# Soccer Data Warehouse Presentation Deck

By: Gary Buckley, Scott Kalich, Brennan Tolman, Zach Zhang
Group 9

# Background

Vision & Objective

Deliver effective team game strategies ideas to the organization based on data. To help organization create and implement efficient game tactics and achieve greater business benefits.

Product & Service

SQL, Talend, DBMS, BI…

# Current Scenario of the Business

Creating a data warehouse for better analysis and discover the relations of the data. To meet the objectives, the data warehouse serves to explore the effect of the locations on the goal rate.

Design team tactics based on scientifically proven information. It can improve the pertinence of training and the effectiveness in actual competition.

# Challenges

- Data sets have outliers and unintuitive structure

- Large amount of data stored in 2 sets

- Unable to easily query questions based on items such as location

# Solution

- Remove outliers prior to the data warehouse and use Talend to split tables

- Split data sets in Talend to fit a dimensional model

- Use ETL processes to create and label zones to easily slice data

# Why a Data Warehouse?

- Connects multiple data sources together for analysis
- Creates consistent quality and information for end-users
- Allows a security hierarchy to block unwanted writes/updates to data
- Standardized naming conventions
- Increases ease of queries
- Decreases latency between query and result

# Create a dataflow that supports end-users

# Business Requirements

Current Requirements

- Where are most shots generated from?
- Where are most goals generated from?
- What time of the game yields the most shots and goals?
- Which players lead the league in shots and goals?

Potential Future Requirements

- Where do key passes (passed that lead to shots) come from?
- Where do assists come from?

# Bus Matrix

| Business Process / Match Strategies | Location on field | Time in Game | Player | Game | Event Type | Date | Notes |
|---|---|---|---|---|---|---|---|
| **Shooting** | cartesian product of coordinates | | | | | | **Team will be captured in player dimension. Type 2 changes** |
| Taking Shots | x | x | x | x | x | | Where should try to take shots from? |
| Defending Shots | x | x | x | x | x | | How do we limit shots in high-conversion zones? |
| Scoring Goals | x | x | x | x | x | | Where are most goals scored from? |
| Assisting Goals | x | x | x | x | x | | How does the area passed from influence goal percentage? |
| | | | | | | | |
| **Set Pieces** | | | | | | | |
| Corners | x | x | x | x | x | | Are out-swinging, in-swinging, or short corners most effective? |
| Free Kicks | x | x | x | x | x | | When should we shoot vs. cross a free kick? |
| | | | | | | | |
| **Player Management** | | | | | | | **Include date to see if injuries happen more often on short rest** |
| Injuries | x | x | x | x | x | x | What conditions lead to injuries more often? |
| Subbing | | x | x | x | x | | How and when should we use our subtitutes? |
| | | | | | | | |
| **Passing** | | | Role-playing for 2-player events | | | | **use views of role-playing player dimension for passer/receiver** |
| Offense | x | x | x | x | x | | what length of pass sequences lead to goals? |

Dimensions

# Dimensional Model Approach

General Strategy

- Dimension tables are created based on attributes
- Attributes uniquely identify instances on dimensions
- Dimensions link together to a fact table
- Fact table holds frequently queried information

# Dimensional Model

# Game Dimension

# Game Time Dimension

dim_game_time
# time_key
seconds_from_half
seconds_from_start
period
period_desc
minute

# Event Dimension (modified)



dim_event
\# event_key
event_type
event _category

# Location Dimension (modified)



dim_location
# location_key
x
y
zone

# Player Dimension (role-playing)

# Event Fact Table

event_fact
# game_key
# time_key
# player_1_key
player_2_key
event_key
location_key

# Dimensional Model

# Extract-Transform-Load



- Grab data from sources (flat file)
- Transform data by removing outliers, noise, and unused data
- Load clean data into warehouse to be used

# Dimension Creation

# Data Loading

# Queries

# Keying

## Location

| | | |
|---|---|---|
| **Type 1 fields** | | |
| zone | | |

**Source keys**

x
y

**Surrogate keys**

| | |
|---|---|
| name | location_key |
| creation | Table max + 1 |
| complement | |

