]

# DEMO



_c0	medal_type	medal_code	medal_date	sport_name	athlete_name
1	Gold Medal	1	2021-07-24 00:00:00.0		KIM Je Deok
2	Gold Medal	1	2021-07-24 00:00:00.0	AN S	AN San
3	Silver Medal	2	2021-07-24 00:00:00.0	SCHLOESSER G	SCHLOESSER Gabriela
4	Silver Medal	2	2021-07-24 00:00:00.0	WIJLER S	WIJLER Steve
5	Bronze Medal	3	2021-07-24 00:00:00.0	ALVAREZ L	ALVAREZ Luis
6	Bronze Medal	3	2021-07-24 00:00:00.0	VALENCIA A	VALENCIA Alejandra
7	Gold Medal	1	2021-07-24 00:00:00.0	CARAPAZ R	CARAPAZ Richard

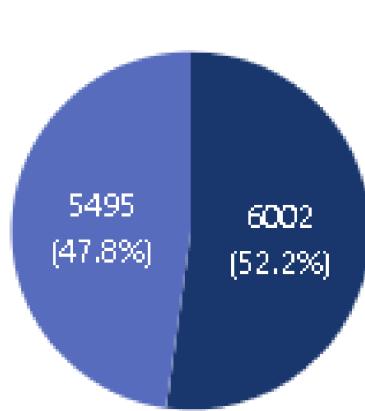


# VISUALIZATION



# TOKYO 2020 OLYMPIC SUMMER GAMES

## Distribution of medals by gender



## Total Country

**206**

Count of country\_code

## Total Athletes

**11.656K**

Count of Athletes\_Id

Medal Type	Country
All	All

## Distribution of medal by name

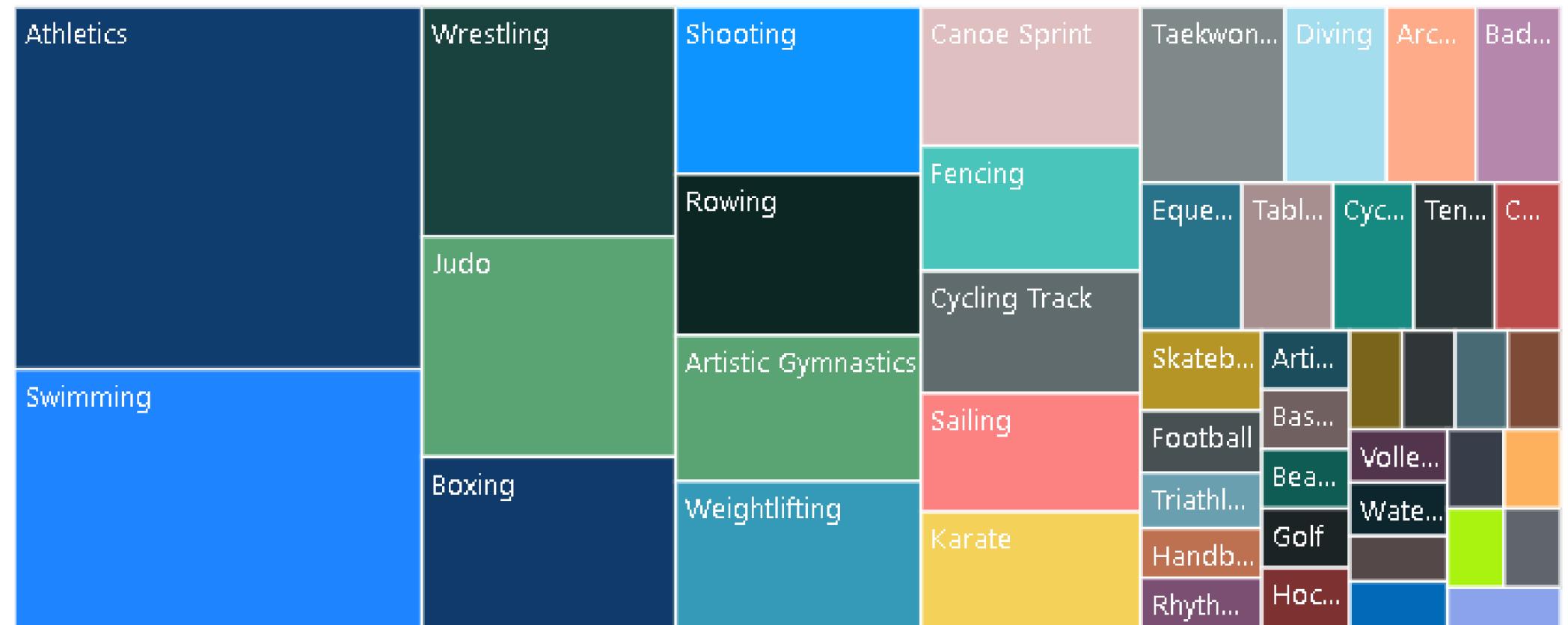
medal\_type ● Gold Medal ● Silver Medal ● Bronze Medal



