

# IPL Cricket Match Data Integration

## Introduction

The Indian Premier League (IPL) is a premier Twenty20 cricket league, celebrated for its thrilling matches and vast global appeal. Its franchise-based model and remarkable talent not only shape modern cricket but also revolutionize its economic landscape, inspiring a new generation of players and fans alike.

This project aims to create an ETL pipeline using SQL Server Integration Services (SSIS) to structure IPL match data for analysis. The pipeline will support informed decision-making and enhance stakeholder engagement, leading to improved strategic insights and operational outcomes.



Submitted By  
Dhruvitsinh Rathod

# Table of Contents

1.	<a href="#">Introduction.....</a>	1
2.	<a href="#">Project Requirements.....</a>	3
3.	<a href="#">Project Overview.....</a>	4
4.	<a href="#">Data Model Design.....</a>	5
5.	<a href="#">Relationships.....</a>	7
6.	<a href="#">ETL Pipeline Architecture.....</a>	10
7.	<a href="#">ETL Process Stages.....</a>	13
8.	<a href="#">ETL Process Steps.....</a>	14
9.	<a href="#">Logging Implementation.....</a>	15
10.	<a href="#">Incremental loading.....</a>	16
11.	<a href="#">Scheduling SSIS Package.....</a>	17
12.	<a href="#">Challenges and Solution.....</a>	19

## Project Requirements

The objective of this project is to develop a robust ETL pipeline using SQL Server Integration Services (SSIS) for the integration and analysis of IPL cricket match data. The pipeline will transform raw ball-by-ball data into an organized and analysis-ready format to support stakeholders in making informed decisions.

### Key Features to Implement

1. **Data Validation**
  - Ensure the accuracy, completeness, and reliability of the incoming data.
2. **Bad Data Handling**
  - Identify and segregate invalid or incomplete records for analysis and resolution.
3. **Data Transformation**
  - Apply business rules to convert raw data into a structured and meaningful format.
4. **Incremental Loading**
  - Enable efficient updates by loading only new or modified data into the target system.
5. **Error and Success Handling**
  - Implement mechanisms to capture and log errors, as well as record successful operations.
6. **Logging**
  - Maintain a detailed log of all ETL processes for monitoring and troubleshooting.
7. **Scheduling**
  - Automate the ETL pipeline to run at predefined intervals for timely data availability.

### Development Guidelines

- Adhere to SSIS development best practices.
- Follow consistent naming conventions for components and processes.
- Optimize the use of SSIS components for performance and scalability.

### Deliverables

1. **SSIS Packages:**
  - Fully functional ETL packages for data extraction, transformation, and loading.
2. **Relational Database Schema:**
  - A well-defined schema representing the IPL cricket match data model.

### Expected Outcomes

The ETL pipeline will provide stakeholders with timely, accurate, and actionable insights by leveraging IPL data, facilitating enhanced decision-making across the cricket ecosystem.

## Project Overview

The IPL Cricket Match Data Integration project leverages SSIS to build an ETL pipeline for Processing and analyzing. ball-by-ball match data. It focuses on transforming raw data into a structured format to support performance evaluation, strategy development, and informed decision-making.

## Scope

- Design a relational data model for IPL match data.
- Develop ETL processes with features like data validation, incremental loading, and error handling.
- Deliver a fully functional SSIS package and database schema.

## Significance

This project showcases how ETL pipelines can efficiently handle large-scale data to deliver actionable insights, enhancing IPL's strategic and analytical capabilities.



## Data Model Design

The data model for this project is designed to effectively store and manage IPL cricket match data. It represents a star schema, optimized for analytical querying and business intelligence use cases. The model consists of fact and dimension tables, ensuring both data normalization and efficient querying for insightful analysis.

### Key Components of the Data Model:

#### 1. Fact Tables

- **fact\_Ball\_by\_Ball**: Captures granular ball-by-ball details, such as the bowler, batsman, ball outcome, and team performance.
- **fact\_Match**: Contains aggregated match-level information, including match results, winners, and key performance indicators.
- **fact\_Batsman\_Scored**: Tracks runs scored by batsmen at the granular over-ball level.
- **fact\_Player\_Match**: Highlights player-specific contributions for individual matches.
- **fact\_Wicket\_Taken**: Logs details of dismissals, including player out, bowler, and dismissal type.

#### 2. Dimension Tables

- **dim\_Player**: Stores player details, including attributes like batting hand, bowling skill, and nationality.
- **dim\_Team**: Contains team names and related details.
- **dim\_Venue**: Represents match venues with attributes such as city and country.
- **dim\_Bowling\_Style & dim\_Batting\_Style**: Categorize players by their respective styles.
- **dim\_Out\_Type**: Defines various dismissal methods.
- **dim\_City and dim\_Country**: Geographic information for matches and venues.
- **dim\_Toss\_Decision**: Captures toss-related details such as the choice to bat or field first.

