

BASEBALL STATS PIPELINE PROJECT  
BY SCOTT SILVERSTEIN

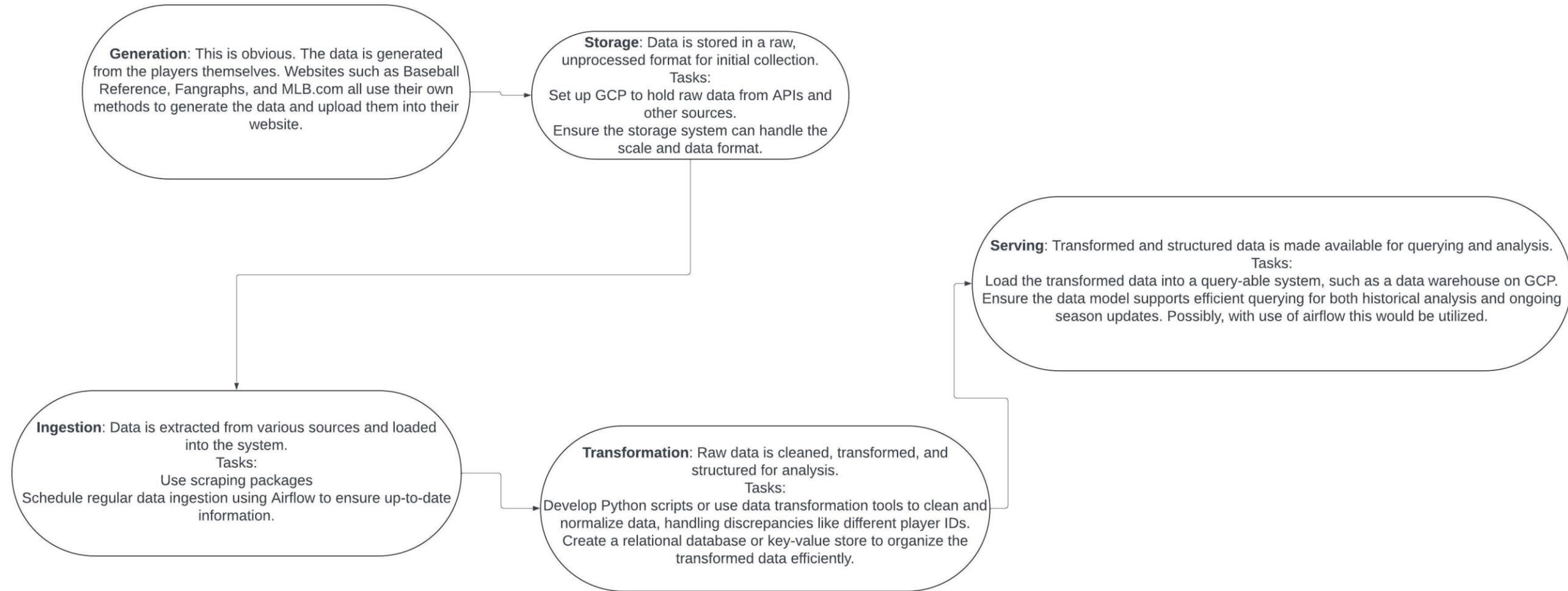


# Why a Baseball Stats Pipeline

Because it's the reason I am here!