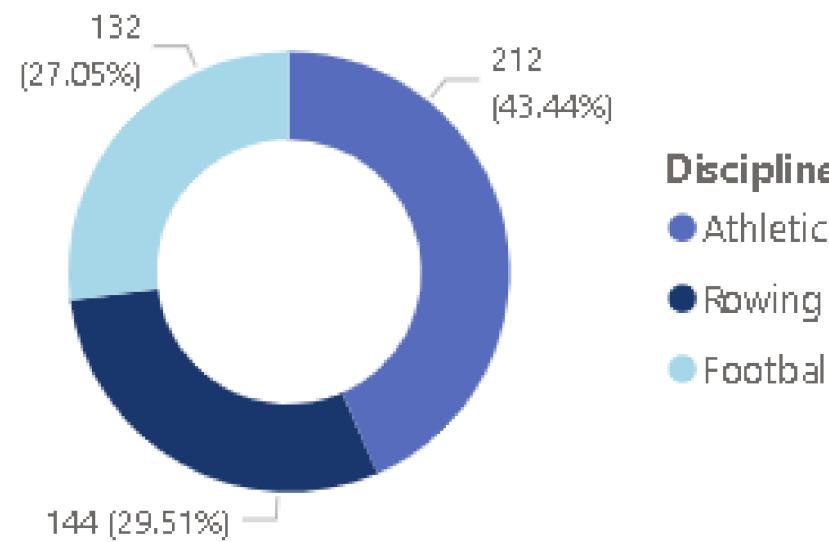
## Top Country

Country	Gold Medal	Silver Medal	Bronze Medal	Total
United States of America	39	41	33	113
People's Republic of China	38	32	18	88
ROC	20	28	23	71
Great Britain	22	21	22	65
Japan	27	14	17	58
Australia	17	7	22	46
Italy	10	10	20	40
Germany	10	11	16	37
Netherlands	10	12	14	36
France	10	12	11	33
Canada	7	6	11	24
Brazil	7	6	8	21
Hungary	6	7	7	20
New Zealand	7	6	7	20
Republic of Korea	6	4	10	20
Ukraine	1	6	12	19
<b>Total</b>	<b>340</b>	<b>338</b>	<b>402</b>	<b>1080</b>

