



ETL PROJECT: FOOTBALL DATA ANALYTICS

MISBAH BIN HOSSAIN

DATA ENGINEER STUDENT AT NACKADEMIN

The ETL workflow consists of:

1. Data Extraction

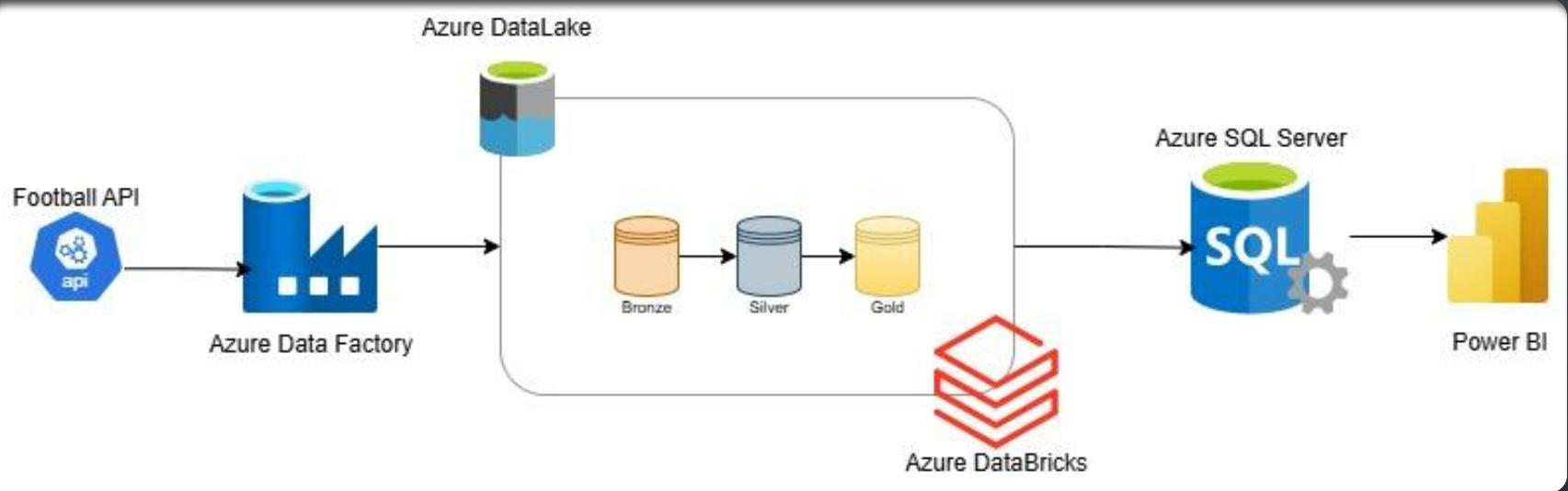
2. Data Transformation

3. Data Loading

4. Data Visualization

Each phase is supported by different Azure and data engineering tools.

ETL WORKFLOW OVERVIEW



ETL WORKFLOW

FOOTBALL API

- WEBSITE - [HTTPS://WWW.FOOTBALL-DATA.ORG/](https://www.football-data.org/)
- MY PARAMETERS:
 - FILTER = **["MATCHES","STANDINGS","TEAMS","SCORERS"]**
 - SEASON= [2022,"2023","2024"]
 - LEAGUE = ["PL","SA","PD","BL1"]
- API CALL EXAMPLE:
 - HTTPS://API.FOOTBALL-DATA.ORG/v4/COMPETITIONS/**PL/MATCHES/?SEASON=2024**
 - HTTPS://API.FOOTBALL-DATA.ORG/v4/COMPETITIONS/**PD/TEAMS/?SEASON=2023**
 - HTTPS://API.FOOTBALL-DATA.ORG/v4/COMPETITIONS/**SA/SCORERS/?SEASON=2022**
 - HTTPS://API.FOOTBALL-DATA.ORG/v4/COMPETITIONS/**BL1/STANDINGS/?SEASON=2023**