#### 3. Relationships and Integrity

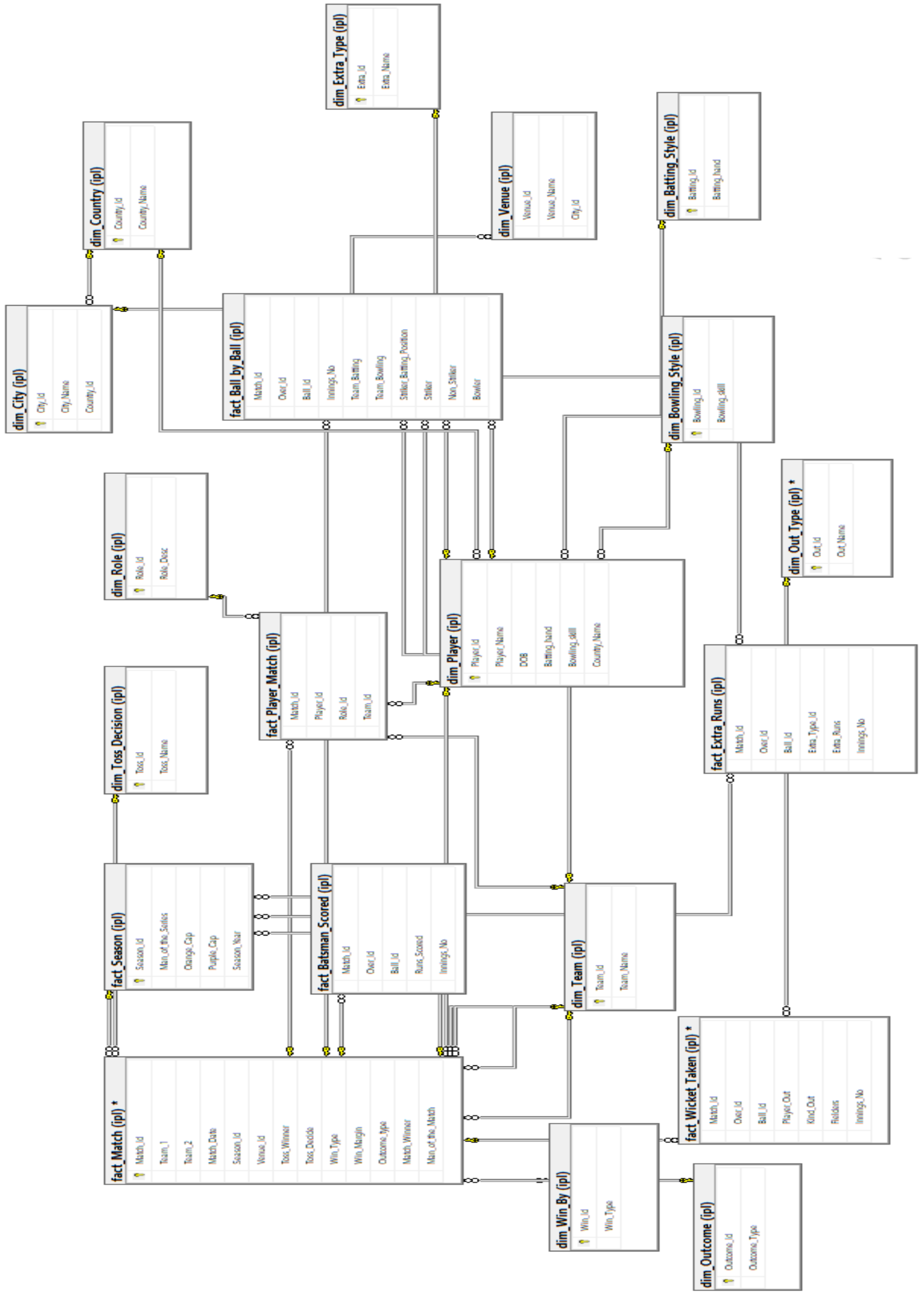
The model establishes relationships between fact and dimension tables using primary and foreign keys, ensuring data integrity and enabling comprehensive analytical queries.

### Highlights

- **Scalability**: The model accommodates additional seasons, teams, and players seamlessly.
- **Granularity**: Offers detailed insights from ball-level data to match-level summaries.
- **Flexibility**: Supports diverse analytical use cases, such as player performance evaluation, team comparisons, and venue-based analysis.

### Conclusion

This data model forms the backbone of the ETL pipeline, ensuring efficient storage, retrieval, and processing of IPL match data. It aligns with industry standards to support advanced analytics and reporting needs.



## Relationships

The IPL Cricket Match Database is structured to maintain data integrity and establish clear relationships between entities using **Primary Keys (PK)** and **Foreign Keys (FK)**. Below is an overview of the key relationships:

1. **Primary Keys (PK):**

Each table has a unique identifier column (e.g., Match\_Id, Player\_Id, Team\_Id) to ensure every record is distinct and easily accessible.

2. **Foreign Keys (FK):**

Foreign keys establish relationships between tables, allowing data to be connected logically. For example:

- Match.Team\_1 and Match.Team\_2 reference Team.Team\_Id.
- Match.Venue\_Id references Venue.Venue\_Id.
- Ball\_by\_Ball.Match\_Id references Match.Match\_Id.

3. **One-to-Many Relationships:**

Commonly seen in the schema, where one entity (e.g., a match) links to multiple related records (e.g., ball-by-ball data or player performances).

4. **Many-to-Many Relationships:**

These are implemented using bridge tables, such as Player\_Match, which connects Player and Match to record individual player participation and roles in matches.

5. **Referential Integrity:**

Foreign keys enforce consistency, ensuring that related records exist before being referenced (e.g., a match cannot reference a non-existent team).

