# DAMG 7275 : Advanced Data Management Systems
## Project P3 Submission

**Topic :** NBA Game Analytics
**Team : T**eam 6 – Azure SQL Multi-model
**Team Members :**
      Krishika Singh- 002194016
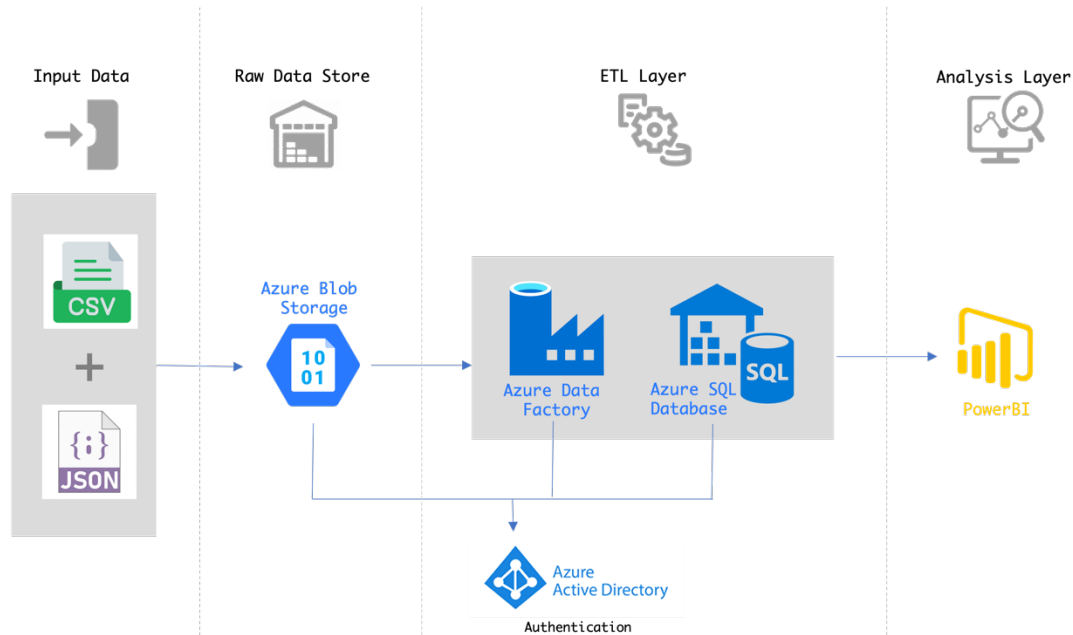      Mohammad Rafi Shaik- 001525707
      Shika Shyam- 002194543
      Shreya Soni – 002146758

---

## Implementation Process :
Our Project architecture is as shown below :



- The data we are using is available in the NBA dataset provided by Kaggle. The data is in the form of multiple csv files for different entities.
- Since these CSV files are not normalized and has a lot of redundant columns we will be needing extensive preprocessing before we can use the data for reporting.

Our Implementation Process consists of three phases :
## Step 1 : Raw Data Layer
- We are using Azure Blob Storage to store the csv files from Kaggle.
- We are also using Blob storage as the staging area for JSON documents.

## Step 2 : ETL Layer

- Azure Data Factory is our ETL tool. We are using Data Copy activity for some entities and full-fledged data flows for others.
- We have also enabled both schedule based and storage event-based triggers.

## Step 3 : Foundation Data Layer
- The sink of our ETL pipelines point to Azure SQL tables.
- We are storing data as tables, nodes and edges for graphs as well as document.

## Raw Data Store :

We have made use of Azure Blob Storage as shown below to store our CSVs. The source for all our ETLs will be pointing to this storage location.

nbadataset is the container where all our source csv files are available.



Additionally, we use documentdata container to store the staging json file that we create from the csv in order to store as document in foundation layer.



## ETL Layer:

Since we have multiple source csv files our first step in ETL Design was to create a **consolidated data dictionary** as shown below where we determine from which file the source data for each target column in a given table comes from.