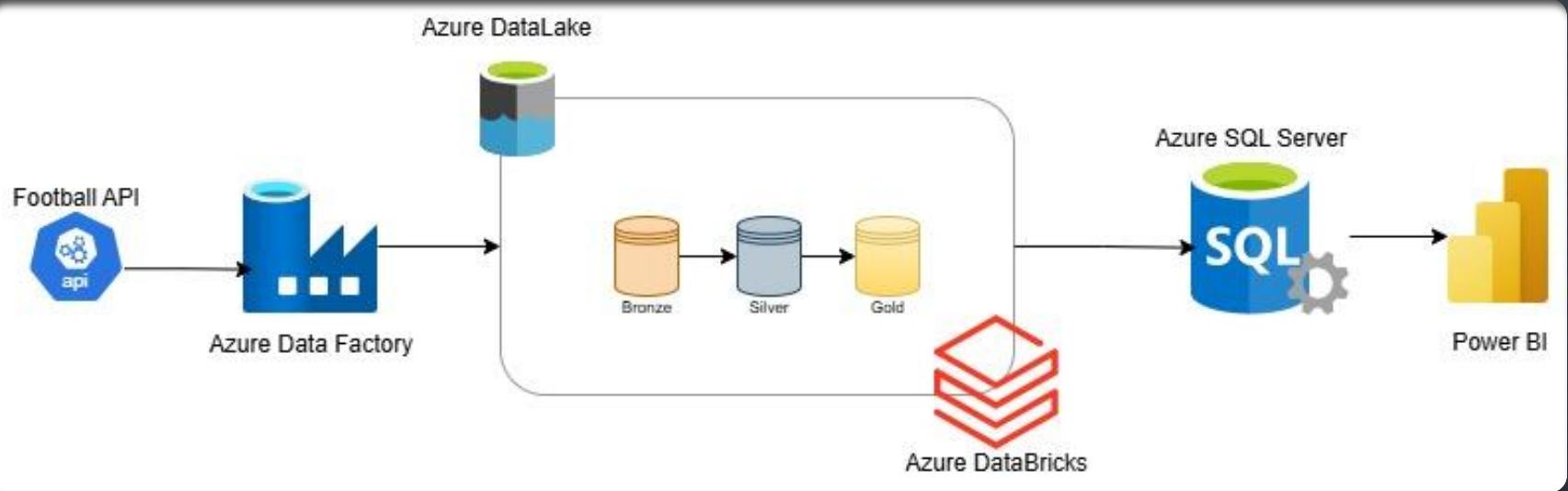
Free
€0,00 /mo

12 competitions
Scores delayed
Fixtures
Schedules delayed
League Tables

—
—
—

10 calls/minute

You have this plan.

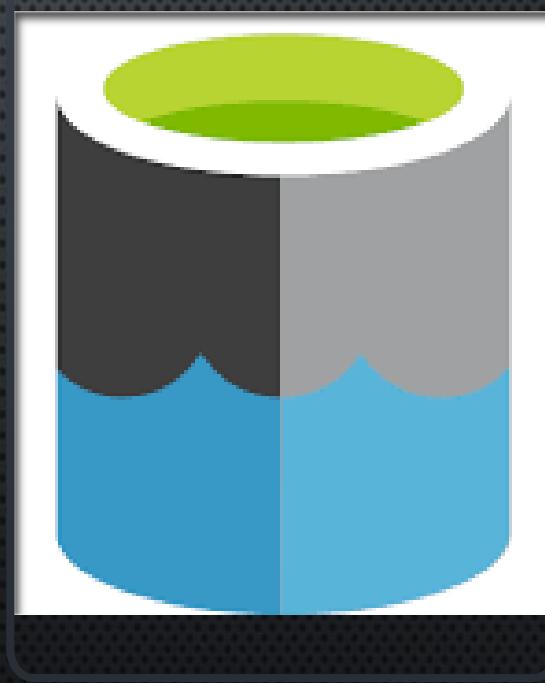


ETL WORKFLOW

DATA EXTRACTION

TOOLS: AZURE DATA FACTORY, AZURE BLOB STORAGE

- - DYNAMIC PIPELINES FETCH FOOTBALL DATA FROM API
- - HANDLES MULTIPLE PARAMETERS (SEASON, LEAGUE, FILTER TYPES)
- - DATA STORED IN PARQUET FORMAT IN AZURE BLOB STORAGE.



General Source Sink Mapping Settings User properties

Source dataset *

	Dynamic import		Open		New		Preview data		Learn more
Dataset properties									
Name									Type
Filter									string
Season									string
League									string

General Source Sink Mapping Settings User properties

Sink dataset *

	Standing		Open		New		Preview data		Learn more
Dataset properties									
Name									Type
Fil									string
Sea									string
lea									string

Copy behavior

Flatten hierarchy

General Source Sink Mapping Settings User

Import schemas

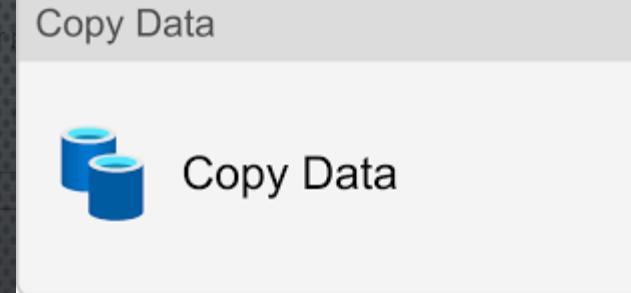
New mapping

Clear

Collection reference

\$['standings'][0]['table']

Map complex values to string



Name	Type	Collection reference	Column name	Type	Include
standings	array	<input type="checkbox"/>			<input checked="" type="checkbox"/>
table	array	<input checked="" type="checkbox"/>			<input checked="" type="checkbox"/>

General Source Sink Mapping Settings User properties

Source dataset *

 Dynamic_import

 Open

 New

 Preview data

 Learn more

Dataset properties 

Name	Value	Type
Filter	@item().filter	string
Season	@item().season	string
League	@item().league	string

General Source Sink Mapping Settings User properties

Sink dataset *

 Standing

 Open

 New

 Learn more

Dataset properties 

Name	Value	Type
Fil	@item().filter	string
Sea	@item().season	string
lea	@item().league	string

