



Anthony Rizzo hugs Kris Bryant, while Mike Montgomery, Javier Baez and Addison Russell celebrate after winning the World Series in Cleveland on Wednesday night.

# At last!

Cubs capture first title in 108 years in extra-inning Game 7 thriller

**PHIL KILGUS**  
Chicago Tribune

**CLEVELAND** — Finally, the most epic drought in sports history broke, and the Cubs won their first World Series title since 1908. After more than a century of waiting, the Cubs won the 2016 World Series with a 5-3 Game 7 victory over the Indians in 10 innings Wednesday night at Progressive Field. The triumph completed their climb back from a 3-1 Series deficit to claim their first championship since 1908.

This is not a dream. The Cubs did it. It was real. And it was the greatest moment in the history of the franchise.

Seven days of intense celebration after the final and in the 1900s, and fans celebrated with the world's largest group hug, remembering all the loved ones who would only imagine what it would be like to experience the triumph of their lives.

The 2016 Cubs are the team that defined the word "underdog" more than any other team in the history of the franchise. It was a team that was once thought to be a joke, but they proved them wrong. They were the team that was once thought to be a joke, but they proved them wrong. They were the team that was once thought to be a joke, but they proved them wrong.

The Cubs' journey was a roller coaster ride. They were the team that was once thought to be a joke, but they proved them wrong. They were the team that was once thought to be a joke, but they proved them wrong. They were the team that was once thought to be a joke, but they proved them wrong.

The Cubs' journey was a roller coaster ride. They were the team that was once thought to be a joke, but they proved them wrong. They were the team that was once thought to be a joke, but they proved them wrong. They were the team that was once thought to be a joke, but they proved them wrong.

The Cubs' journey was a roller coaster ride. They were the team that was once thought to be a joke, but they proved them wrong. They were the team that was once thought to be a joke, but they proved them wrong. They were the team that was once thought to be a joke, but they proved them wrong.



# WAR and PCA

Predicting Player Value in Baseball

**Presented by Team: The Matplotlibs**

Addison Naylor, Eleni Bovalis, Sahil Sangani,  
Tanvee Athavale, Annika Luthe







# Table of contents



**01**

## Project Overview

What is WAR? What problem do we want to solve?

**02**

## Exploratory Data Analysis

Our data cleaning and EDA process.

**03**

## Modeling

The models we used and why.

**04**

## Challenges

Roadblocks we faced and how we dealt with them.

**05**

## Conclusion

Our results and next steps!





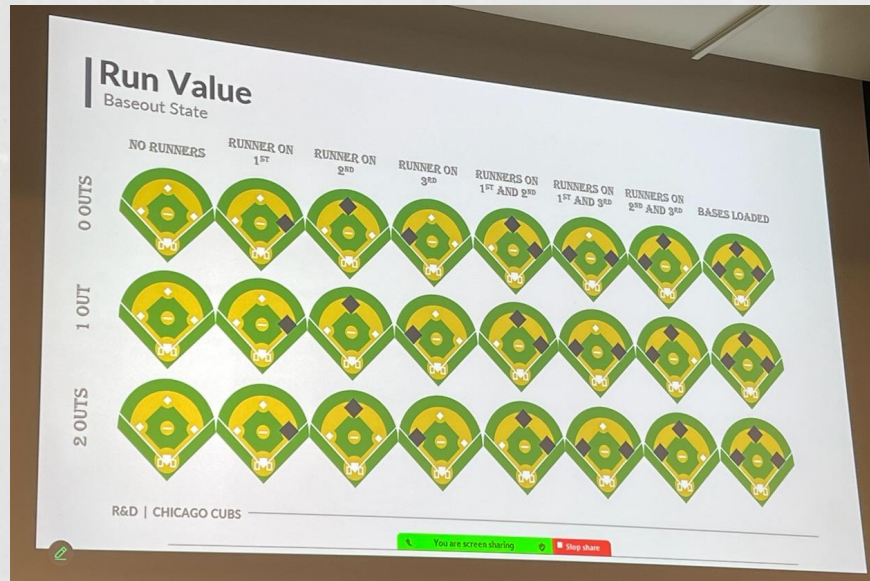
# 01

# Project Overview



# Initial Idea

Our first idea for this project was to generate a **lineup optimization order using Markov chains**. However, after attending a Quantitative Psychology baseball seminar with Cubs Assistant General Manager Dr. Ehsan Bokhari, we realized that we would have to complete **too many simulations** to create our model, which would require **a lot of computational power**.