Microsoft Excel
Worksheet

We created this for all columns in every single table in our Foundation layer.

| | A | B | C | D |
|---|---|---|---|---|
| 1 | Table Name : | **TeamDimension** | | |
| 2 | | | | |
| 3 | Table Column | Available in | Column Name | Comments |
| 4 | team_ID | team_details.csv | | |
| 5 | teamName | team_info_common.csv | CONCAT(team_city,team_name) | |
| 6 | City | team_details.csv | city | |
| 7 | State | team.csv | State | Available |
| 8 | Arena | team_details.csv | arena | |
| 9 | Head Coach | team_details.csv | head_coach | |
| 10 | Arena Capacity | team_details.csv | arena_capacity | |
| 11 | conference | team_info_common.csv | team_conference | |
| 12 | division | team_info_common.csv | team_division | |
| 13 | Year Found | team.csv | Year_Founded | New |
| 14 | Wins | team_info_common.csv | w | New |
| 15 | Losses | team_info_common.csv | l | New |
| 16 | PCT | team_info_common.csv | pct | In basketball, "PCT" stands for "Winning Percentage" or "Win-Loss Percentage".PCT = Wins / (Wins + Losses) |
| 17 | div_rank | team_info_common.csv | Division_rank | |
| 18 | | | | |
| 19 | | | | |
| 20 | | | | |
| 21 | | | | |
| 22 | | | | |
| 23 | | | | |
| 24 | | | | |
| 25 | | | | |
| 26 | | | | |
| 27 | | | | |

BoxScore | TeamDimension | PlayerDimension | GameDimension | Officials | Draft History | +

## ETL Pipeline in Azure Data Factory:

We have created the below ETL pipeline for our project. It loads each entity in the order of execution, and runs the required stored procedures in target after loading.
Due to the cost constraints that come with using heavy clusters for ADF pipeline executions, we have traded off on the parallel runs and decided to run each component of the ETL sequentially. In an ideal scenario with unlimited costs and unlimited compute resources available – we could possibly run most of our pipeline activities in parallel.
Another reason to run the jobs sequentially is due to the foreign key constraints between tables that require data to be available in some tables before it Is available in other tables.

## Pipeline execution triggers :

We have set up two types of triggers for our pipeline –
- There is a Scheduled trigger called Daily Run – which is scheduled to run everyday at 12AM UTC
- There is a Storage Events Trigger called BucketUpdateTrigger- which will run every time there is a new file added or updated in our Raw Data Storage Container

**Activity runs**

Pipeline run ID 1bb969e8-b6d5-4ae6-aeef-0453b72f8d40

All status ∨                                                                                            ⬇ Export to CSV

Showing 1 - 10 items

| Activity name ↑↓ | Status ↑↓ | Activity type ↑↓ | Run start ↑↓ | Duration ↑↓ | Log | Integration runtime ↑↓ | User properties ↑↓ | Run ID ↑↓ |
|---|---|---|---|---|---|---|---|---|
| LoadJSONToTable | ✓ Succeeded | Stored procedure | 4/6/2023, 4:07:50 AM | 00:00:03 | | AutoResolveIntegrationRunti | | 468c568e-38ef-41d1-k |
| DraftCombineStatsJSON | ✓ Succeeded | Copy data | 4/6/2023, 4:07:37 AM | 00:00:11 | | AutoResolveIntegrationRunti | | abc48077-158a-4f2f-a |
| InsertPlaysInTeam_Game | ✓ Succeeded | Stored procedure | 4/6/2023, 4:07:26 AM | 00:00:09 | | AutoResolveIntegrationRunti | | 80a93d88-6da4-4848- |
| InsertPlaysForPlayer_Team | ✓ Succeeded | Stored procedure | 4/6/2023, 4:07:20 AM | 00:00:04 | | AutoResolveIntegrationRunti | | a8759722-1e67-492c-l |
| BoxScore | ✓ Succeeded | Data flow | 4/6/2023, 4:06:51 AM | 00:00:27 | | AutoResolveIntegrationRunti | | f72b4cc4-91aa-44fd-9 |
| GameDimension | ✓ Succeeded | Data flow | 4/6/2023, 4:05:54 AM | 00:00:56 | | AutoResolveIntegrationRunti | | c00627fe-baea-4dce-k |
| TeamDimensionTable | ✓ Succeeded | Data flow | 4/6/2023, 4:05:18 AM | 00:00:35 | | AutoResolveIntegrationRunti | | a9b9f385-af98-436f-8 |
| PlayerDimensionTable | ✓ Succeeded | Data flow | 4/6/2023, 4:04:51 AM | 00:00:24 | | AutoResolveIntegrationRunti | | 97adf483-5ee5-4464-! |
| OfficialsTable | ✓ Succeeded | Data flow | 4/6/2023, 4:01:02 AM | 00:03:47 | | AutoResolveIntegrationRunti | | 0e195100-d331-40e3- |
| Draft_History Table | ✓ Succeeded | Copy data | 4/6/2023, 4:00:38 AM | 00:00:24 | | AutoResolveIntegrationRunti | | 849ac5c1-d618-412a-: |