# What I did?



Generation

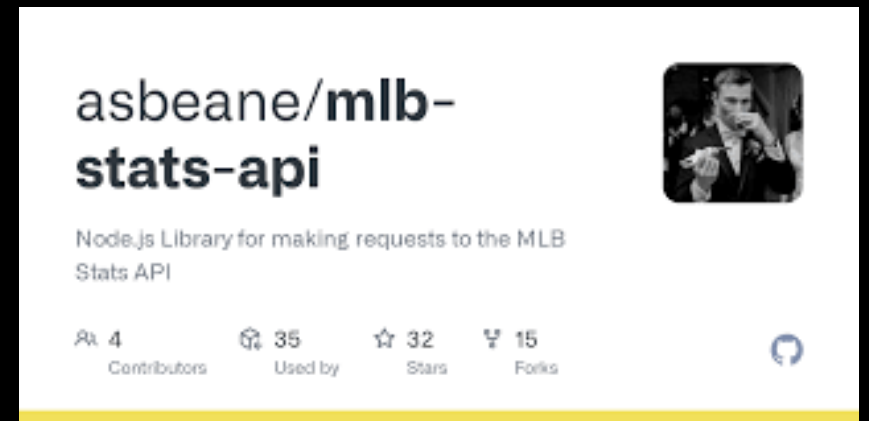
# **FANGRAPHS**

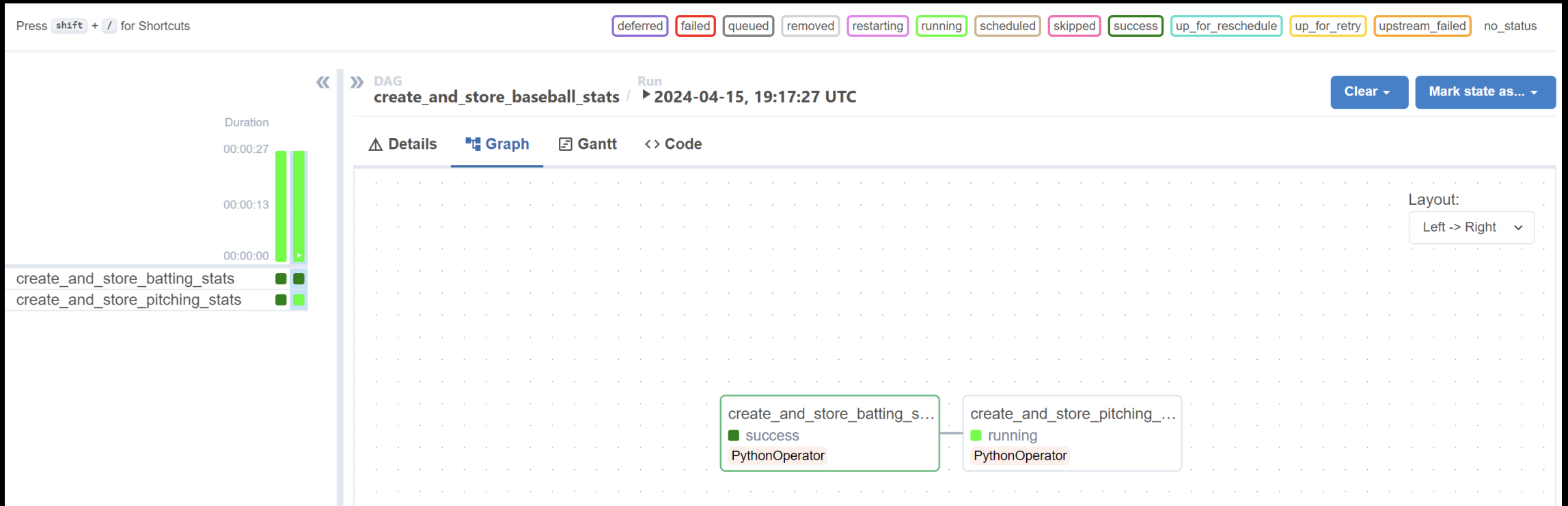
## 2022 Top 50 Free Agents



# Pybaseball Package

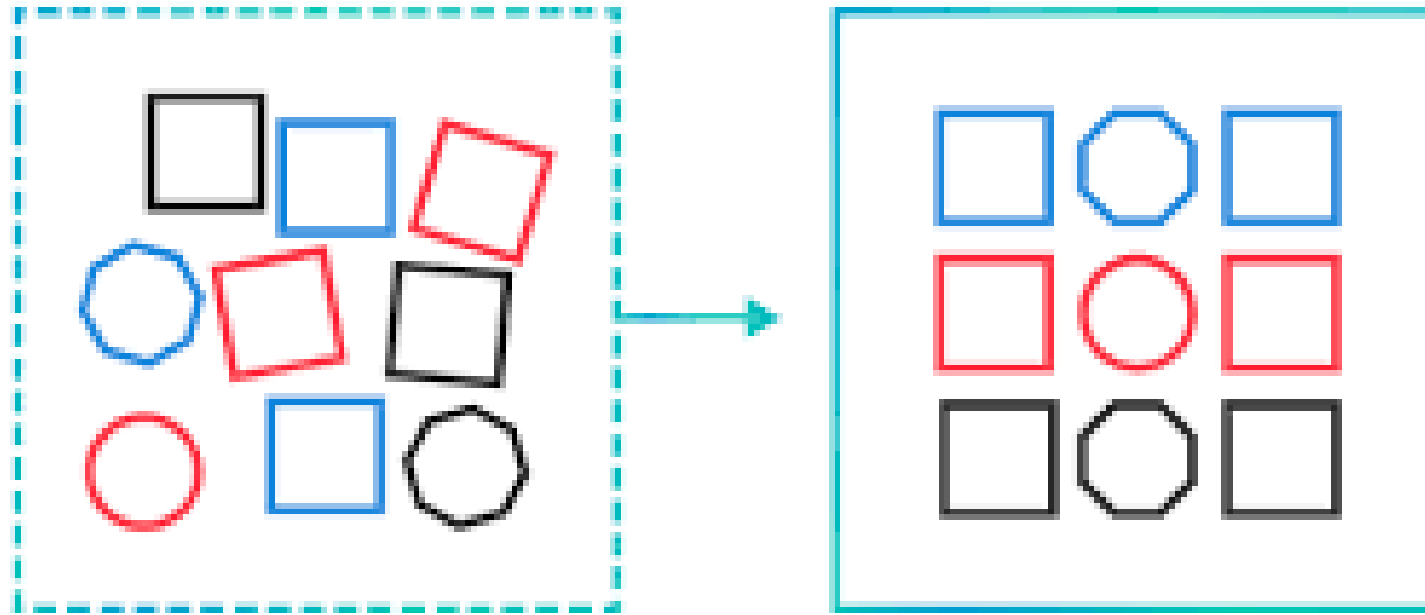