**15** usable players in  
**9** possible positions  
would need  
**1,816,214,400+**  
simulations.





# Our Project

Our **goal** was to build a model to predict how valuable a baseball player will be to their team in the upcoming season based on their current season's statistics.

How this project can be used:

- 👤 Help baseball managers and executives make trading decisions
- 👤 Decide which players to bring up from the minor leagues
- 👤 Fans can use it to predict which players will contribute the most to their favorite team in the upcoming season!





# What is WAR?

In baseball, "WAR" stands for **Wins Above Replacement** which measures a player's overall value to a team based on how many wins they can earn compared to a "replacement-level" player.

Example: A player with a **WAR of 6.0** contributes **6 more wins** to their team compared to a different player who could be acquired at minimal cost (like a minor league player).

Fun fact: Babe Ruth has the highest career WAR of all time: 182.6



# 02 Exploratory Data Analysis







# Data Cleaning

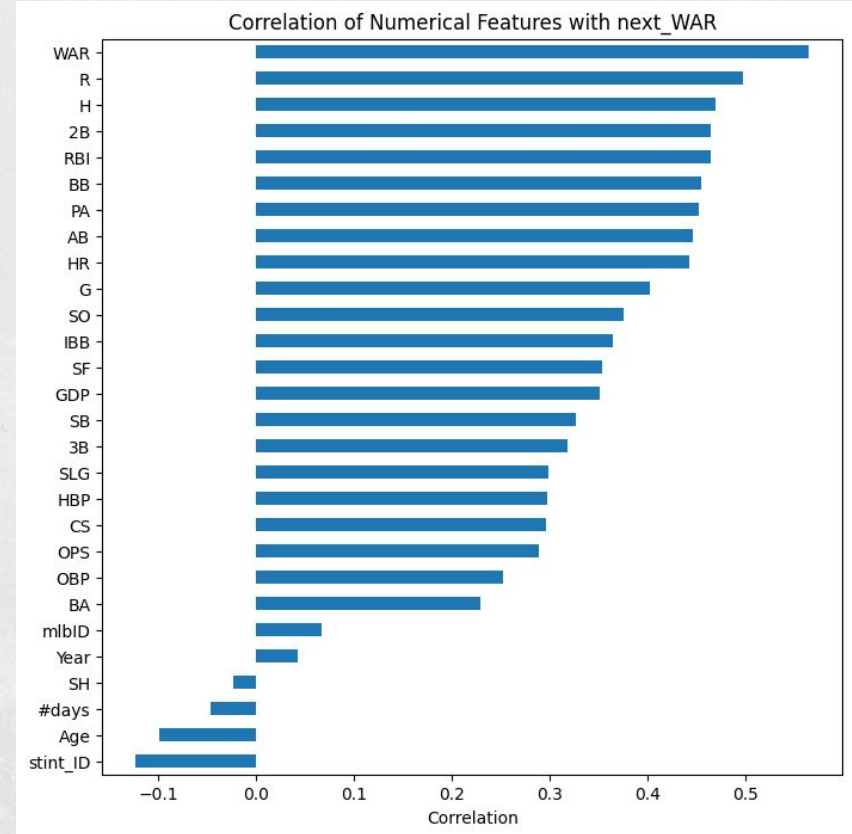
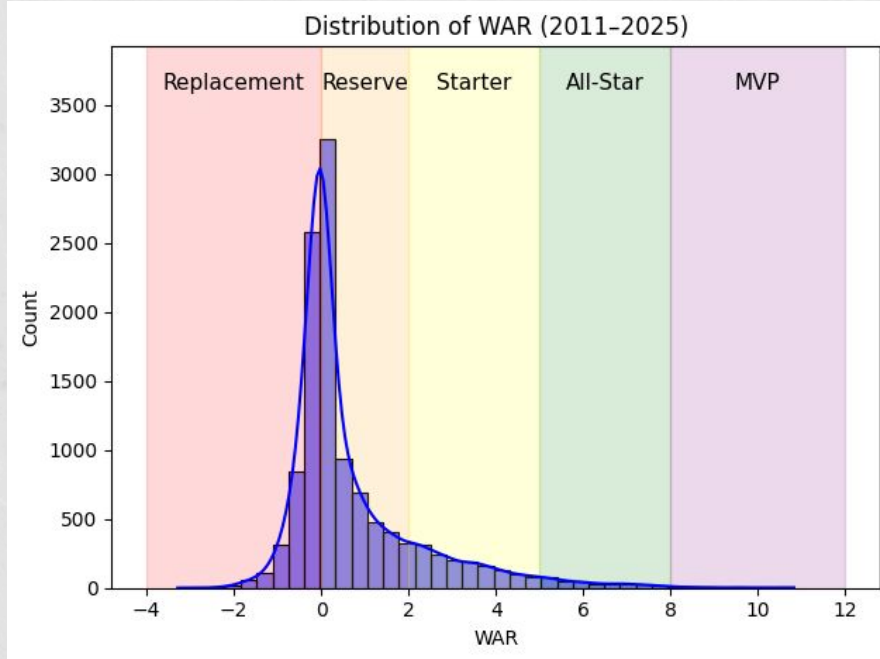