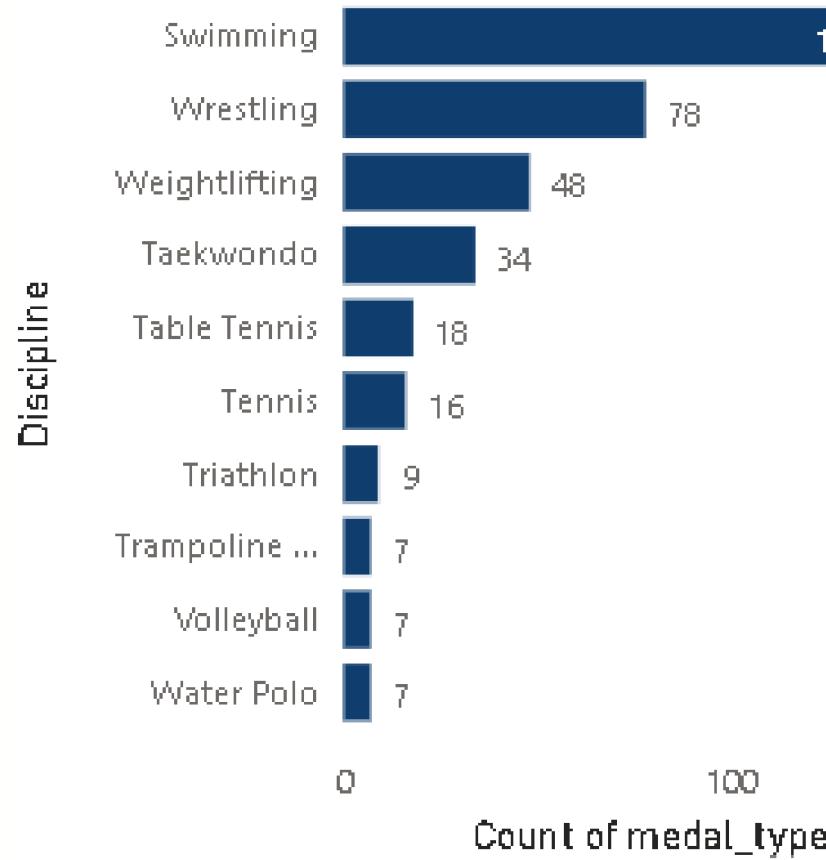
## Distribution of medals by discipline



## Top 3 Discipline with most Athletes

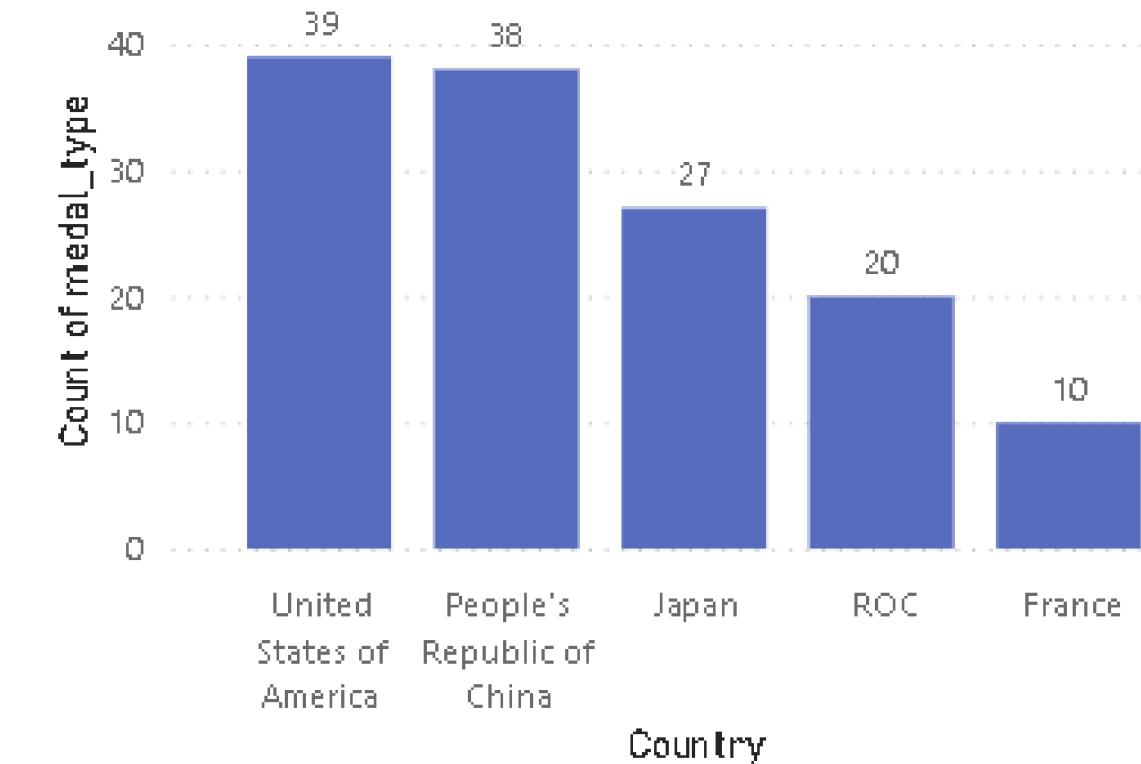


## Top 10 Discipline with most Medals



## Top 5 Country with