- I used Pybaseball package to scrape data
- Many options
  - Baseball\_scraper
  - Mlb Stats API
- MLB API is both great and terrible





# Ingestion

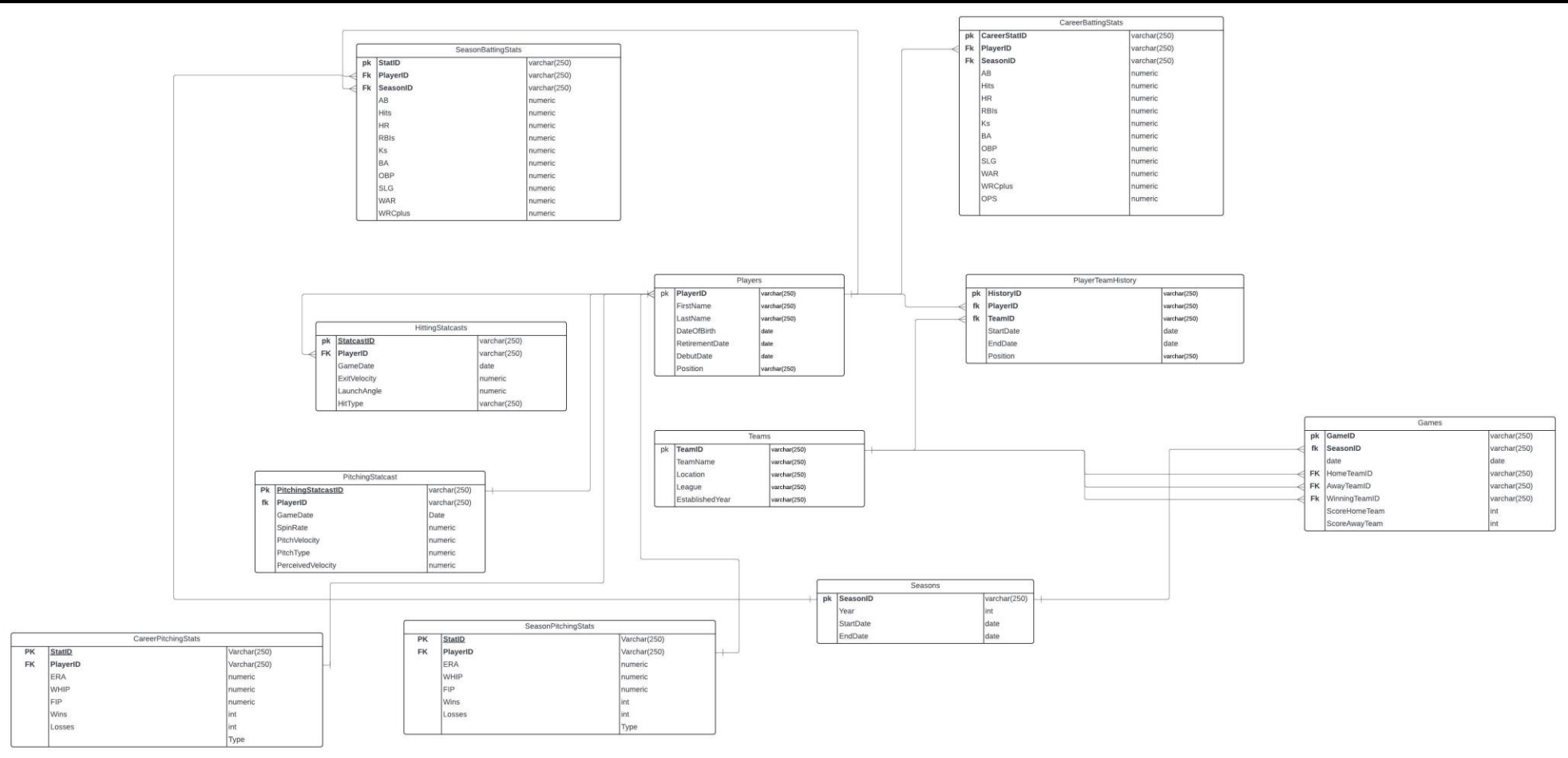
- This was done by creating tables and using functions to fetch data from pybaseball scraper
- Eventually took these functions and put them into an airflow for easier future automation