We used data from all MLB teams from **2010-2023 to train** our model and data from **2024-2025 to test** our model. We started with about **13000 rows of data**.



After filtering out unneeded variables, we found **1206 missing values**, and given our large dataset, decided to drop them.

# EDA



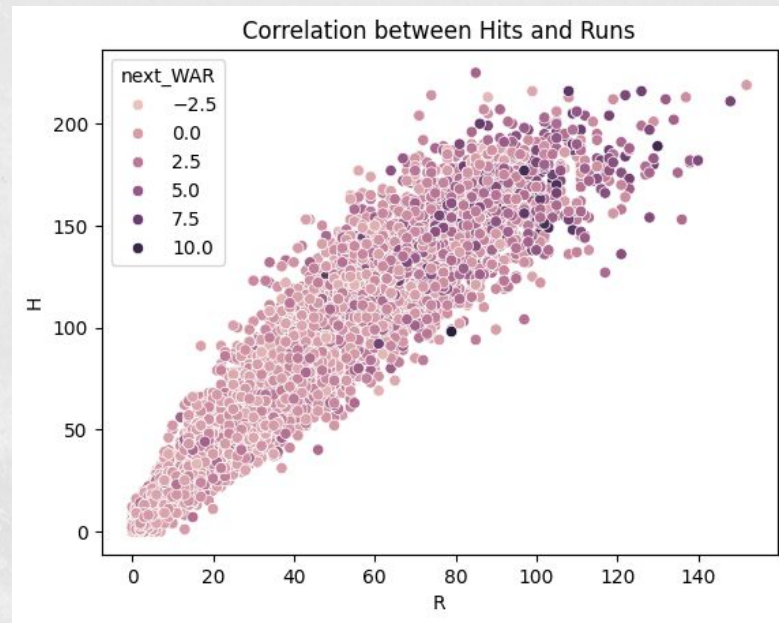
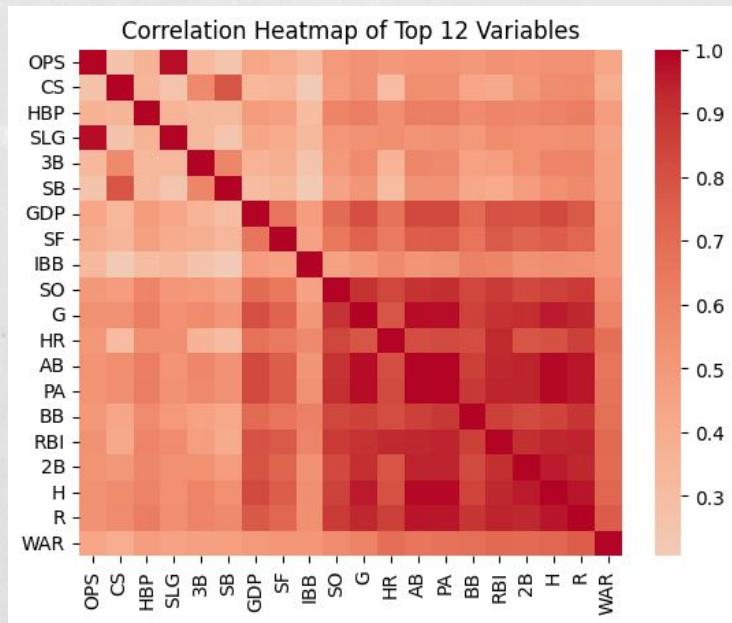