Copy behavior 

Flatten hierarchy

General Source Sink Mapping Settings User properties

 Import schemas

 New mapping

 Clear 

 Delete

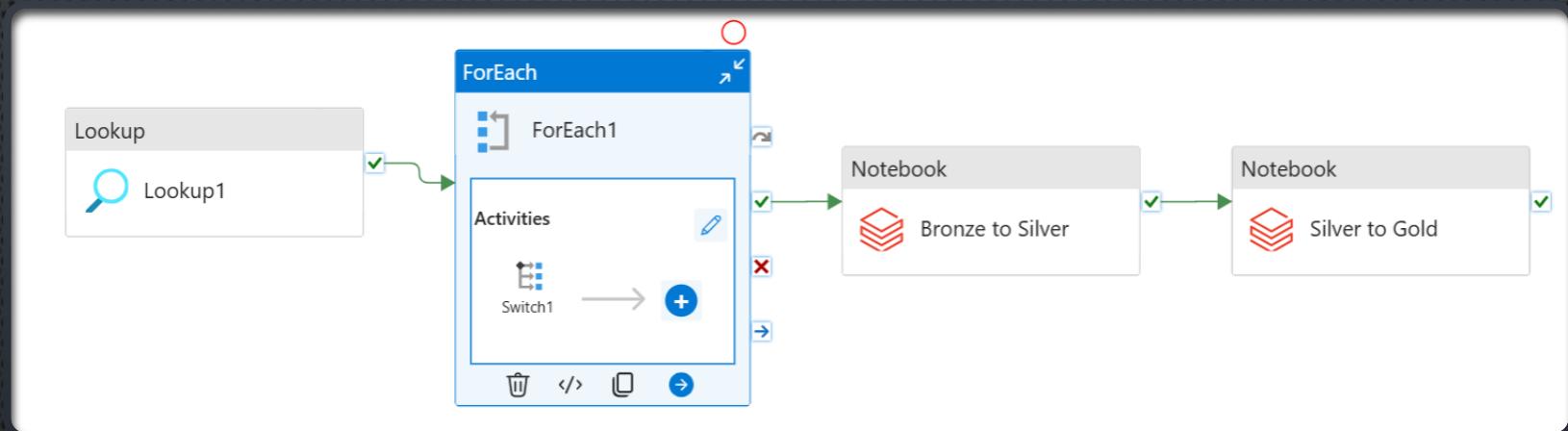
 Advanced editor

Collection reference 

\$['standings'][0]['table']

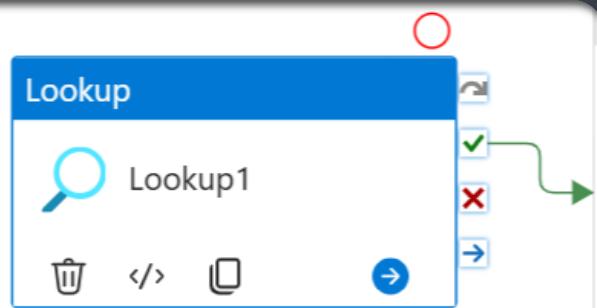
Map complex values to string

Name	Type	Collection reference	Column name	Type	<input checked="" type="checkbox"/> Include
standings	[] array	<input type="checkbox"/>			<input checked="" type="checkbox"/> Include
table	[] array	<input checked="" type="checkbox"/>			<input checked="" type="checkbox"/> Include



ADF PIPELINE

LOOKUP ACTIVITY



General **Settings** User properties

Source dataset *

AzureSqlTable1 Open

First row only

Use query

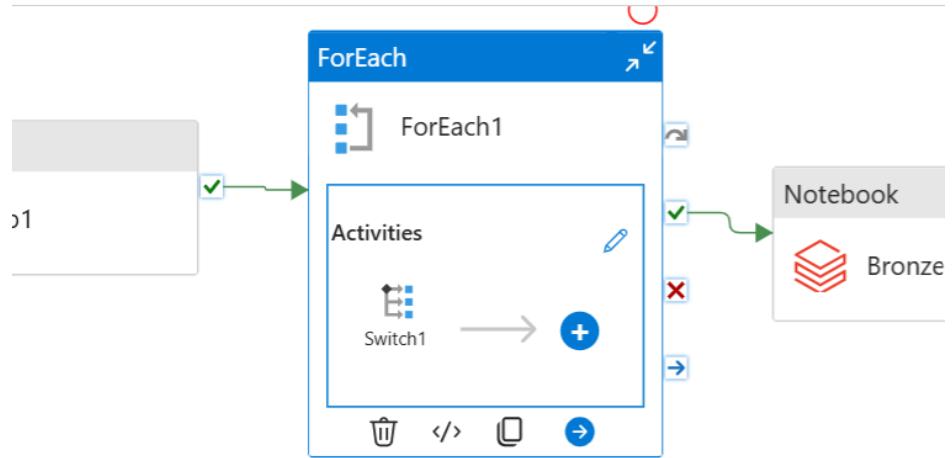
Table Query Stored procedure

Query *

```
select filter, league, season  
from [dbo].[FootballParameters]
```

Edit

FOREACH ACTIVITY



General **Settings** Activities (1) User properties

Sequential

Items

@activity('Lookup1').output.value

SWITCH ACTIVITY

General Activities (8) User properties

Expression ⓘ @item().filter

+ Add case

Case ⓘ

	Activity
Default	No activities
standings	Wait_Stand ing Copy data -Stan...
scorers	Wait_Scoreres Copy data -Top ...
teams	Wait_Teams Dataflow_Teams
matches	Wait_Matches Copy data -Mat...



