# Why I Moved My Power BI Solution to Microsoft Fabric

**Power BI Solution**

**Fabric System**

Get Current Season Data

Get Live Schedule

Bronze

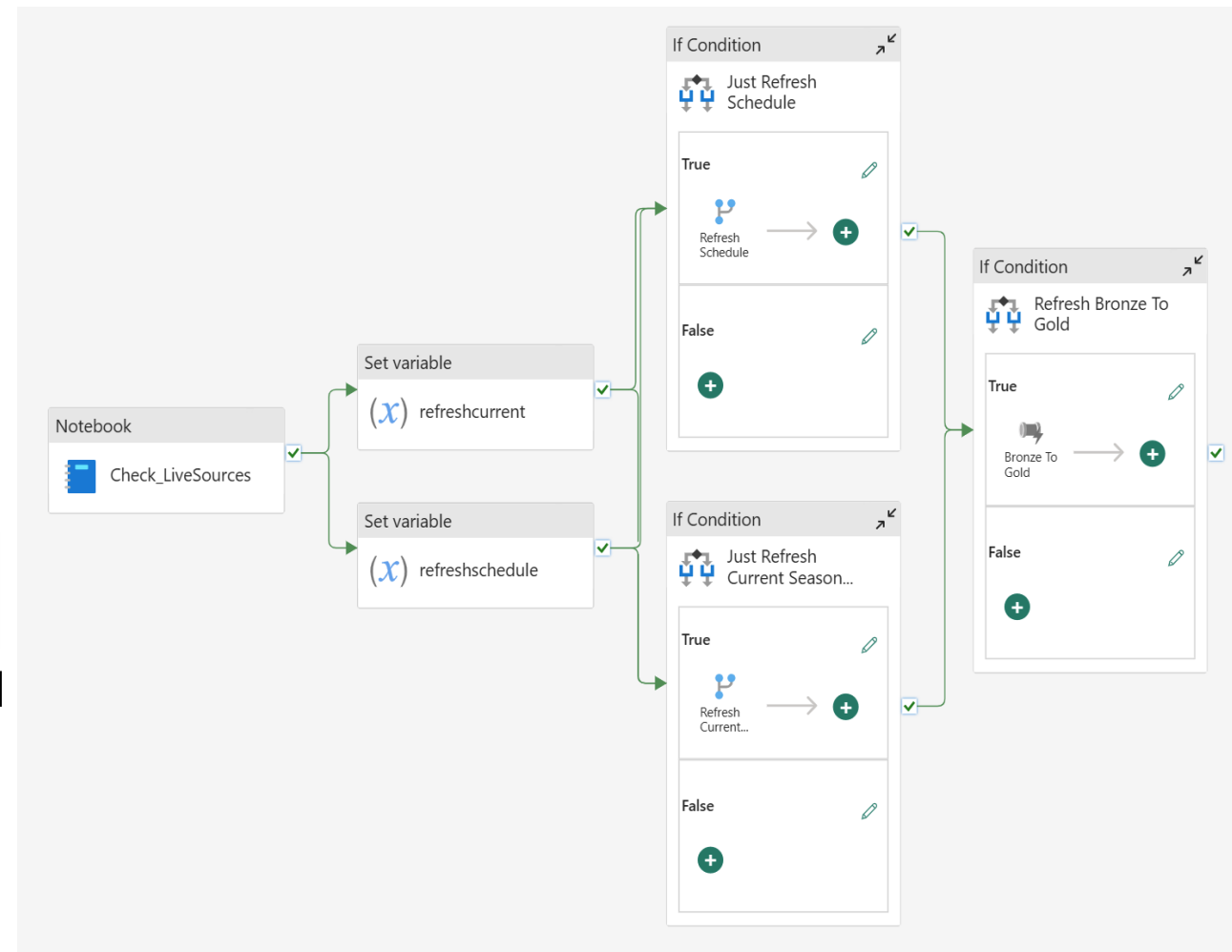Silver 1

Silver 2

Gold

# 1. Resilience & Robustness

I am ingesting data from two web sources.