Table	Columns	Data Type	Key (Reference: Tablename .TableColumn)
Match	Match_Id	int	Primary Key
	Team_1	int	Foreign key (Team.Team_Id)
	Team_2	int	Foreign key (Team.Team_Id)
	Match_Date	Date	
	Season_Id	int	Foreign key (Season.Season_Id)
	Venue_Id	int	Foreign key (Venue.Venue_Id)
	Toss_Winner	int	Foreign key (Team.Team_Id)
	Toss_Decide	int	Foreign key (Toss_Decision.Toss_Id)
	Win_Type	int	Foreign key (Win_By.Win_Id)
	Win_Margin	int	
	Outcome_type	int	Foreign key (Outcome.Outcome_Id)
	Match_Winner	int	Foreign key (Win_By.Win_Id)
	Man_of_the_Match	int	Foreign key (Player.Player_Id)
Ball_by_Ball	Match_Id	int	Foreign key (Match.Match_Id)

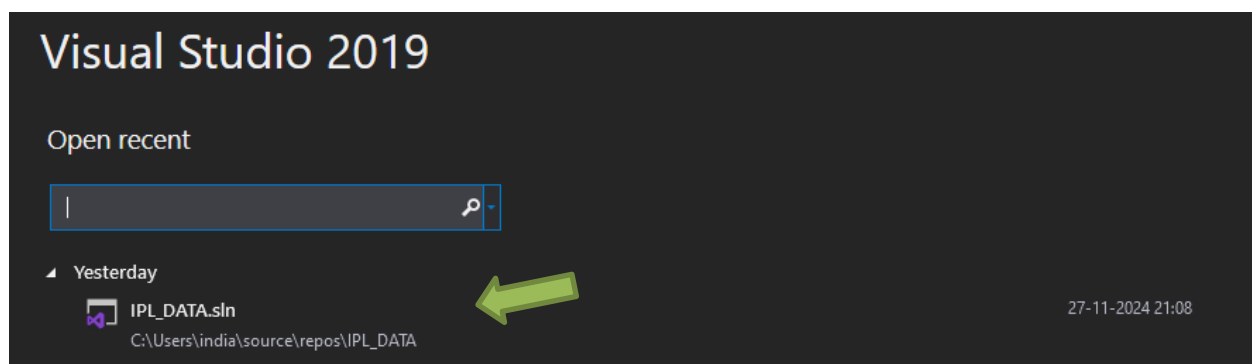
	Over_Id	int	
	Ball_Id	int	
	Innings_No	int	
	Team_Batting	int	Foreign key (Team.Team_Id )
	Team_Bowling	int	Foreign key (Team.Team_Id )
	Striker_Batting_Position	int	
	Striker	int	Foreign key (Player.Player_Id)
	Non_Striker	int	Foreign key (Player.Player_Id)
	Bowler	int	Foreign key (Player.Player_Id)
Batsman_Scored	Match_Id	int	Foreign key (Match.Match_Id)
	Over_Id	int	
	Ball_Id	int	
	Runs_Scored	int	
	Innings_No	int	
Extra_Runs	Match_Id	int	Foreign key (Match.Match_Id)
	Over_Id	int	
	Ball_Id	int	
	Extra_Type_Id	int	
	Extra_Runs	int	
	Innings_No	int	
Player_Match	Match_Id	int	Foreign key (Match.Match_Id)
	Player_Id	int	Foreign key (Player.Player_Id)
	Role_Id	int	Foreign key (Role.Role_Id)
	Team_Id	int	Foreign key (Team.Team_Id)
Player	Player_Id	int	Primary Key
	Player_Name	varchar	
	DOB	Date	
	Batting_hand	int	Foreign key (Batting_Style.Batting_Id)
	Bowling_skill	int	Foreign key (Bowling_Style.Bowling_Id)
	Country_Name	int	Foreign key (Country.Country_Id)
Batting_Style	Batting_Id	int	Primary Key
	Batting_hand	varchar	
Bowling_Style	Bowling_Id	int	Primary Key
	Bowling_skill	varchar	
Country	Country_Id	int	Primary Key
	Country_Name	varchar	
City	City_Id	int	Primary Key
	City_Name	varchar	
	Country_id	int	Foreign key (Country.Country_Id)
Venue	Venue_Id	int	Primary Key
	Venue_Name	varchar	



	City_Id	int	Foreign key (City.City_Id)
Out_Type	Out_Id	int	Primary Key
	Out_Name	varchar	
Outcome	Outcome_Id	int	Primary Key
	Outcome_Type	varchar	
Role	Role_Id	int	Primary Key
	Role_Desc	varchar	
Season	Season_Id	int	Primary Key
	Man_of_the_Series	int	Foreign key (Player.Player_Id)
	Orange_Cap	int	Foreign key (Player.Player_Id)
	Purple_Cap	int	Foreign key (Player.Player_Id)
	Season_Year	int	
Team	Team_Id	int	Primary Key
	Team_Name	varchar	
Toss_Decision	Toss_Id	int	Primary Key
	Toss_Name	varchar	
Extra_Type	Extra_Id	int	Primary Key
	Extra_Name	varchar	
Wicket_Taken	Match_Id	int	Foreign key (Match.Match_Id)
	Over_Id	int	
	Ball_Id	int	
	Player_Out	int	Foreign key (Player.Player_Id)
	Kind_Out	int	Foreign key (Out_Type.Out_Id)
	Fielders	int	Foreign key (Player.Player_Id)
	Innings_No	int	
Win_By	Win_Id	int	Primary Key
	Win_Type	varchar	

## ETL Pipeline Architecture

The **IPL Cricket Match Data Integration ETL pipeline** was designed and implemented in **Visual Studio 2019** using SQL Server Integration Services (SSIS). This pipeline automates the process of extracting raw IPL data, transforming it into a structured format, and loading it into a relational database for further analysis. The implementation started with the creation of an SSIS project, as shown below.