most Gold Medal

medal\_type ● Gold Medal



GOLD

340

SILVER

338

BRONZE

402

Medals

All

Country

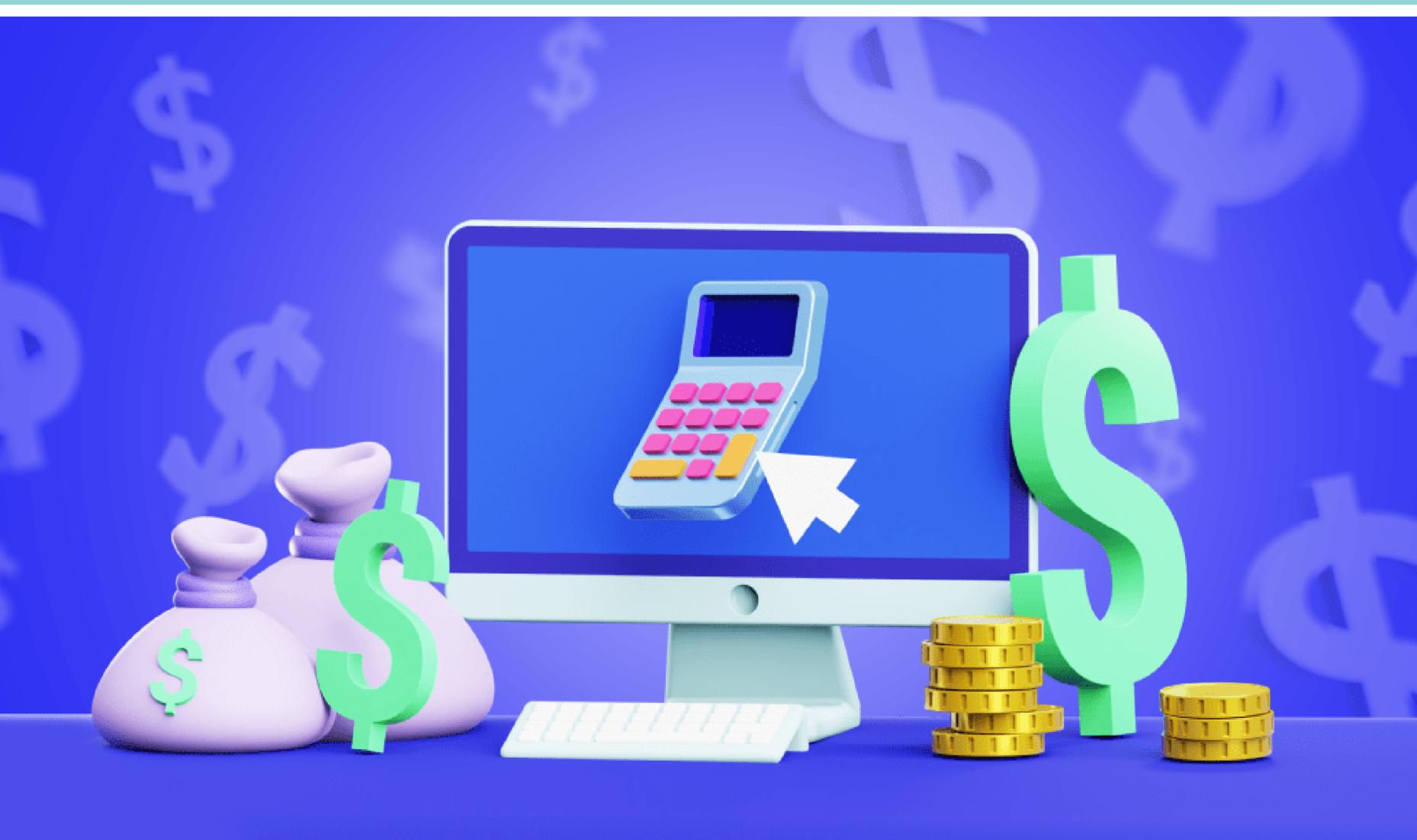
All

Medals ● Bronze Medal ● Gold Medal ● Silver Medal

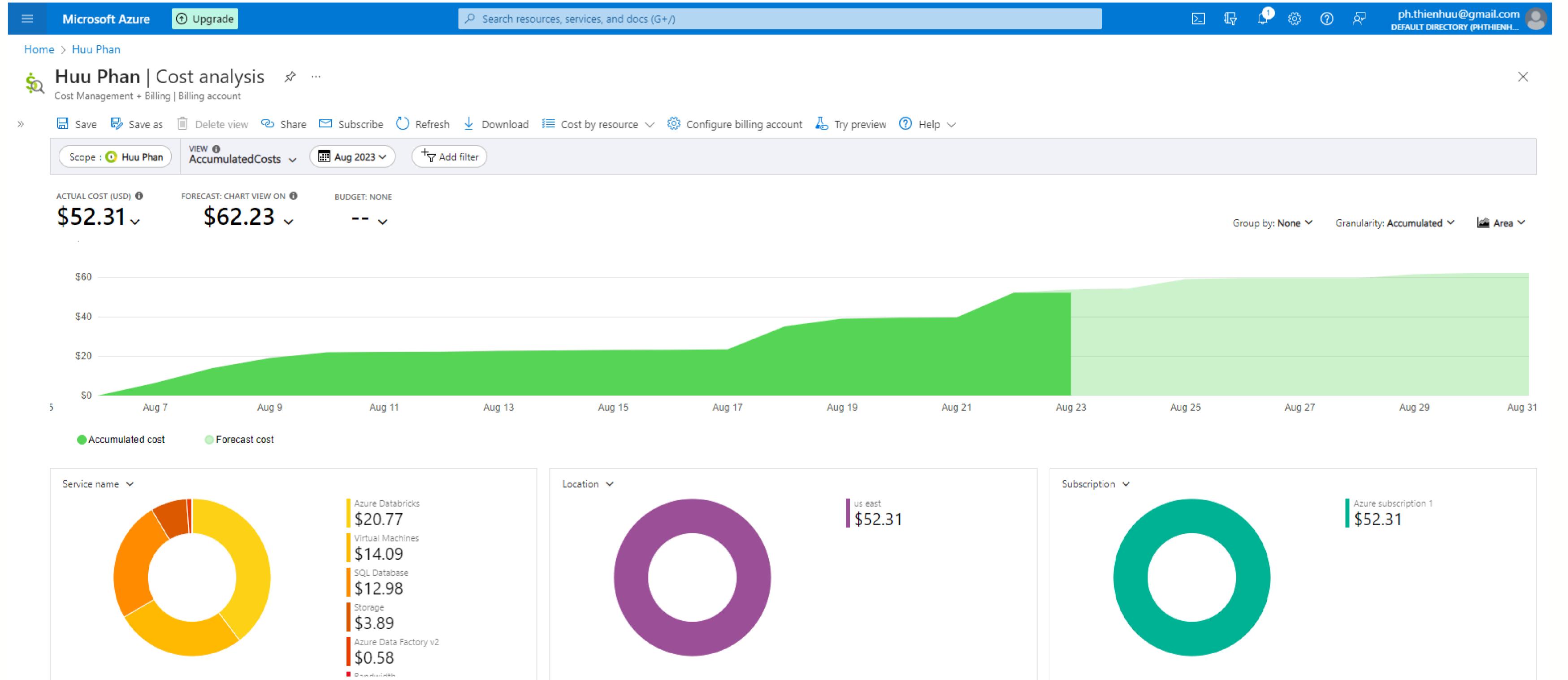