source1

Import data from
Dynamic_import

flatten2

Unrolling arrays from
body.teams to with columns
'competition', 'count', 'filters',
'season', 'address', 'area'

Reference:
1

SinkRunningCom

Columns:
4 total

Export data to
RunningCompetition_Fla

flatten2

Unrolling arrays from
body.teams to with columns
'competition', 'count', 'filters',
'season', 'address', 'area'

flatten1

Unrolling arrays from squad to
with columns 'id', 'name',
'nationality', 'position',
'dateOfBirth'

SinkSquad

Export data to Squad_fla

DATAFLOW

Unrolling arrays from
body.teams to with columns
'competition', 'count', 'filters',
'season', 'address', 'area'

Unroll Team

Export data to
Teams_Flatten_Dataflow

Parameters

Settings

+ New

Delete

Name

Fil

Sea

Lea

Type

abc string

abc string

abc string

Default value

Enter expression...

ANY



Enter expression...

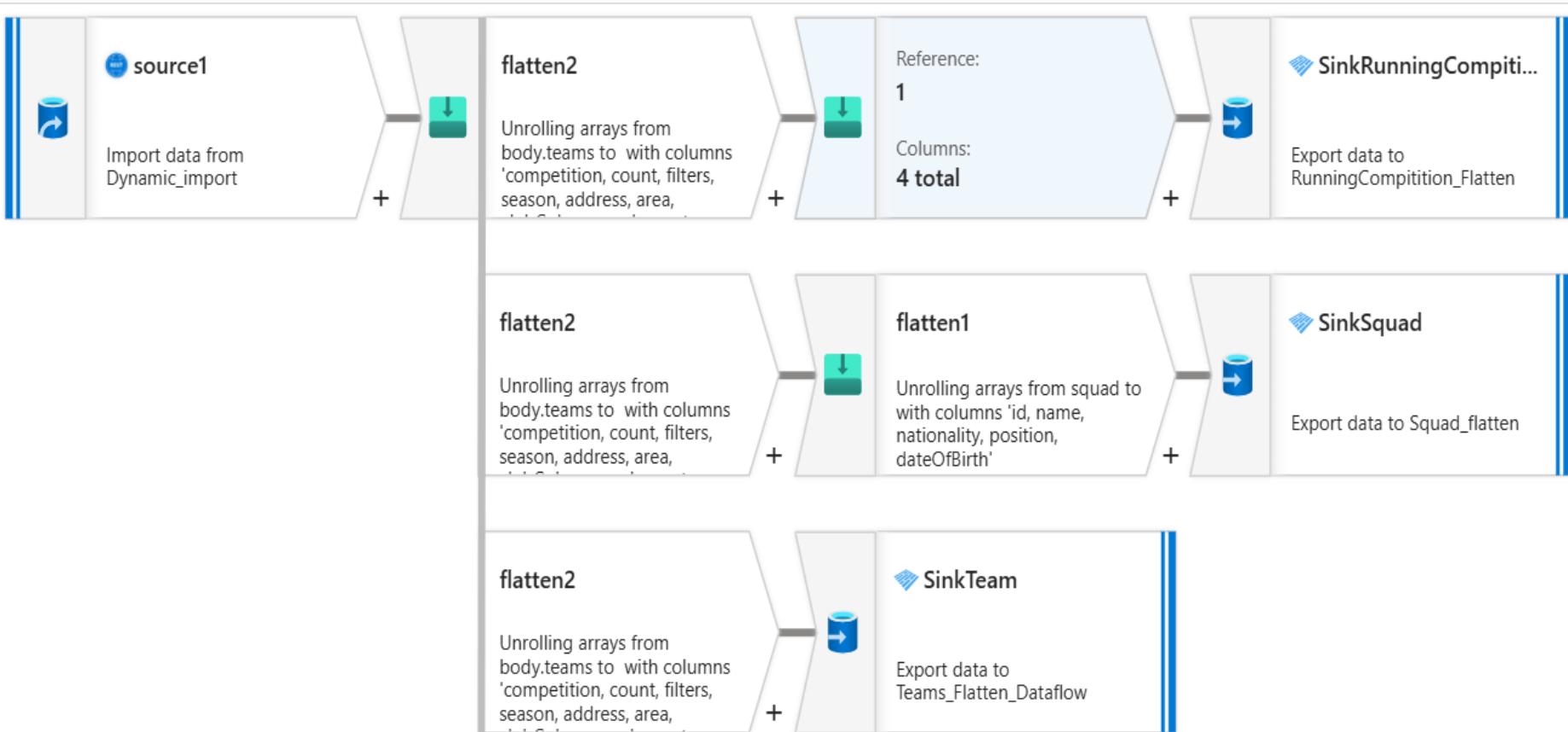
ANY



Enter expression...