# PCA

n\_components = **0.90**

From 34 variables to just **14 components**

- 🛡 Mitigation of multicollinearity
- 🛡 Faster training times
- 🛡 Reduced overfitting



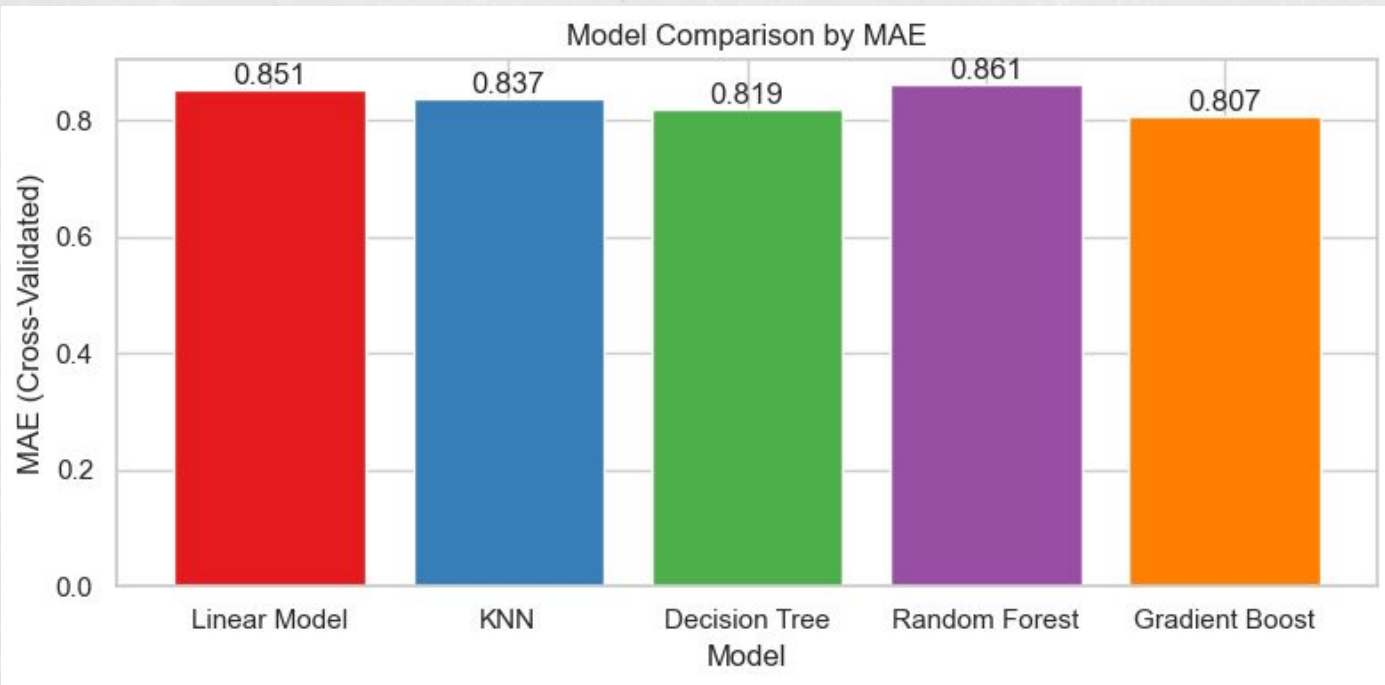


# 03

# Modeling



# Models





# Best Model

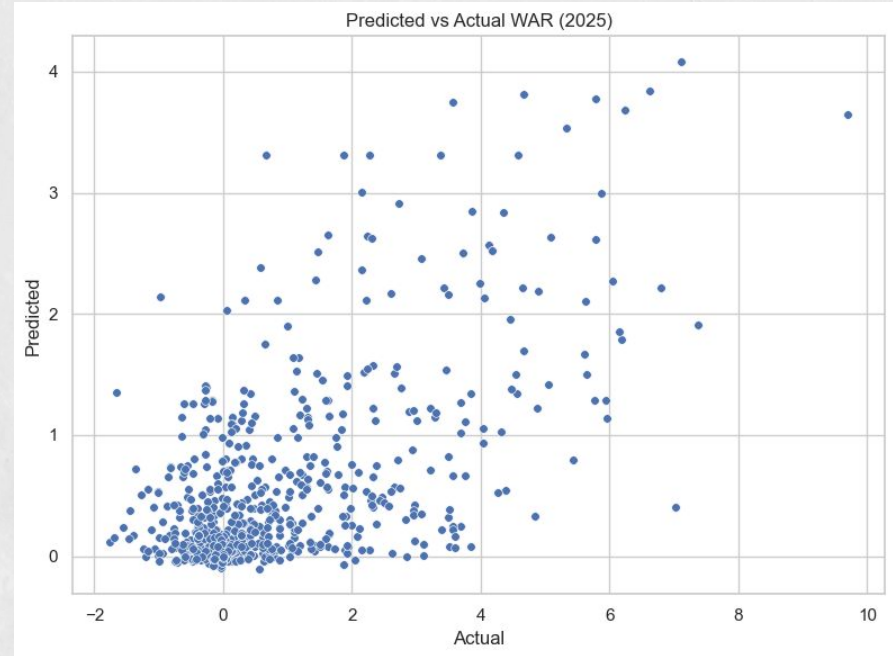
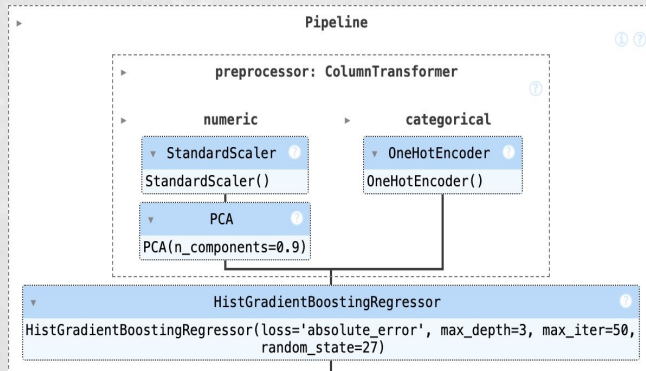