# Transformation

- This was done by converting all the datasets generated in pybaseball and formatting them to fit my logical model.

```
Number of columns in SeasonBattingStats: 317
```



# Challenges



Teams		
pk	<u>TeamID</u>	varchar(250)
	TeamName	varchar(250)
	Location	varchar(250)
	League	varchar(250)
	EstablishedYear	varchar(250)

PitchingStatcast		
Pk	<u>PitchingStatcastID</u>	varchar(250)
fk	<u>PlayerID</u>	varchar(250)
	GameDate	Date
	SpinRate	numeric
	PitchVelocity	numeric
	PitchType	numeric
	PerceivedVelocity	numeric

SeasonPitchingStats		
PK	<u>StatID</u>	Varchar(250)
FK	<u>PlayerID</u>	Varchar(250)
	ERA	REAL
	WHIP	REAL
	FIP	REAL
	Wins	REAL
	Losses	REAL
	OTHER STATS CONTINUED	REAL

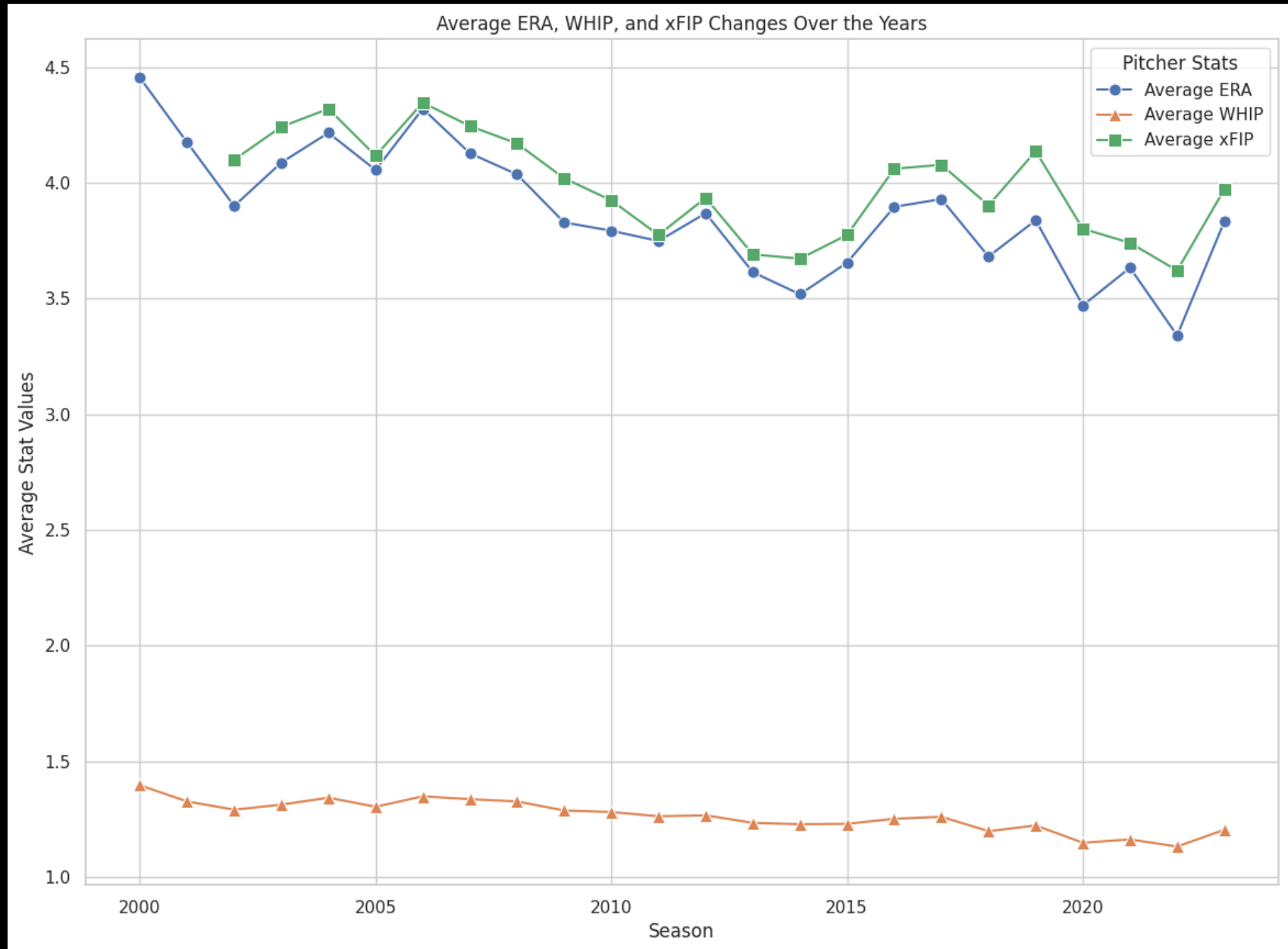
Players		
pk	<u>PlayerID</u>	varchar(250)
	FirstName	varchar(250)
	LastName	varchar(250)
	DateOfBirth	date
	RetirementDate	date
	DebutDate	date
	Position	varchar(250)
FK	<u>TeamID</u>	INTEGER