**Key Problem:**
What if one or both data sources are unavailable due to downtime or maintenance?

**Solution:**
- I can use a **notebook** to return two **exit values** indicating whether each web source is accessible and providing valid data.
- If a web source fails validation, the exit value is "No." If it passes validation, the exit value is "Yes."
- Using **Fabric Pipelines**, these exit values are stored as **variables** and applied in "If Condition" activities to **orchestrate the pipeline**. This ensures that it skips refreshing any source that is down or providing invalid data.
- This ensures the pipeline remains **resilient** and continues processing available sources while defaulting any problematic ones.

Additionally, I enabled the **retry property** in the general settings of pipeline activities to handle temporary server errors (like HTTP 500). This allows the pipeline to automatically attempt to fetch data again if it encounters transient errors, enhancing overall **robustness**.

**Key Problem:**

Team Names are inconsistent between the two sources (i.e. one source could identify Manchester United as "Man United" whilst the other as "Man Utd").

Furthermore, we need safeguards to prevent unexpected or unrecognized team names from appearing

**Solution:**

We can use **notebooks** to **validate** and standardize team names between the two sources.

**(1)**: Create a mapping dictionary where each standardized team name is the key, and the value is a list of possible variations.

```
1   # Define the team mapping dictionary
2   team_mapping = {
3       "Arsenal": ["Arsenal"],
4       "Aston Villa": ["Aston Villa"],
5       "Barnsley": ["Barnsley"],
6       "Birmingham City": ["Birmingham", "Birmingham City"],
7       "Blackburn Rovers": ["Blackburn", "Blackburn Rovers" ],
8       "Blackpool": ["Blackpool"],
9       "Bolton": ["Bolton"],
10      "Bournemouth": ["Bournemouth"],
11      "Bradford": ["Bradford"],
12      "Brentford": ["Brentford"],
13      "Brighton and Hove Albion": ["Brighton", "Brighton and Hove Albion"],
14      "Burnley": ["Burnley"],
15      "Cardiff": ["Cardiff"],
16      "Charlton": ["Charlton"],
```

```
52      "Wimbledon": ["Wimbledon"],
53      "Wolverhampton Wanderers": ["Wolves", "Wolverhampton Wanderers"]
54  }
55
56  # Convert team mapping to a list of tuples (standard_name, variation)
57  team_mapping_data = [(standard, variation) for standard, variations in team_mapping.items() for variation in variations]
58  team_mapping_df = spark.createDataFrame(team_mapping_data, ["Standard_Team", "Variation"])
```
✓ 12 sec - Session ready in 9 sec 364 ms. Command executed in 3 sec 406 ms by Krishan Patel on 3:58:06 PM, 10/05/24

**(2)**: Implement a **validation step** to check if the team names from both sources match the standardized dictionary. Flag any unrecognized names for review..

```
1   from pyspark.sql import functions as F
2
3   # Identify teams in Silver 1 Data Sources that are not found in team_mapping.
4   unexpected_teams = silver_1_teams.alias("dt").join(
5       team_mapping_df.alias("tm"),
6       F.col("dt.Home_Team") == F.col("tm.Variation"),
7       "left_anti"
8   )
9
10  # Show teams not expected in data source (if any)
11  display(unexpected_teams)
```
[7]  ✓  4 sec - Command executed in 3 sec 601 ms by Krishan Patel on 4:21:10 PM, 10/05/24

> ⊟ Spark jobs (10 of 10 succeeded)   ▥ Resources   ▤ Log

No data available

```
1   if unexpected_teams.isEmpty():
2       # Code to proceed
3       print("Teams are as expected. Proceeding...")
4   else:
5       raise ValueError("Unexpected teams appear in Silver 1. Please check your data.")
```
[8]  ✓  2 sec - Command executed in 1 sec 621 ms by Krishan Patel on 4:21:35 PM, 10/05/24

> ⊟ Spark jobs (10 of 10 succeeded)   ▥ Resources

Teams are as expected. Proceeding...