ANY





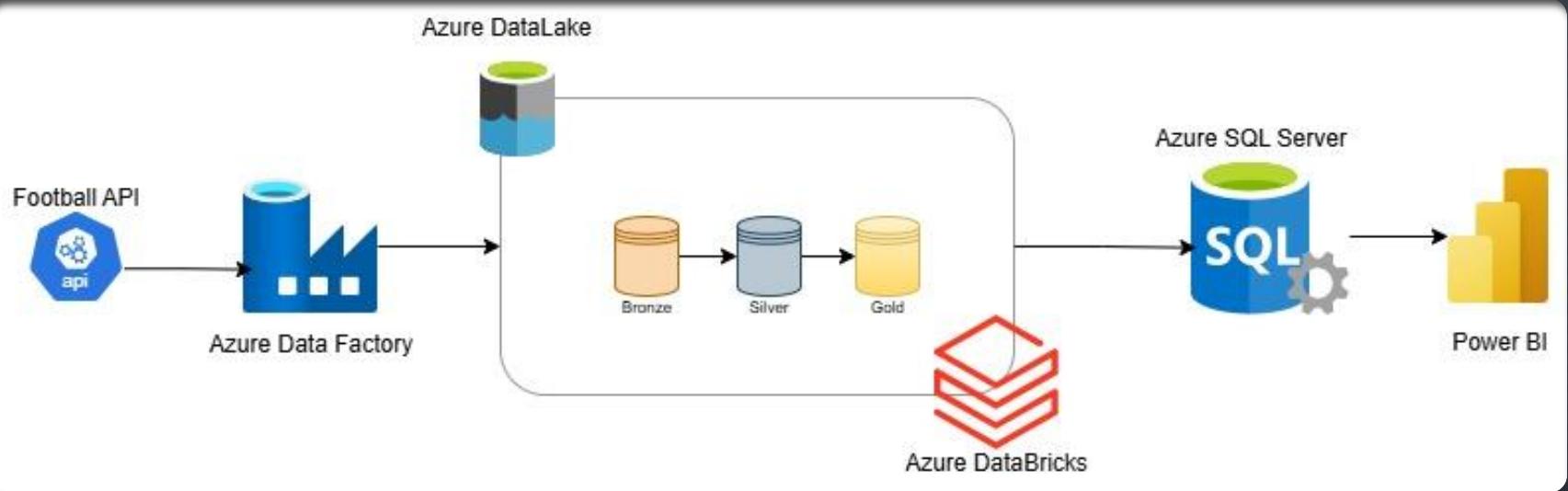
Parameters

Settings

+ New

Delete

	Name	Type	Default value	
<input type="checkbox"/>	Fil	abc string	Enter expression... ANY	
<input type="checkbox"/>	Sea	abc string	Enter expression... ANY	
<input type="checkbox"/>	Lea	abc string	Enter expression... ANY	



ETL WORKFLOW

Storage accounts

misbahetblob | Storage browser

Storage account

Default Directory (..)

+ Create ⌂ Restore ...

Filter for any field...

Name ↑

- dbstorageucivez4da5h64
- misbahetblob

Overview

Activity log

Tags

Diagnose and solve problems

Access Control (IAM)

Data migration

Events

Storage browser

Partner solutions

Data storage

Security + networking

Networking

Access keys

Shared access signature

Encryption

Microsoft Defender for Cloud

Search

misbahetblob

Favorites

Recently viewed

Blob containers

- \$logs
- esco
- football-data
- synapse-fs

View all

File shares

Queues

Tables

Add Directory ⌂ Upload ⌂ Refresh ⌂ Delete ⌂ Copy ⌂ Paste ⌂ Rename ...

Blob containers > football-data

Authentication method: Access key (Switch to Microsoft Entra user account)

Search blobs by prefix (case-sensitive) Only show active objects

Showing all 4 items

<input type="checkbox"/>	Name	Last modified	Access tier	Size
<input type="checkbox"/>	Bronze	11/7/2024, 1:17:45 PM		
<input type="checkbox"/>	Gold	11/9/2024, 10:35:42 AM		
<input type="checkbox"/>	Silver	12/3/2024, 9:50:50 PM		
<input type="checkbox"/>	\$_azurertmp...	11/18/2024, 12:03:31 AM		

< Page 1 < of 1 >

The screenshot shows the Azure Storage Explorer interface. On the left, a sidebar lists storage accounts: 'dbstorageucivez4da5h64' and 'misbahetblob'. The 'misbahetblob' account is selected. The main pane displays the 'Storage browser' for the 'misbahetblob' account. It shows a tree view with 'misbahetblob' at the root, containing 'Blob containers' which include '\$logs', 'esco', 'football-data' (selected), and 'synapse-fs'. Below this is a 'View all' section for 'File shares', 'Queues', and 'Tables'. At the top right, there are navigation and action buttons like '+ Add Directory', 'Upload', 'Refresh', 'Delete', 'Copy', 'Paste', 'Rename', and more. The 'football-data' container is expanded, showing four blob items: 'Bronze', 'Gold', 'Silver', and '\$_azurertmp...'. A search bar at the top right allows searching by blob prefix. The bottom left shows pagination controls ('Page 1 of 1').