Histogram-Based Gradient Boosting  
Regressor from Scikit-Learn

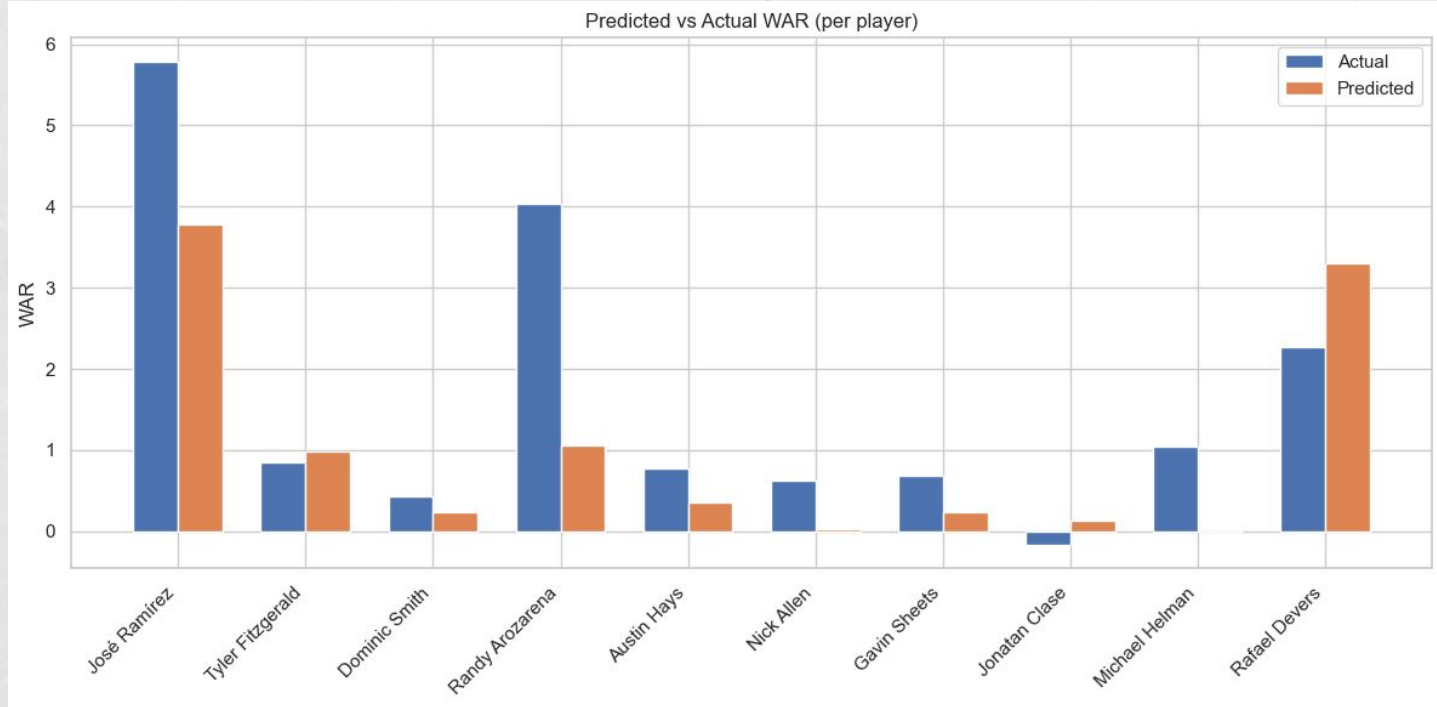
**5 Fold CV MAE: 0.807**

**Test RMSE: 1.368**

**Test MAE: 0.931**



# Prediction Accuracy





**04**

# Challenges





# Roadblocks & Breakthroughs



Figuring out our project (initial Markov chain idea)



Looking for an efficient way to import our data



Finding fixes for variable selection and multicollinearity



We asked for advice (Dr. Bokhari's presentation)



Got help from our Project Lead who introduced us to pybaseball



Implemented PCA for dimensionality reduction





05

# Conclusion



# Model Reliability

Our best model has a **MAE** of **0.807**, indicating that it is **reasonably trustworthy**.

It can be used for predicting whether a player is valuable, but shouldn't be used to determine if a player needs to be let go.

Example: If a player's WAR for next season is predicted to be 3.0, a baseball manager can be reasonably certain (within about a 2.2 – 3.8 range) that this player will make valuable addition to their team.







# Future



## Further analysis could include...



Predicting the value of pitchers  
(different than batters)