# 2. Validation & Quality Checks

**(3)**: If all teams are recognized, then team names in sources from the Silver 1 Lakehouse are standardized according to the mapping dictionary. If any team is unrecognized, the process is halted, and the issue is flagged for investigation.

```python
1    #Function To Clean Team Names
2
3    from pyspark.sql import functions as F
4    from pyspark.sql.types import StringType
5
6    # Broadcasting the team mapping dictionary for efficiency
7    team_mapping_broadcast = spark.sparkContext.broadcast(team_mapping)
8
9    # Define the function to map team variations to standardized names and raise an error if not found
10   def map_team_name(team):
11       mapping = team_mapping_broadcast.value
12       for standard, variations in mapping.items():
13           if team in variations:
14               return standard
15       # Raise an error if the team is not found
16       raise ValueError(f"Team name '{team}' not found in team mapping")
17
18   # Register the function as a UDF (user-defined function) in Spark
19   map_team_udf = F.udf(map_team_name, StringType())
```
✓ <1 sec - Command executed in 342 ms by Krishan Patel on 9:15:05 PM, 10/08/24

```python
1    # Transform Team Columns
2    df_cleaned_schedule = df_schedule.withColumn("Home_Team", map_team_udf(F.col("Home_Team"))) \
3                          .withColumn("Away_Team", map_team_udf(F.col("Away_Team")))
4
5    df_cleaned_historical = df_historical.withColumn("Home_Team", map_team_udf(F.col("Home_Team"))) \
6                          .withColumn("Away_Team", map_team_udf(F.col("Away_Team")))
7
8    df_cleaned_current = df_current.withColumn("Home_Team", map_team_udf(F.col("Home_Team"))) \
9                          .withColumn("Away_Team", map_team_udf(F.col("Away_Team")))
10
11   df_cleaned_penalties = df_penalties.withColumn("Team", map_team_udf(F.col("Team")))
```
✓ <1 sec - Command executed in 246 ms by Krishan Patel on 9:15:06 PM, 10/08/24
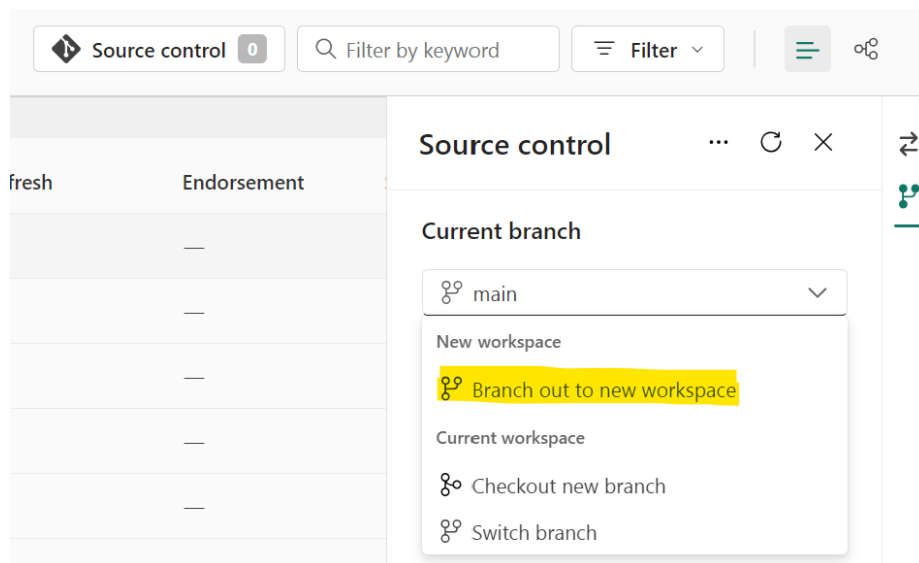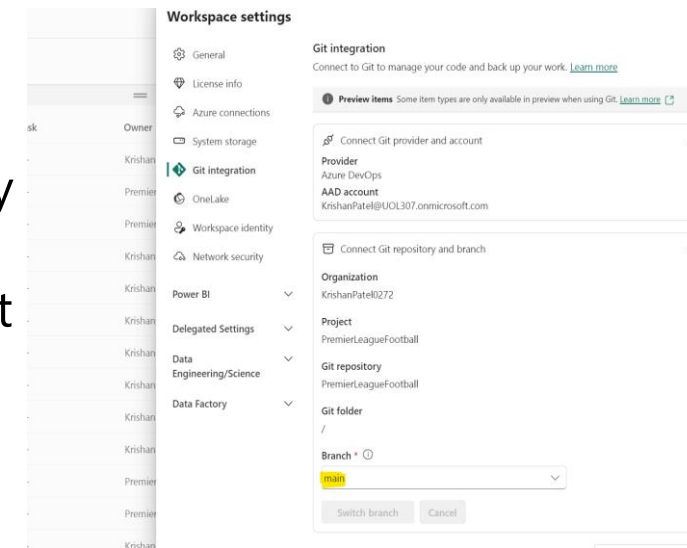
# 3. Enhanced Collaboration