### Creating the Project

The first step was setting up the SSIS environment by creating a new solution in Visual Studio 2019:

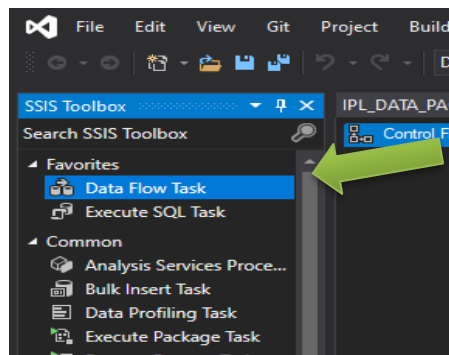
**Solution Name:** IPL\_DATA.sln

**Location:** Defined in a local repository to organize project files.

**Purpose:** This solution serves as the foundation for the ETL pipeline, integrating all SSIS packages and configurations.

### Adding a Data Flow Task to Control Flow

The **Data Flow Task** is a crucial component in the ETL pipeline, enabling the extraction, transformation, and loading of data. In the Control Flow, this task sets the foundation for processing IPL match data.



### Key Actions:

1. Added the **Data Flow Task** from the SSIS Toolbox to the Control Flow.
2. Renamed the task for clarity.
3. Configured the task to handle source, transformation, and target definitions in the Data Flow tab.

### Configuring Data Flow Tasks for Tables

In the ETL pipeline, individual **Data Flow Tasks** were created for each table to ensure a structured and efficient data integration process. The tasks were divided into two distinct stages:

#### Loading Dimension Tables:

Data Flow Tasks for all dimension tables were configured first, as dimension tables provide contextual information for the fact tables.

Each task was appropriately named to maintain clarity and align with best practices.

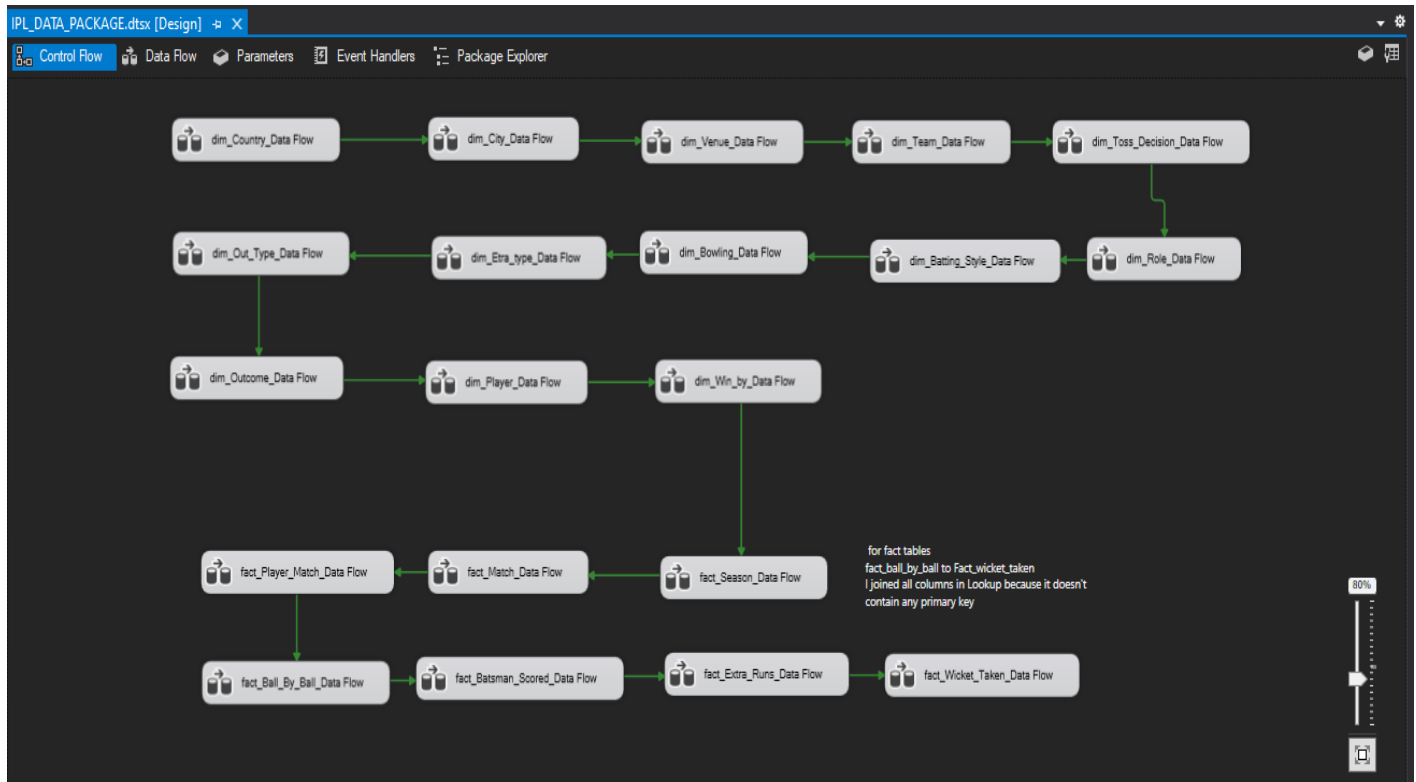
#### Loading Fact Tables:

Following the dimensions, Data Flow Tasks for fact tables were added. These tasks handle larger datasets with transactional information, such as match details, ball-by-ball data, and player statistics.

Fact table tasks were also clearly labeled.

#### Purpose and Organization

This sequential approach ensures that dimension data is loaded before fact data, maintaining referential integrity between tables. Each Data Flow Task is modular, promoting better management, debugging, and scalability of the pipeline.



## ETL Pipeline Architecture

### Pipeline Overview

- The ETL pipeline will outline the process of moving data from raw source files to the target database.
- The ETL pipeline is designed to automate the flow of IPL match data from raw source files into a structured relational database. It uses SSIS for data extraction, transformation, and loading. The key stages of the pipeline ensure the data is validated, transformed according to business rules, and loaded efficiently into the target schema, which supports advanced reporting and analytics.