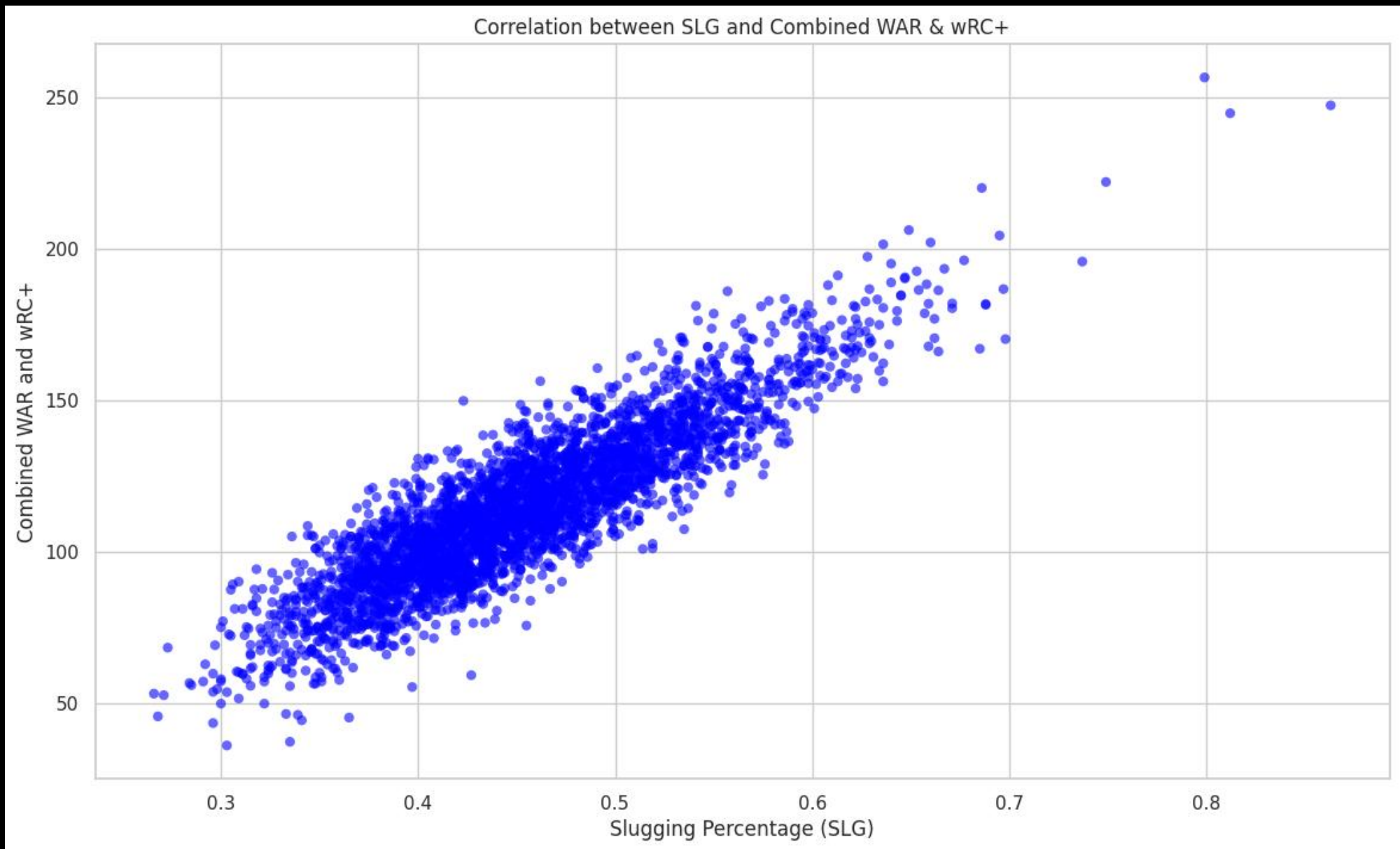
SeasonBattingStats		
pk	<u>StatID</u>	INTEGER
Fk	<u>PlayerID</u>	INTEGER
	Hits	REAL
	HR	REAL
	RBIs	REAL
	Ks	REAL
	BA	REAL
	OBP	REAL
	SLG	REAL
	WAR	REAL
	WRCplus	REAL
	OTHER STATS CONTINUED	REAL

HittingStatcasts		
pk	<u>StatcastID</u>	varchar(250)
FK	<u>PlayerID</u>	varchar(250)
	GameDate	date
	ExitVelocity	numeric
	LaunchAngle	numeric
	HitType	varchar(250)