**Type 2 fields**

**Versioning**

| type | name | creation | comple |
|---|---|---|---|
| start | scd_start | Job start time | |
| end | scd_end | NULL | |
| version | scd_version | | |
| active | scd_active | | |

## Event

**Type 1 fields**

event_category

**Source keys**

event_type

**Type 2 fields**

**Surrogate keys**

| | |
|---|---|
| name | event_key |
| creation | Table max + 1 |
| complement | |

**Versioning**

| type | name | creation | comple |
|---|---|---|---|
| start | scd_start | Job start time | |
| end | scd_end | NULL | |
| version | scd_version | | |
| active | scd_active | | |

## Time

**Type 1 fields**

minute
period_desc
seconds_from_start

**Source keys**

period
seconds_from_half

**Type 2 fields**

**Surrogate keys**

| | |
|---|---|
| name | time_key |
| creation | Table max + 1 |
| complement | |

**Versioning**

| type | name | creation | comple |
|---|---|---|---|
| start | scd_start | Job start time | |
| end | scd_end | NULL | |
| version | scd_version | | |
| active | scd_active | | |

# Keying



**Player 1**

| Source keys |
|---|
| player_name_1 |

| Surrogate keys | |
|---|---|
| name | player_1_key |
| creation | Table max + 1 |
| complement | |

| Type 2 fields |
|---|
| player_team_1 |

| Versioning | | | | |
|---|---|---|---|---|
| type | name | creation | comple |
| start | team_start | Job start time | |
| end | team_end | NULL | |
| ☐ version | scd_version | | |
| ☐ active | scd_active | | |

**Player 2**

| Source keys |
|---|
| player_name_2 |

| Surrogate keys | |
|---|---|
| name | player_2_key |
| creation | Table max + 1 |
| complement | |

| Type 2 fields |
|---|
| player_team_2 |

| Versioning | | | | |
|---|---|---|---|---|
| type | name | creation | comple |
| start | team_start | Job start time | |
| end | team_end | NULL | |
| ☐ version | scd_version | | |
| ☐ active | scd_active | | |

**Game**

| Source keys |
|---|
| game_date |
| home_team |
| away_team |

| Surrogate keys | |
|---|---|
| name | game_key |
| creation | Table max + 1 |
| complement | |

# Sample Data

# Sample Data

# Fact Table Creation

# Metabase Set Up

- Allows for drilling filters
- Easy graphics on dashboard
- Updates when new data enters