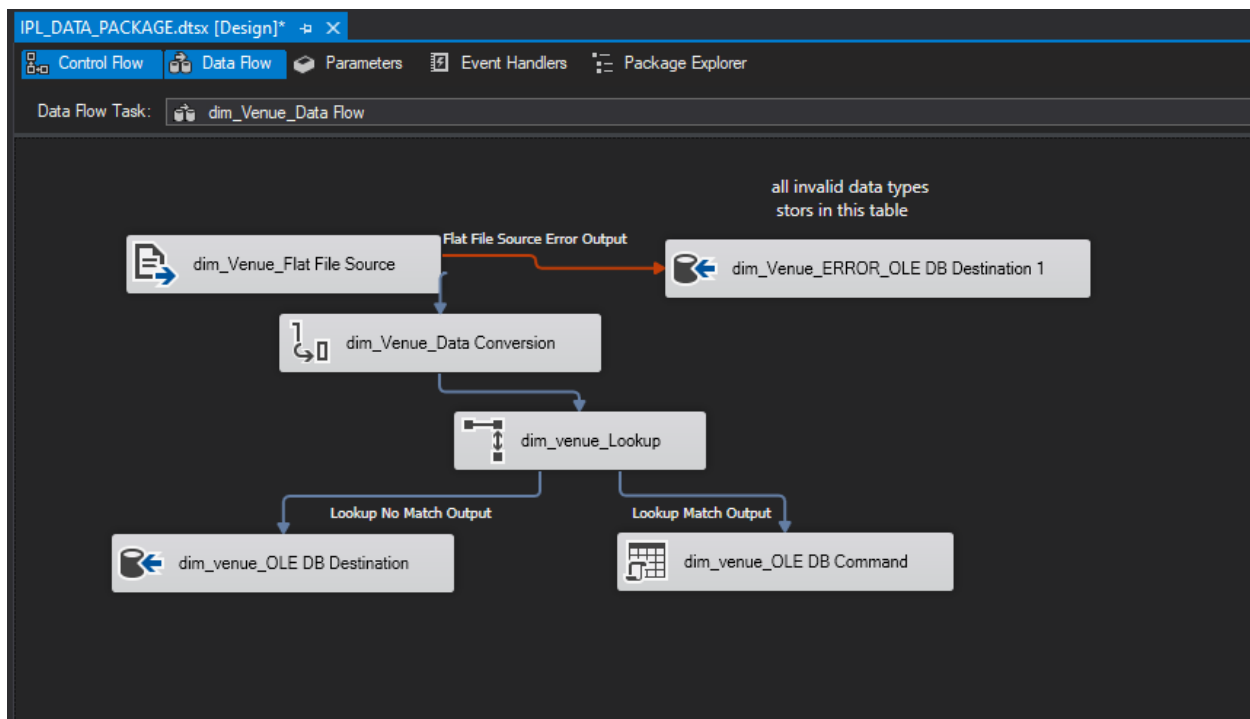
### Data Flow Description

This section explains how data moves through each stage in the ETL process.

1. **Extract:** Raw data is extracted from CSV files containing match, player, and team information.
2. **Validate:** The extracted data is validated to check for completeness, accuracy, and format correctness.
3. **Transform:** The data is transformed by applying business rules, converting data types, and performing necessary calculations.
4. **Load:** The cleaned and transformed data is loaded into the target database, using incremental loading techniques to ensure only new or updated records are processed.

## ETL Process Stages

This section will detail each of the ETL process stages—**Data Extraction**, **Data Validation**, **Bad Data Handling**, **Data Transformation**, **Incremental Loading**, and **Data Loading**—with an explanation of each stage and examples where necessary.



## ETL Process Steps

The ETL (Extract, Transform, Load) process for the IPL Cricket Match Database is designed to ensure efficient, accurate, and scalable data integration. Below is a professional overview of its key components:

1. Data Extraction:

Data is sourced from structured files (e.g., CSV) and external databases. This process extracts raw data about matches, teams, players, and match events into the SSIS workflow using the Flat File Source component, preparing it for transformation and validation.

2. Data Validation:

Validates data to ensure accuracy, completeness, and conformity with predefined rules. It includes checks for missing values, data type mismatches, and reference integrity to confirm that all foreign key dependencies are valid.

3. Bad Data Handling:

Non-compliant data is managed through rejection or rerouting to a designated "bad data" repository. This enables detailed review and corrective action without interrupting the ETL workflow.

4. Data Transformation:

Applies business logic to the extracted data, converting it into a format suitable for analysis and reporting. Transformations include type conversions, calculations (e.g., aggregations), and data enrichment to enhance usability.

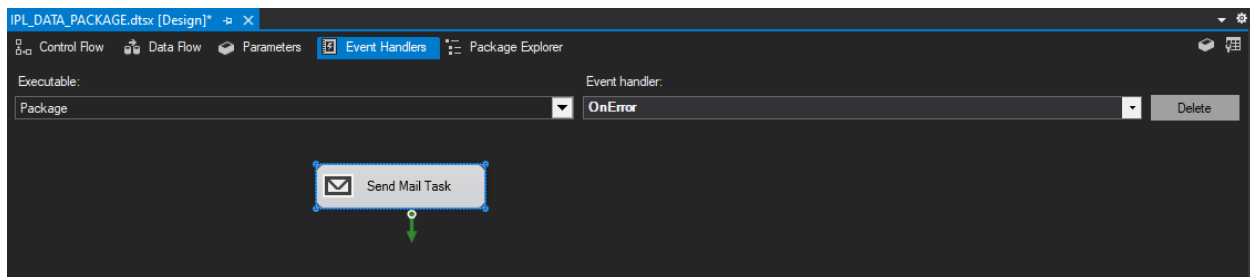
5. Data Loading:

The final stage involves inserting or updating data in the target database tables using the OLE DB Destination component. The process ensures data integrity while maintaining synchronization with the source systems.