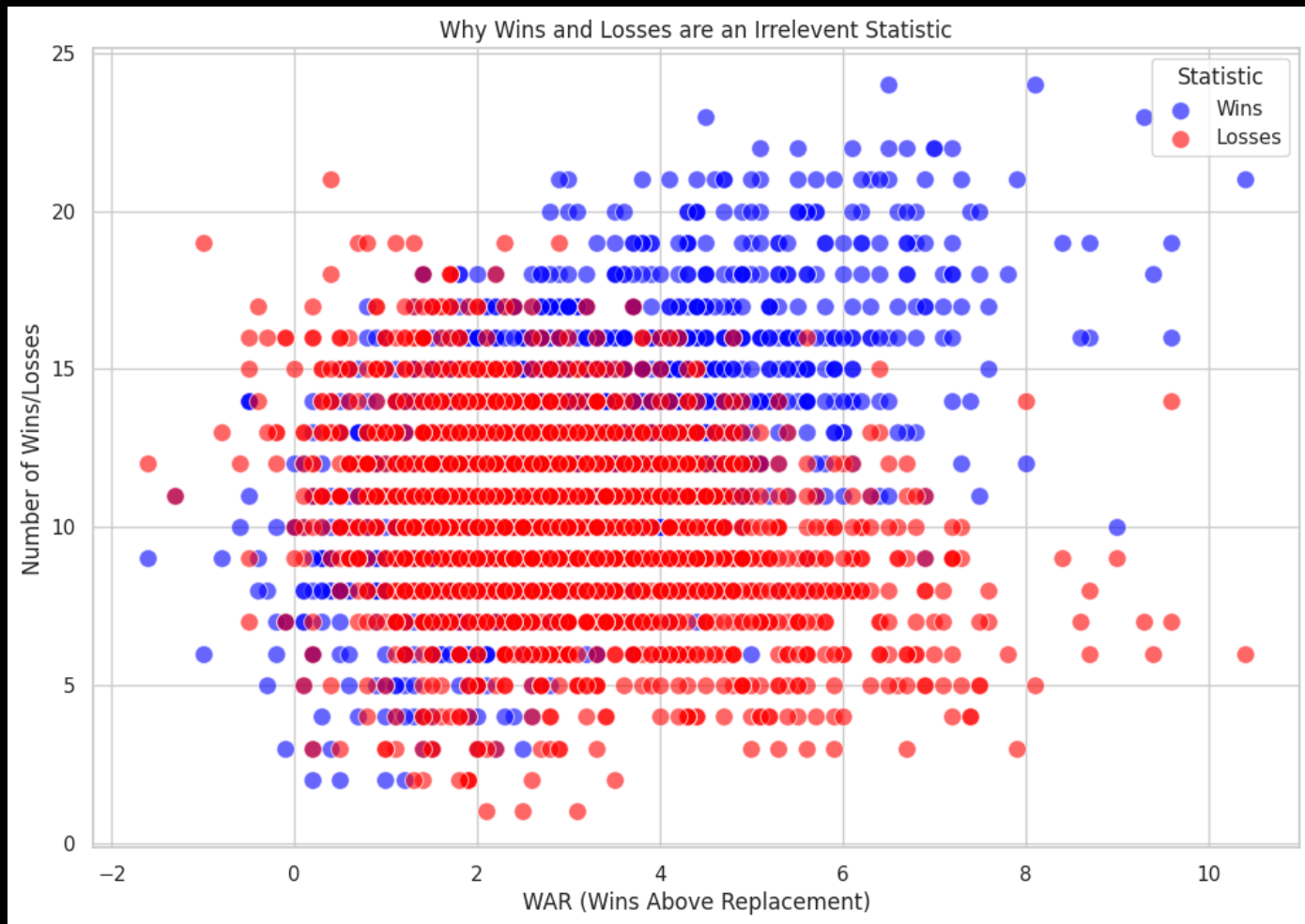
Final ERD

# Serving





	Parameter	Wins	Losses
0	Coefficient	1.1576772471390062	-0.5858824832461207
1	P-value	3.230995391795463e-131	2.709156638605293e-41
2	R-squared	0.2814402070241847	0.09587028552839727
3	Adj. R-squared	0.28104056309038283	0.09536743251701152



# Storage

OBJECTS

CONFIGURATION

PERMISSIONS

PROTECTION

LIFECYCLE

OBSERVABILITY

INVENTORY REPORTS

OPERATIONS

Folder browser

baseball\_project12

path/

to/

destination/

Buckets > baseball\_project12 > path > to > destination

UPLOAD FILES

UPLOAD FOLDER

CREATE FOLDER

TRANSFER DATA

MANAGE HOLDS

EDIT RETENTION

DOWNLOAD

DELETE

Filter by name prefix only

Filter

Filter objects and folders

Show

Live objects only

	Name	Size	Type	Created	Storage class	Last modified	
	baseball_stats.db	3.8 MB	application/octet-stream	Apr 15, 2024, 2:16:50 PM	Standard	Apr 15, 2024,	



# Google Cloud

## Future Plans

---

- I now have a working database and an airflow that can easily automatically update my data!
- Try to figure out issues with different ID's for different websites
- Make a more cohesive ERD with information on player contracts, player history, etc.

**500+**  
**BASEBALL**  
**BLOG NAME**  
**IDEAS**