**Key Problem:**

Without an integrated platform, collaboration on an analytic solution becomes challenging, leading to difficulties in managing version control, tracking changes, and maintaining code quality.
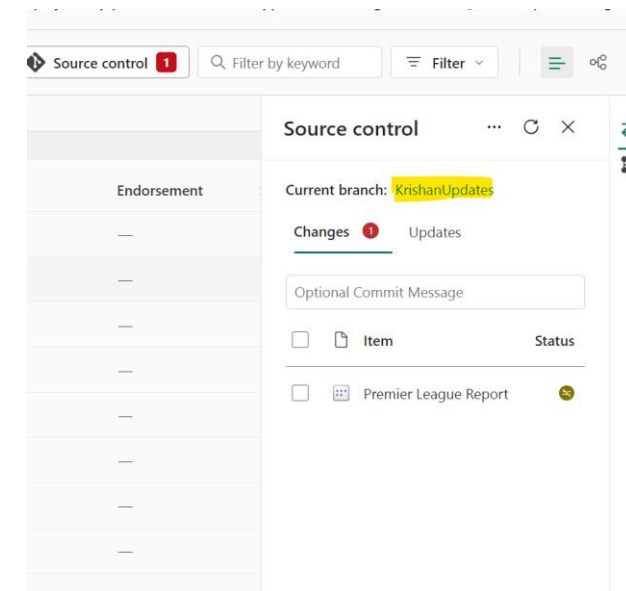
**Solution:**

I enabled **Git integration** using Microsoft Fabric by connecting to an **Azure DevOps repository**. This setup enhanced collaboration on the analytic solution by streamlining version control, facilitating efficient change tracking, and enabling code reviews, ensuring consistent code quality and effective teamwork throughout the project.

An example of how Git Integration works in Fabric:

**1:**

To prevent interfering with others, each team member should select "Branch out to new workspace" to work in an isolated workspace until ready to share changes.

**2:**

Once I'm satisfied with my changes, I can commit them to a separate branch, "KrishanUpdates."

Continued: An example of how Git Integration works in Fabric:



**3:**

In the Git repo, I created a **pull request** to merge the KrishanUpdates branch with the main branch.



**4:**

A **pull request** enables code reviews, protects the main branch with reviewer approval, promotes collaboration, and ensures traceability, allowing safe integration of changes without disrupting the shared workspace.

**5:**

After merging my changes, I sync my workspace with the main branch to stay updated with the team's progress.

# 4. Enhanced Analytical Capabilities

Power BI is great for visualization and reporting, but by adopting Microsoft Fabrics, we can leverage:
- Notebooks for interactive, explorative analysis, which can extend to using advanced statistical techniques
- T-SQL (Warehouse) offers more flexibility and control for complex data analysis and exploration.

## Querying Gold-Layer Warehouse using T-SQL:

## Visualize data in a Spark notebook