The above image shows all steps of our ADF pipeline has run and completed successfully.

Our pipeline consists of three major types of ADF activities :

1. **DataCopy Activity :**
   i. We use Data Copy activity to populate DraftHistory table : Since this table requires a direct pull with the omission of a few columns from one single source csv file, we can leverage Copy Data Activity Here.

General    Source    **Sink**    Mapping    Settings    User properties

**Sink dataset** *                          DraftHistorySQLTable   ⌄        ✎ Open   + New        Learn more ↗

**Write behavior**                 ○ Insert   ◉ Upsert   ○ Stored procedure

**Use TempDB** ⓘ                   ☐

**Select user DB schema** ⓘ        dbo

**Key columns** ⓘ                  + New   |   🗑 Delete   ⟳ Refresh

                                   ☐

                                   ☐   ᴬⁿʸ  Person_id                    ⌄

                                   ☐   ᴬⁿʸ  season                       ⌄

                                   ☐   ᴬⁿʸ  team_id                      ⌄

**Bulk insert table lock** ⓘ       ○ Yes   ◉ No

ii.    We also use CopyData activity for converting the csv file into JSON format for our
       DocumentStore table – PlayerDraftStatistics



General    **Source**    Sink    Mapping    Settings    User properties

**Source dataset** *                    DraftStatsCSV   ⌄        ✎ Open   + New   👓 Preview data   Lea

**File path type**              ◉ File path in dataset   ○ Prefix   ○ Wildcard file path   ○ List of files ⓘ

Copy data

DraftCombineStatsJSON

General | Source | **Sink** | Mapping | Settings | User properties

| | |
|---|---|
| **Sink dataset *** | DraftStatsJSON  ✎ Open  + New  Lea |
| **Copy behavior** ⓘ | None |
| **Max concurrent connections** ⓘ | |
| **Block size (MB)** ⓘ | |
| **Metadata** ⓘ | + New |
| **File pattern** | Array of objects |

## 2. Data Flow Activity:

We have ADF Dataflow components for OfficialsTable PlayerDimensionTable, TeamDimension ,GameDimension and BoxScore.

We have made use of various mapping components available within AWS such as – Cast, Derived Columns, Lookups, Filters, Select, Aggregation component etc.

i.    Officials Table:



ii.    Player Table :

In the above Data Flow we are creating two sink datasets – PlayerDimension table and PlayerBridge Table. PlayerBridge table acts as an associative entity to link players to teams – we use this same information in order to define edges between Player and Team in our Graph model.

iii.    TeamDimension Table :



iv.    GameDimension Table:



v.    BoxScore Table :



3.  **Stored Procedure Execution:**

We use the Stored Procedure Execution component for executing three different stored procedures. We will explain the purpose of each Stored Proc later in this document :

i.    Usp_insert_plays_for Stored Proc



| General | Settings | User properties |

Linked service * ⓘ          [ 🗄 AzureSqlDatabaseConnection      ∨ ]   ⚡ Te

Stored procedure name *      [dbo].[usp_Insert_Plays_For]
                             ☑ Edit ⓘ