AZURE BLOB

DATA TRANSFORMATION

TOOLS: AZURE DATABRICKS (PYSPARK)

PROCESS:

- ****BRONZE TO SILVER**:** CLEAN AND NORMALIZE RAW DATA
- ****SILVER TO GOLD**:** AGGREGATE AND PREPARE DATA FOR INSIGHTS

OUTCOME:

HIGH-QUALITY, ANALYSIS-READY DATASETS.



databricks

Bronze to Silver Python

Last edit was 2 days ago

Run all Connecting Schedule Share

1: Install necessary Packages

```
# Install necessary packages %pip install azure-storage-blob pandas pyarrow
```

2: Merge all files to fixed tables

```
%python import os from azure.storage.blob import BlobServiceClient import pandas ...
```

3: Observing DataFrames

```
from azure.storage.blob import BlobServiceClient from io import BytesIO # Azure ...
```

4: Fix running competition

```
df_teams_running_competition = load_and_display(file_paths[3]).drop_duplicates(s ...
```

5: Fix Teams table

```
%python # Assuming df_teams_team is a pandas DataFrame # 1. Delete specified co ...
```

6: More necessary Packages

```
%python %pip install fsspec
```

7: Fix Season table (Explode from Teams)

```
# Set up the Azure Blob Storage account details account_name = 'misbahetblob' a ...
```

8: Fix Matches Table

```
# Function to rename columns by removing the prefix def rename_columns(df, prefix ...
```

9: Fix Coach table (Explode from Teams)

```
# Extract dictionary values from 'coach' column and create separate columns for ...
```

10: Fix Squad table (Explode from Teams)

```
# Extract dictionary values from 'squad' column and create separate columns for ...
```

11: Fix Teams table

```
df_teams_team['season_id'] = df_teams_team['season'].apply(lambda x: x['id'] if ...
```

Silver to gold

Python



File Edit View Run Help Last edit was 6 days ago

Run all

Terminated

Schedule

Share



1: Install necessary packages
Install necessary packages %pip install azure-storage-blob pandas pyarrow

2: Display DF
import os from azure.storage.blob import BlobServiceClient import pandas as pd f ...

3: Fix all Date types (for all date type)
for name, df in dataframes.items(): for column in df.columns: if "da ..."

4: Fix df_scorers for null values
df_scorers = dataframes["df_scorers"].drop(columns=["position", "shirtNumber"]). ...

5: Checking DF for cleaning
Function to check for missing values and data types in each dataframe def chec ...

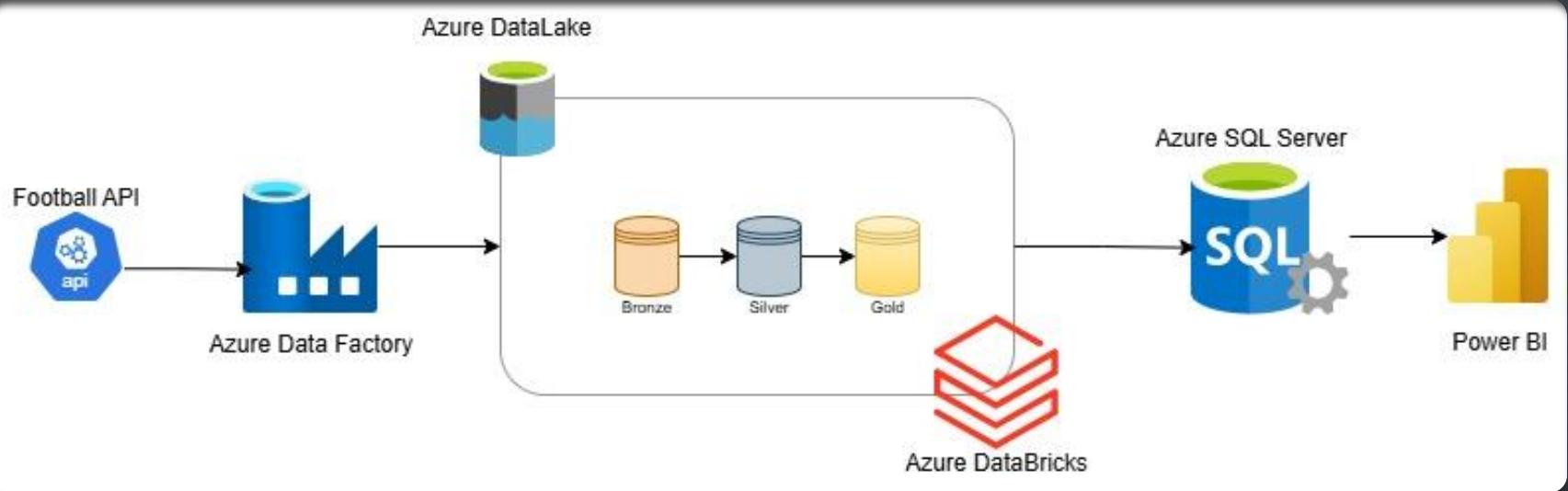
6: Fix Primary Keys
List of dataframes to save to SQL Server dataframes_to_save = { "df_season" ...

7: Install Python
%python

8: Export df to GOLD
Azure Blob Storage details account_name = 'misbahetlblob' connection_string = ...

9: Export to SQL Server
Define the JDBC URL and connection properties jdbc_url = "jdbc:sqlserver://foo ..."

[Shift+Enter] to run and move to next cell



ETL WORKFLOW

DATA LOADING

TOOLS: AZURE DATA BRICKS, SQL SERVER

- - TRANSFORMED DATA LOADED INTO SQL SERVER VIA JDBC
- - ENSURES EFFICIENT STORAGE AND ACCESS FOR ANALYSIS.



SQL databases

Default Directory (misbahbinhossain111@gmail.on...)

+ Create Reservations ...

Filter for any field...

Name ↑↓

football (footballapi/football) ...