6. Error and Success Handling:

During the ETL process, robust logging mechanisms are implemented to track errors and successful executions:

**Error Handling:** Logs issues such as data validation failures, transformation errors, or database connection problems for troubleshooting and resolution.



**Success Tracking:** Logs successful completion of each ETL stage, enabling process monitoring and ensuring smooth tracking of the workflow

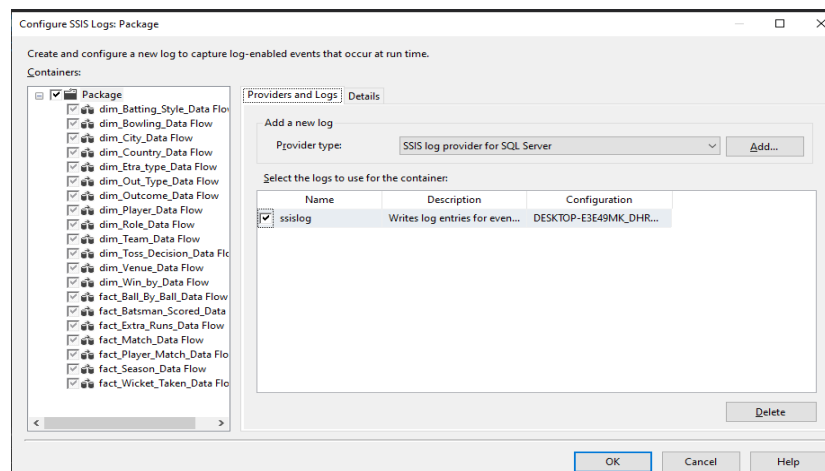
## Logging Implementation

### Objective

To track key events during ETL execution, ensuring process transparency, quick issue identification, and data integrity.

### Tracked Events

1. **OnError:** Logs error details like description, code, and source task for troubleshooting.
2. **OnPostExecute:** Records task completion status, name, and execution time.
3. **OnPostValidate:** Ensures configurations are validated before execution.
4. **OnPreExecute:** Captures task initiation details, including start time and context.
5. **OnQueryCancel:** Tracks user or system-initiated query cancellations.
6. **OnTaskFailed:** Logs failed tasks with reasons and timestamps.
7. **OnWarning:** Captures warnings for potential issues without halting execution.



# Incremental loading

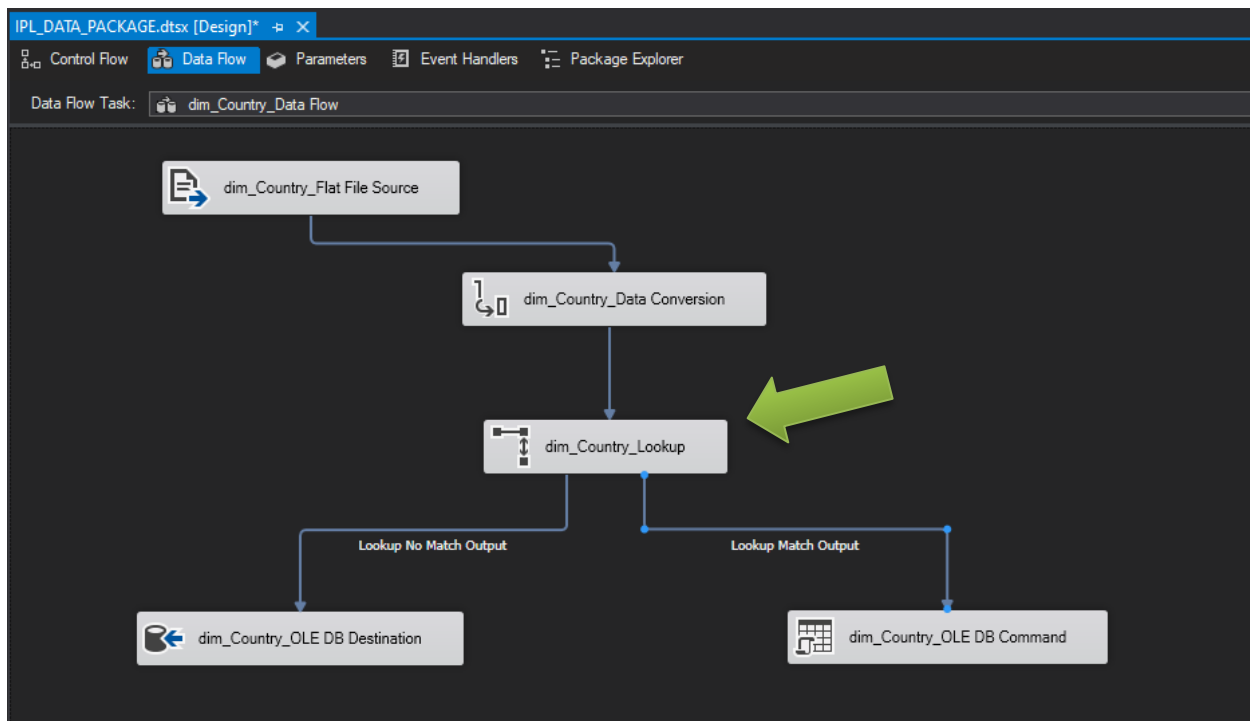
## Lookup Configuration

### 1. No Match Output

- **Purpose:** Identifies records where incoming data does not match the reference data.
- **Configuration:**
  - If no match is found, the data is directed to a destination table using the **OLE DB Destination** component.
  - This process ensures that new records are inserted into the database.