∨ Stored procedure parameters ⓘ

⤆ Import    ＋ New

ii.    Usp_insert_plays_in Stored Proc



| General | Settings | User properties |

Linked service * ⓘ          [ 🗄 AzureSqlDatabaseConnection      ∨ ]   ⚡ Test

Stored procedure name *      [dbo].[usp_Insert_Plays_In]
                             ☑ Edit ⓘ

∨ Stored procedure parameters ⓘ

⤆ Import    ＋ New

iii.   JsonImport Stored Proc

## Foundation Data Layer

Like described in the ERD, we have three types of Data Store within Azure SQL:

- Relational Data : Draft_History, Officials, BoxScore
- Document Data : PlayerDraftStatistics
- Graph Data : PlayerDimension, TeamDimension, GameDimension

<u>The Relational Tables</u> are created with the below DDLs and are loaded via ETL through Azure Data Factory as described above.

```sql
create table BoxScore
(
    Game_id                              int not null primary key,
    team_id_home                         int,
    team_id_away                         int,
    Minutes_played                       int,
    Points_hometeam                      int,
    Points_away                          int,
    field_goal_attempts_home             int,
    field_goal_Made_home                 int,
    field_goal_attempts_away             int,
    field_goal_Made_away                 int,
    three_point_field_goal_attempts_home int,
    three_point_field_goal_attempts_away int,
    three_point_field_goal_made_home     int,
    three_point_field_goal_made_away     int,
    free_throw_attempts_Home             int,
    free_throw_attempts_away             int,
    free_throw_made_Home                 int,
    free_throw_made_away                 int,
    reb_away                             int,
    reb_home                             int,
```

```
    assist_home                          int,
    assist_away                          int,
    Steal_home                           int,
    Steal_away                           int,
    Pf_away                              int,
    Pf_home                              int,
    Blocked_shots_away                   int,
    Blocked_shots_home                   int
)

create table Draft_history
(
    Person_id      int not null,
    Player_name    varchar(50),
    season         int not null,
    Round_number   int,
    round_pick     int,
    team_id        int not null,
    primary key (Person_id, season, team_id)
)

create table Official
(
    Official_id int not null,
    Name        varchar(50),
    Game_id     int not null,
    Jerseynum   int,
    primary key (Official_id, Game_id)
)
```

**The Document model** is created with the below DDL, and loaded using the stored proc below that pulls the JSON data from the Azure Blob Storage which is defined as a DATA SOURCE within our nba SQL Server

```
create table PlayerDraftStatistics
(
    JsonData nvarchar(max)
)
```
The below stored proc loads the PlayerDraftStatistics Document table. The Stored Proc is executed through the ADF pipeline.

```
CREATE PROCEDURE JsonImport
AS
BEGIN

TRUNCATE TABLE [dbo].[PlayerDraftStatistics]
BULK INSERT [dbo].[PlayerDraftStatistics]
```

```
FROM 'documentdata/draft_combine_stats.json'
WITH ( DATA_SOURCE = 'MyAzureBlobStorage');

END
```

The Graph Data Model is created using 3 different entities as shown in the Graph ERD-
PlayerDimension, TeamDimension and GameDimension. The DDLs are below:



```sql
create table GameDimention
(
    Game_ID                         int not null primary key,
    Date                            date,
    home_team_id                    int,
    away_team_id                    int,
    Season                          varchar(50),
    Score                           varchar(50)
) AS NODE


create table Player_Dimension
(
    Player_id                       int not null primary key,
    PlayerName                      varchar(50),
    Position                        varchar(50),
    Height                          varchar(10),
    Weight                          int,
    Birthdate                       date,
    College                         varchar(50),
    Country                         varchar(50),
    season_experience               int,
    Jersey                          int,
    from_year                       int,
    to_year                         int,
    nba_flag                        varchar(5)
) AS NODE
```

```sql
create table Team
(
    TeamID                              int not null primary key,
    TeamName                            varchar(255),
    City                                varchar(255),
    State                               varchar(255),
    Arena                               varchar(255),
    ArenaCapacity                       int,
    HeadCoach                           varchar(255),
    Conference                          varchar(255),
    TeamDivision                        varchar(255),
    YearFounded                         int,
    Wins                                int,
    Losses                              int,
    PCT                                 decimal(5, 3),
    Division_Rank                       int
) AS NODE

create table Player_Team_Bridge
(
    Player_id int,
    Team_id   int
)
```