Search

Overview Activity log Tags Diagnose and solve problems **Query editor (preview)**

Mirror database in Fabric (preview)

> Settings

> Data management

> Integrations

> Power Platform

> Security

> Intelligent performance

> Monitoring

> Automation

> Help

Page 1 of 1 >

football (footballapi/football) | Query editor (preview)

SQL database

Login New Query Open query Feedback Getting started

Showing limited object explorer here. For full capability please click here to open Azure Data Studio.

Tables

- > dbo.df_matches
- > dbo.df_scorers
- > dbo.df_season_details
- > dbo.df_squad_details
- > **dbo.df_standings**
- > dbo.df_teams_team
- > dbo.FootballParameters

Views

Stored Procedures

Query 2

Run Cancel query Save query Export data as Show only Editor Open Copilot

```
1 SELECT TOP (1000) * FROM [dbo].[df_standings]
```

Results

Standings_id	position	team_id	team_r
59	5	94	Villarre
1	1	5	FC Bay
176	2	110	SS Lazi
60	6	90	Real Re

DATA VISUALIZATION

TOOL: POWER BI

- - TEAM STANDINGS AND RANKINGS
- - MATCH RESULTS AND SUMMARIES
- - TOP SCORER STATISTICS
- - SQUAD-LEVEL DETAILS



Football Data Analysis

34

Max of points

5

Min of points

379

Total Goals

20

Total Teams

Competition Name

- Bundesliga
- Premier League
- Primera Division
- Serie A

Season

- 2022
- 2023
- 2024

14

Current Matchday

Position	Team Name	Points	PG	Won	Lost	Draw	GF	GA	GD	FormIcons
1	Liverpool FC	34	13	11	1	1	26	8	18	✓✓✓✓□
2	Arsenal FC	25	13	7	2	4	26	14	12	✓✓□×□
3	Chelsea FC	25	13	7	2	4	26	14	12	✓✓□□✓
4	Brighton & Hove Albion FC	23	13	6	2	5	22	17	5	□□□□□
5	Manchester City FC	23	13	7	4	2	22	19	3	✗✗✗✗✓
6	Nottingham Forest FC	22	13	6	3	4	16	13	3	✓✗✗✓✓
7	Tottenham Hotspur FC	20	13	6	5	2	28	14	14	□□□□□
8	Brentford FC	20	13	6	5	2	26	23	3	✓✓□✗✓
9	Manchester United FC	19	13	5	4	4	17	13	4	✓□✓□✗
10	Fulham FC	19	13	5	4	4	18	18	0	□□□□□
11	Newcastle United FC	19	13	5	4	4	14	14	0	□□□□□
12	Aston Villa FC	19	13	5	4	4	19	22	-3	✗✗✗□
13	AFC Bournemouth	18	13	5	5	3	20	19	1	✓✗✗✓□
14	West Ham United FC	15	13	4	6	3	17	24	-7	✗✓□✗✓
15	Everton FC	11	13	2	6	5	10	21	-11	✗□□✗□
16	Leicester City FC	10	13	2	7	4	16	27	-11	✗✗✗□✗

Season

- 2022
- 2023
- 2024

Competition Name

- Bundesliga
- Premier League
- Primera Division
- Serie A

Nationality

- Select all
- Albania
- Algeria
- Andorra
- Angola

TeamName

- 1. FC Heidenheim 1846
- 1. FC Köln
- 1. FC Union Berlin
- 1. FSV Mainz 05
- AC Milan
- AC Monza
- ACF Fiorentina
- AFC Bournemouth
- Arsenal FC
- AS Roma
- Aston Villa FC
- Atalanta BC
- Athletic Club
- Bayer 04 Leverkusen
- Bologna FC 1909
- Borussia Dortmund
- Borussia Mönchengladbach

name	Age	nationality	position	UI
Aarón	29	Spain	Goalkeeper	UI
Aaron Ciammaglichella	19	Italy	Central Midfield	Tc
Aaron Cresswell	35	England	Left-Back	WB
Aaron Hickey	22	Scotland	Right-Back	Br
Aaron Keto-Diyawa	21	England	Left-Back	WB
Aaron Martín	27	Spain	Left-Back	1.
Aaron Martín	27	Spain	Left-Back	Ge
Aaron Ramsdale	26	England	Goalkeeper	Ar
Aaron Ramsdale	26	England	Goalkeeper	Sc
Aaron Ramsey	21	England	Attacking Midfield	Bu
Aaron Seydel	28	Germany	Offence	S\
Aaron Wan-Bissaka	27	England	Right-Back	M
Aaron Wan-Bissaka	27	England	Right-Back	W
Aaron Zehnter	20	Germany	Left Midfield	FC