### 2. Match Output

- **Purpose:** Identifies records where incoming data matches the reference data based on specified keys.
- **Configuration:**
  - If a match is found, the data is routed to an **OLE DB Command**.
  - The OLE DB Command executes a **stored procedure** to update existing records in the database.
  - This ensures that any updates to matched records are handled dynamically and efficiently.





## stored procedure

The stored procedure used in the OLE DB Command is responsible for updating the matched records in the database, ensuring that existing data is refreshed with the latest information from the incoming data stream. This process optimizes data maintenance and minimizes manual updates.

```
SQLQuery1.sql - DE...3E49MK\india (58)) *  X
1  go
2  CREATE PROCEDURE spUpdateIplCityById1
3      @city_id INT,                -- Primary key or unique identifier for the city
4      @new_city_name VARCHAR(255), -- New name to update
5      @country_id int
6  AS
7  BEGIN
8      SET NOCOUNT ON;
9
10     -- Update the city name based on the city_id
11     UPDATE [ipl].[dim_City]
12     SET city_name = @new_city_name,
13     Country_Id = @country_id
14     WHERE city_id = @city_id;
15 END;
16 GO
17
```

## Scheduling in SQL Server: Automating the IPL Data Process

To ensure a streamlined and automated workflow, the **IPL Data Schedule** job was created and configured in SQL Server Agent. This job leverages SQL Server's scheduling capabilities to execute critical tasks at predefined intervals, minimizing manual intervention and ensuring consistency in data handling.

### Key Details:

1. **Job Name:** IPL\_DATA\_SCHEDULE
2. **Purpose:**

- Automates the extraction, transformation, and loading (ETL) processes related to IPL data.
  - Ensures timely and accurate updates to the database.
3. **Components:**
- **SQL Server Agent:** Handles the scheduling and execution of the job.
  - **SSIS Package (IPL\_DATA):** Executes the defined logic for data integration and processing.

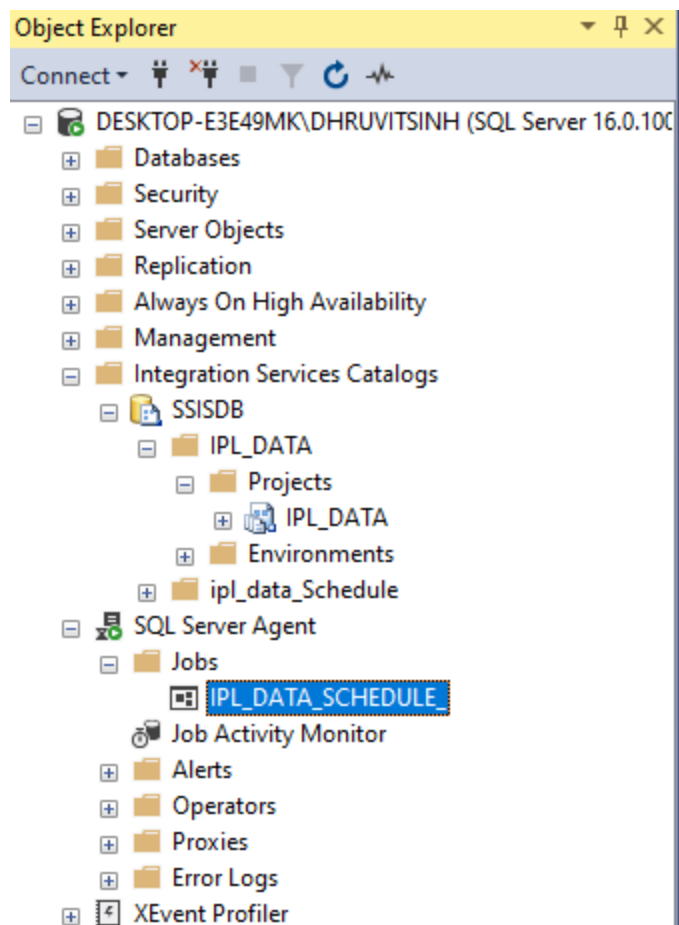
#### Execution Log:

As shown in the **Job History (Figure 1)**, the schedule is configured to run successfully, as indicated by the status "The job succeeded." This ensures that all tasks within the workflow are executed without errors, maintaining data integrity.