# Metabase Dashboard

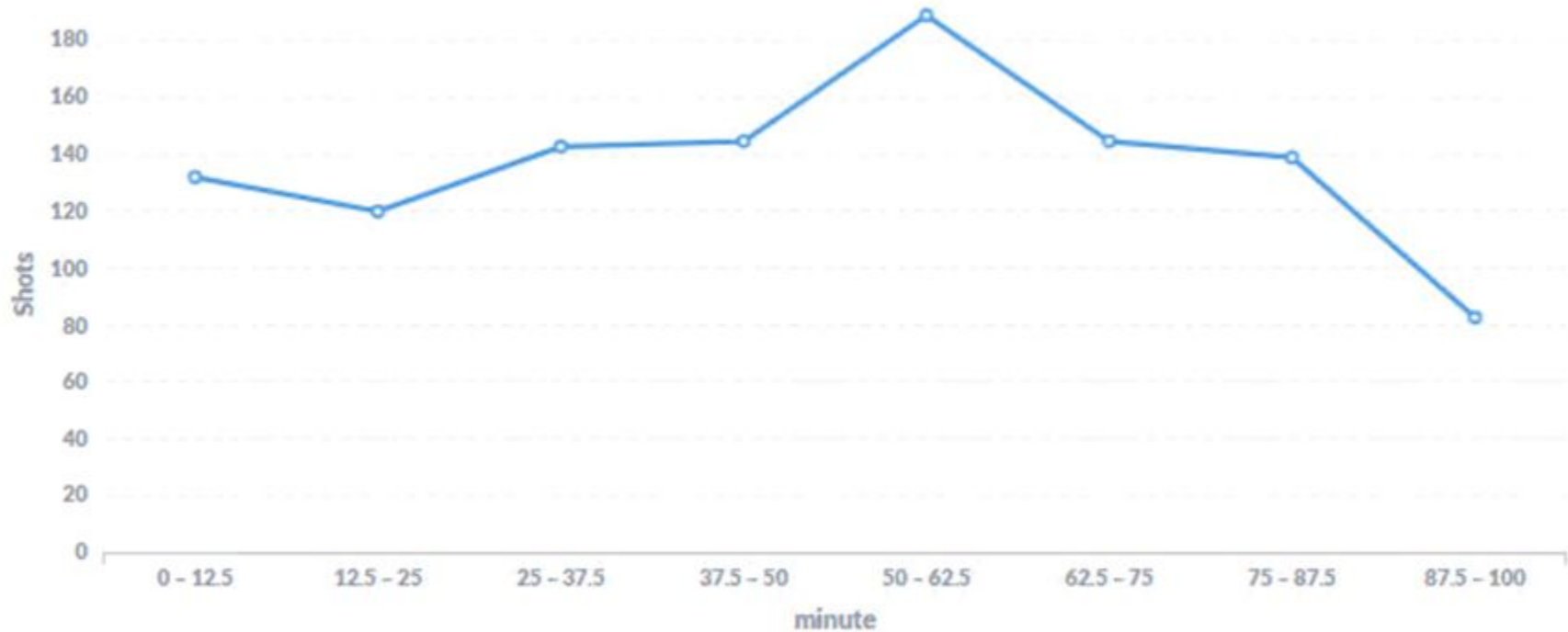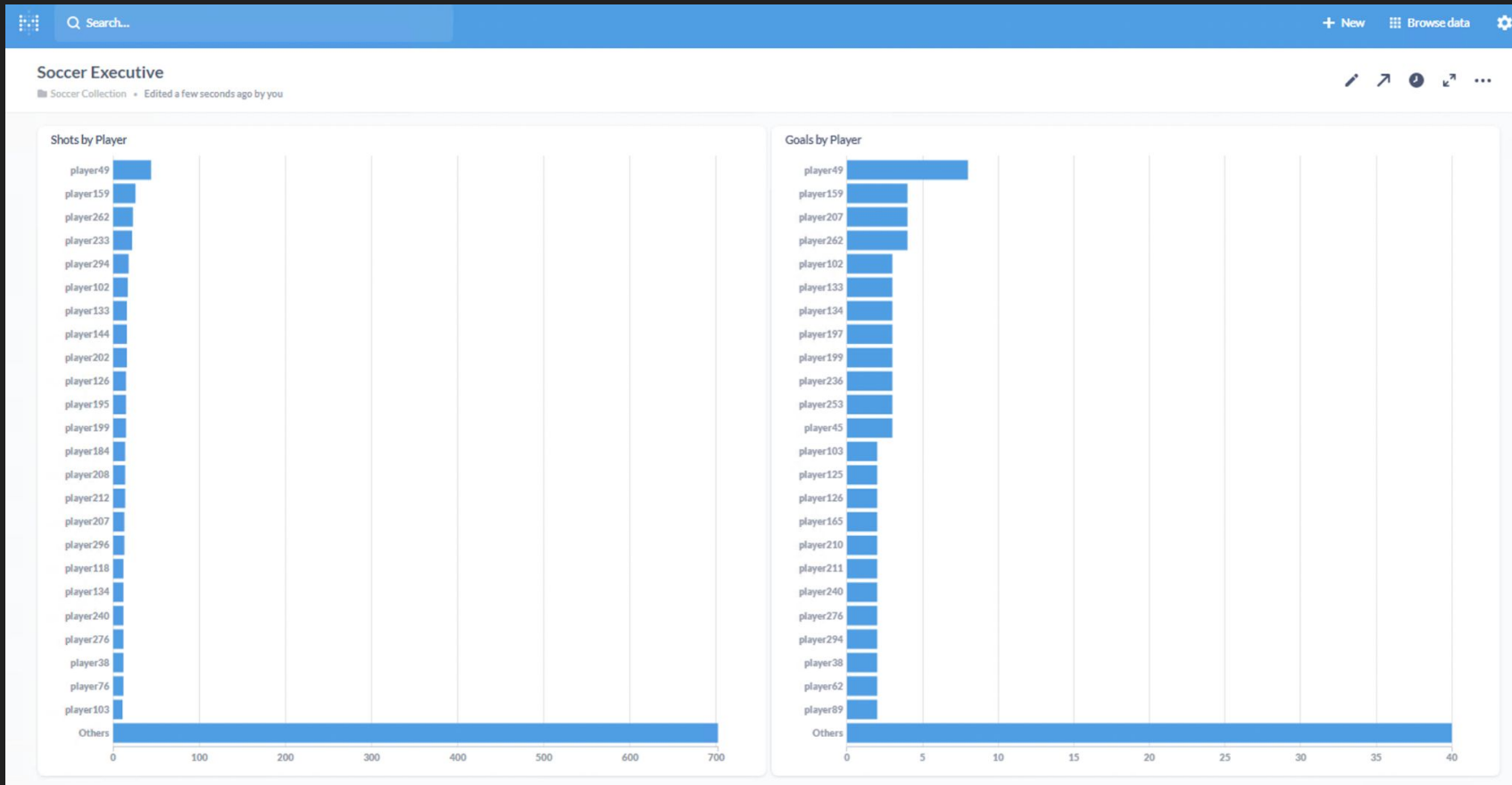# Shots by Zone
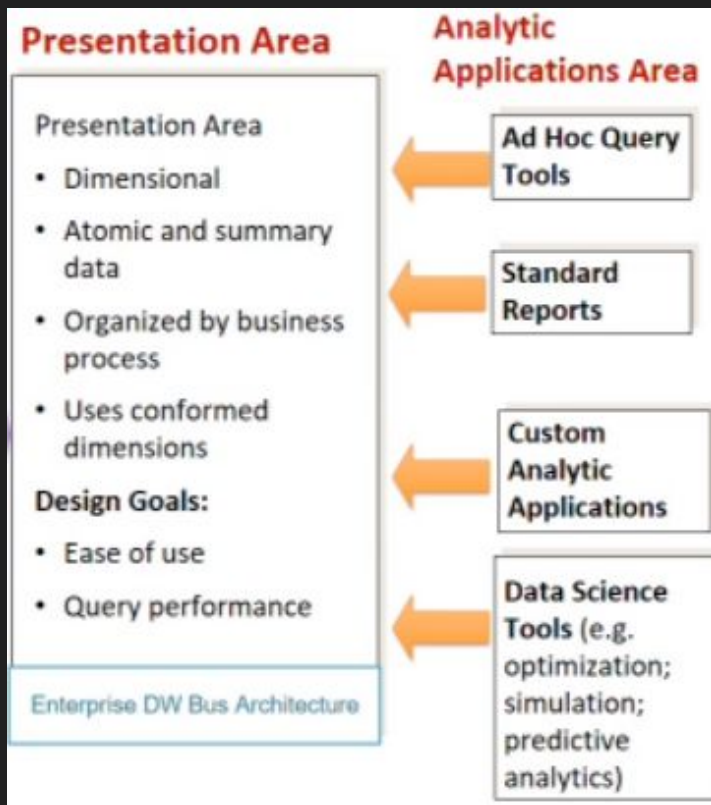
# Goals by Zone

# Shots by Minute

# Metabase Executive Dashboard

# Analytics Ease

# Future Objectives Requirements

- Where do key passes (passed that lead to shots) come from?
- Where do assists come from?

# Conclusion

The business can now:

- Easily load data into a uniform system
- Provide consistent data to analysts
- Manage permissions for specific groups
- Create queries at the speed of thought
- Implement dashboards to easily see how data trends change over time

# Presentation Video

https://youtu.be/lwbl8Ni4TPI