Football Data Analysis

26.01

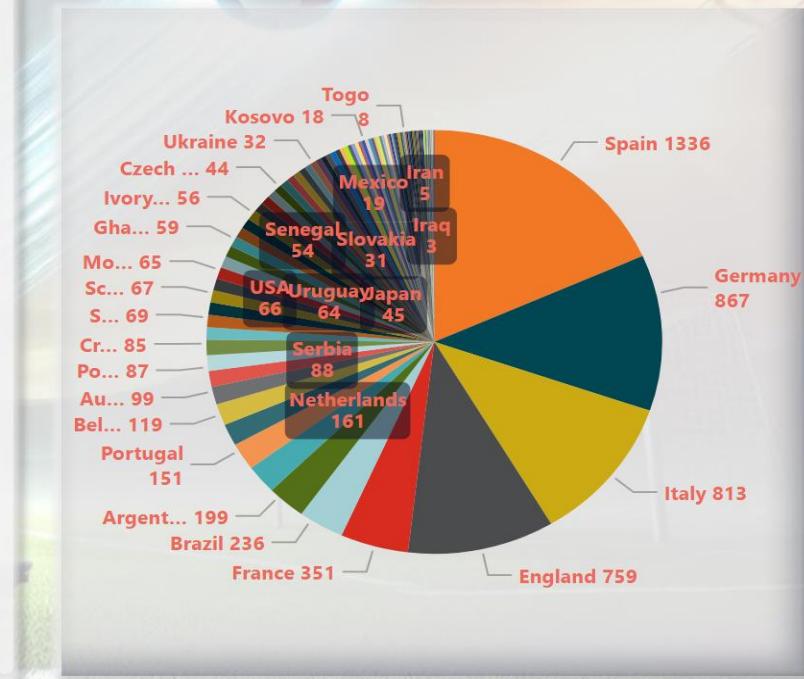
112

1336

Average of Age

No. of Nationality

Max Nationality



Football Data Analysis

Teams

- Select all
- Brentford FC
- Brighton & Hove Albion FC

Match Day

- 1
- 2
- 3

Season

- 2022
- 2023
- 2024

Competition Name

- Bundesliga
- Premier League
- Primera Division
- Serie A

Home Team Name	Score Home	Game Time	Score Away	Away Team Name
Chelsea FC	3	12/1/2024 1:30:00 PM	0	Aston Villa FC
Liverpool FC	2	12/1/2024 4:00:00 PM	0	Manchester City FC
Manchester United FC	4	12/1/2024 1:30:00 PM	0	Everton FC
Nottingham Forest FC	1	11/30/2024 3:00:00 PM	0	Ipswich Town FC
Brentford FC	4	11/30/2024 3:00:00 PM	1	Leicester City FC
Brighton & Hove Albion FC	1	11/29/2024 8:00:00 PM	1	Southampton FC
Crystal Palace FC	1	11/30/2024 3:00:00 PM	1	Newcastle United FC
Tottenham Hotspur FC	1	12/1/2024 1:30:00 PM	1	Fulham FC
Wolverhampton Wanderers FC	2	11/30/2024 3:00:00 PM	4	AFC Bournemouth
West Ham United FC	2	11/30/2024 5:30:00 PM	5	Arsenal FC

Service name ▾



Azure Databricks

kr151.45

Azure Data Factory v2

kr95.65

Virtual Machines

kr73.75

Storage

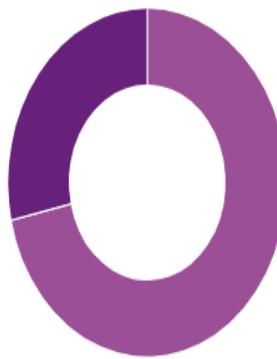
kr24.75

Bandwidth

kr0.07

View details

Location ▾



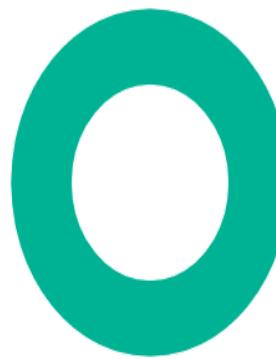
eu north

kr247.28

se central

kr98.38

Subscription ▾



Azure subscription 1

kr345.66

COST ANALYSIS

DEMO TIME

CONCLUSION: APPLYING SKILLS IN THE ETL PROJECT

- **DATABASE DESIGN AND SQL:** APPLIED ADVANCED QUERYING AND MODELING TECHNIQUES FOR EFFICIENT DATA STRUCTURING IN SQL SERVER.
- **PYTHON PROGRAMMING:** USED PYSPARK IN DATABRICKS FOR DATA CLEANING AND TRANSFORMATION.
- **DEVOPS:** AUTOMATED WORKFLOWS AND ENSURED RELIABILITY USING AZURE DATA FACTORY AND CLOUD TOOLS.
- **CLOUD SOLUTIONS:** UTILIZED AZURE PLATFORMS FOR DATA STORAGE, PROCESSING, AND MANAGEMENT.
- **ETL/ELT EXPERTISE:** DESIGNED A ROBUST PIPELINE FOR API INTEGRATION, DATA MIGRATION, AND TRANSFORMATION.
- **SYSTEM UNDERSTANDING AND BUSINESS ALIGNMENT:** DELIVERED ACTIONABLE INSIGHTS VIA A POWER BI DASHBOARD TAILORED TO BUSINESS NEEDS.
- THIS PROJECT DEMONSTRATES MY ABILITY TO INTEGRATE SKILLS FROM THE PROGRAM INTO SOLVING REAL-WORLD DATA CHALLENGES.

A black and white photograph of a soccer ball. A hand has written the words "MERCI POUR VOTRE" in cursive across the center of the ball. Below this, the word "ATTENTION" is printed in a large, bold, sans-serif font. At the bottom of the ball, the question "QUESTAINS ?" is written in a smaller, italicized serif font.

THANK YOU FOR YOU
ATTENTION

QUESTAINS ?