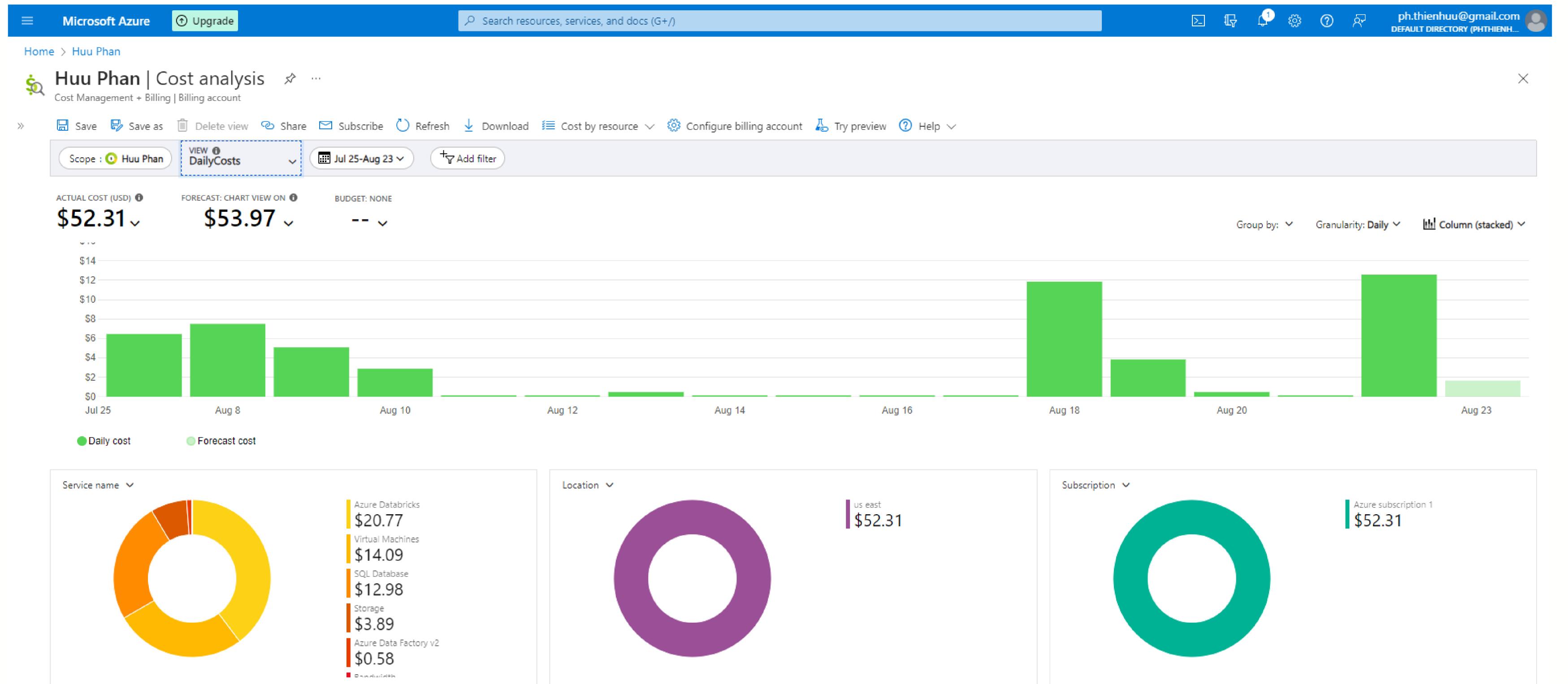
## Top Athletes

Athletes Name	Bronze Medal	Gold Medal	Silver Medal	Total
McKEON Emma	3	4		7
DRESSEL Caeleb		5		5
LEDECKY Kathleen		2	2	4
McKEOWN Kaylee	1	3		4
SCOTT Duncan		1	3	4
TITMUS Ariarne	1	2	1	4
ZHANG Yufei		2	2	4
AN San		3		3
BATSARASHKINA Vitalina		2	1	3
CAMPBELL Cate	1	2		3
CARRINGTON Lisa		3		3
<b>Total</b>	<b>25</b>	<b>45</b>	<b>22</b>	<b>92</b>



# COST MANAGEMENT





# THANK YOU