The edge tables are also defined using the below DDL and two stored procs are leveraged to define the relationship between the nodes based on the source data

```sql
CREATE PROCEDURE usp_Insert_Plays_For
AS
BEGIN
    DROP TABLE plays_for
    CREATE TABLE plays_for AS EDGE
    INSERT INTO plays_for
    SELECT p.$node_id, t.$node_id
    FROM Player_dimension p
    INNER JOIN Player_Team_Bridge pt ON p.Player_id = pt.Player_id
    INNER JOIN Team t ON t.TeamID = pt.Team_id
END



CREATE PROCEDURE usp_Insert_Plays_In
AS
BEGIN
```

```
    DROP TABLE plays_in

    CREATE TABLE plays_in (H_A varchar(10)) AS EDGE

    INSERT INTO plays_in
    SELECT n1, n2, h_a
    FROM (
        SELECT t.$node_id n1, g_h.$node_id n2, 'Home' h_a
        FROM GameDimention g_h
        JOIN team t ON g_h.home_team_id = t.TeamID
        UNION
        SELECT t.$node_id, g_h.$node_id, 'Away'
        FROM GameDimention g_h
        JOIN team t ON g_h.away_team_id = t.TeamID
    ) a;
END
```

## Sample Queries:

We have executed the below sample queries on our Foundation data layer – specifically on the document and graph DB to ensure that data accuracy and integrity is in place.

1. Document DB query:

```
1    SELECT TOP(3) JSON_VALUE(JsonData, '$[0].player_id') AS Player_ID,
2            JSON_VALUE(JsonData, '$[0].position') AS PositionPlayed,
3            JSON_VALUE(JsonData, '$[0].weight') AS Weight
4    FROM PlayerDraftStatistics;
```

**Results**   Messages

🔍 Search to filter items...

| Player_ID | PositionPlayed | Weight |
|---|---|---|
| 2124 | PF-C | 271.0 |
| 12019 | SG-SF | 219.0 |
| 12020 | PF-C | 235.5 |

2. Graph DB query:

```
1    SELECT t.TeamName, g.Game_ID, g.Season, plays_in.H_A
2    FROM Team t, plays_in, GameDimention g
3    WHERE MATCH(t-(plays_in)->g)
4    and t.TeamName = 'Los Angeles Lakers'
5    and g.season = 2018
6    and plays_in.H_A = 'Away';
```

**Results**   Messages

🔍 Search to filter items...

| TeamName | Game_ID | Season | H_A |
|---|---|---|---|
| Los Angeles Lakers | 21800016 | 2018 | Away |
| Los Angeles Lakers | 21800060 | 2018 | Away |
| Los Angeles Lakers | 21800082 | 2018 | Away |
| Los Angeles Lakers | 21800093 | 2018 | Away |
| Los Angeles Lakers | 21800132 | 2018 | Away |

```
8     SELECT TOP (5) p.PlayerName, t.TeamName
9     FROM Player_dimension p, plays_for, Team t
10    WHERE MATCH(p-(plays_for)->t)
11    and t.TeamName = 'Los Angeles Lakers';
```

**Results**   Messages

🔍 Search to filter items...

| PlayerName | TeamName |
|---|---|
| Cliff Anderson | Los Angeles Lakers |
| D.J. Augustin | Los Angeles Lakers |
| Elgin Baylor | Los Angeles Lakers |
| Malik Beasley | Los Angeles Lakers |
| Mario Bennett | Los Angeles Lakers |

## Next Steps:

As a part of next steps, specifically for P4 submission, we aim to automate the raw data store – i.e. figure out a way for the CSV to be uploaded to Azure Blob Storage automatically as and when the data changes in actual source – ie. Kaggle.

We also plan on creating a few reporting based views in the foundation layer for easy access for non-tech folks who do not want to wait for the reporting layer data – just as it would be in a real-world scenario.

**\*\*\*\*\*\*\*\*\*\***