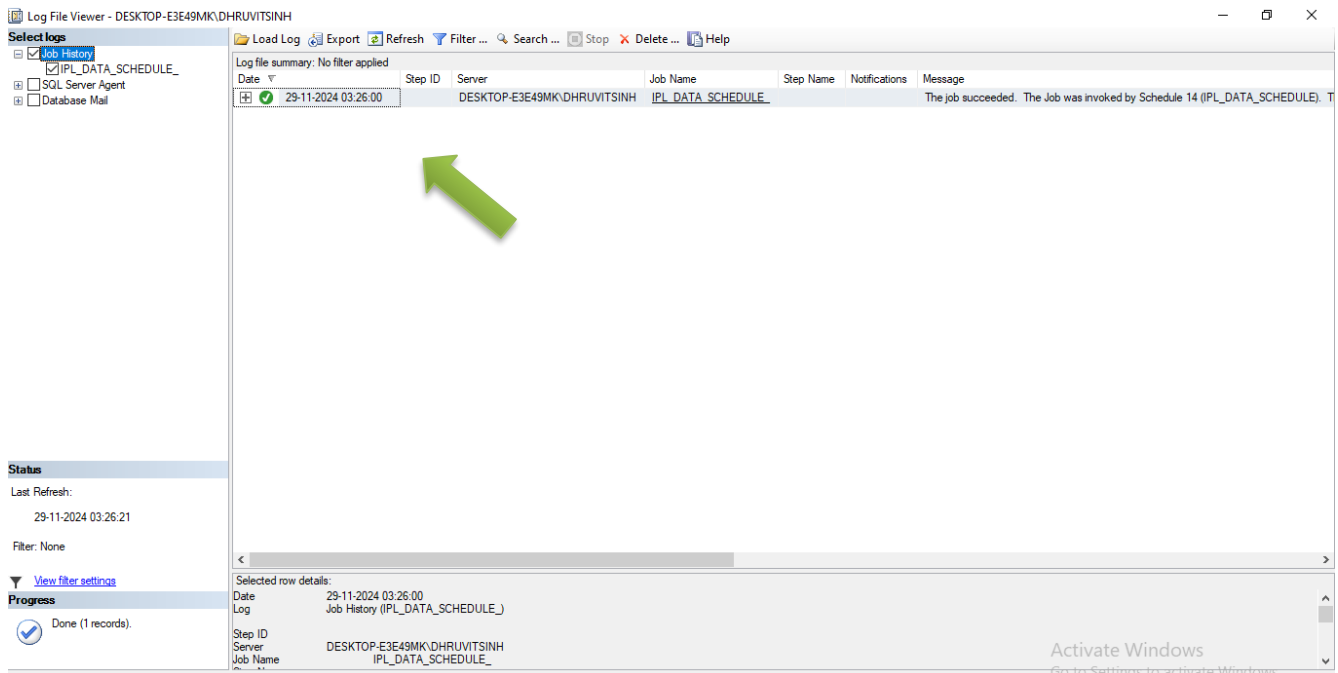
#### Object Explorer Overview:

The configuration of the job can be observed in the **SQL Server Management Studio Object Explorer (Figure 2)**. Key aspects include:

- **SSISDB Integration Services Catalogs:** Houses the IPL\_DATA project, which contains the SSIS package responsible for the ETL operations.
- **SQL Server Agent – Jobs:** Lists the IPL\_DATA\_SCHEDULE job, which is triggered according to the predefined schedule.
- **Efficiency:** Automates repetitive tasks, saving time and resources.
- **Accuracy:** Reduces the chances of human error in data processing.
- **Reliability:** Ensures that critical processes occur at consistent intervals, improving overall system reliability.



## Schedule



## Challenges and Solutions

During the development of the **IPL Data ETL Process**, several challenges were encountered that required innovative problem-solving to maintain the accuracy and integrity of the data pipeline. Below are two key challenges faced during the project and the steps taken to resolve them.

### Challenge 1: Issue with Redirect Row for Bad Data Handling Due to Delimiter in Venue Name

#### Problem Description:

During the ETL process, some venue names in the source data contained unexpected delimiters (e.g., commas or special characters). These delimiters caused mismatches between the Venue\_Name in the Venue table and the Venue\_Id in the Match table, resulting in foreign key violations during the data load.

#### Solution:

1. **Identifying Bad Data:** Detected rows with problematic delimiters that caused referential integrity issues.
2. **Disabling Foreign Key Constraints:** Temporarily disabled the foreign key constraint on the Match table to allow data insertion while isolating bad data.
3. **Redirecting Bad Data:** Used the **Redirect Row** feature in SSIS to capture problematic rows and redirect them to a separate error log table for further review.

4. **Client Review and Data Correction:** Allowed the client to review the error log and fix issues such as removing incorrect delimiters.
5. **Re-enabling Foreign Key Constraints:** Restored the foreign key constraints after cleaning the data.
6. **Reprocessing Cleaned Data:** Reprocessed the cleaned data to ensure that Venue\_Id and Venue\_Name matched correctly in all tables.
7. **Continuous Monitoring and Logging:** Implemented validation checks and logging mechanisms to prevent similar issues in future ETL runs.

**Outcome:**

This approach ensured the ETL process could handle data inconsistencies effectively without compromising data integrity, allowing for seamless reprocessing of corrected records.

---

## Challenge 2 : Violation Unique Key for Primary Key in Certain Tables

**Problem Description:**

Several source tables, including Player\_Match, Ball\_by\_Ball, Batsman\_Scored, Extra\_Runs, and Wicket\_Taken, lacked a unique key. This made it difficult to identify individual records for lookups and transformations. Even composite keys were insufficient for uniquely identifying rows, which complicated data validation and transformation processes.

**Solution:**

1. **Using Lookup Transformation for Uniqueness:** Employed the **Lookup Transformation** in SSIS to match all relevant columns (e.g., Match\_Id, Over\_Id, Ball\_Id) instead of relying on a single key or composite key.
2. **Incremental Load – Insert-Only Strategy:** Adopted an insert-only incremental load approach, focusing exclusively on inserting new records rather than updating existing ones.
3. **Ensuring Data Consistency:** Validated data to filter out duplicate or inconsistent records during the load process, ensuring that only unique records were inserted.
4. **Performance Optimization:** Optimized the lookup process by reducing the number of columns involved, improving efficiency for large datasets.
5. **No Updates for Historical Data:** Avoided updates to historical data, simplifying the transformation process and eliminating potential errors.

**Outcome:**

This solution ensured that the ETL process handled non-unique source tables effectively, maintaining data integrity and improving performance while simplifying the overall workflow.