Incorporating minor league and  
college data



Predicting WAR further into the  
future



Time series analysis to examine  
trends across a player's whole career



The slide features a light gray background with a subtle, mottled texture. Red five-pointed stars are scattered in the corners: three in the top-left, three in the top-right, three in the bottom-left, and three in the bottom-right.

# Thank You!

**CREDITS:** This presentation template was created by [Slidesgo](#), and includes icons by [Flaticon](#), and infographics & images by [Freepik](#)

# Sources

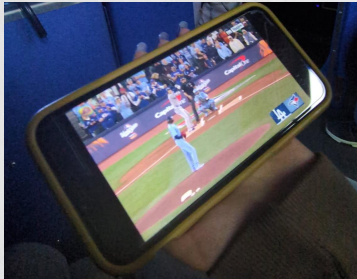
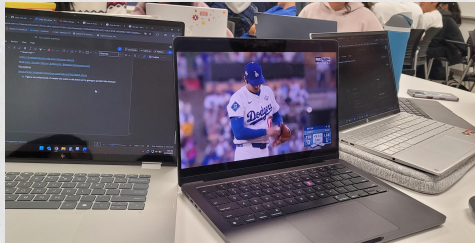


Image 1 (upper left): <https://www.businessinsider.com/kris-bryant-smiling-chicago-cubs-world-series-champions-2016-11>

Image 2 (bottom left):

<https://cloudfront-us-east-1.images.arcpublishing.com/gray/ZVM2P6PZVBMHLHDC2O4Z3LLWOJE.jpg>

Image 3 (upper middle):

<https://ca-times.brightspotcdn.com/dims4/default/d4ed453/2147483647/strip/true/crop/2048x1152+0+0/resize/1200x675!/quality/75/?url=https%3A%2F%2Fcalifornia-times-brightspot.s3.amazonaws.com%2Fb7%2F27%2F7d3a7e9705d8dcb949891fb6e519%2Fa-cxcxschilken-1478188604-snap-photo>

Image 4 (bottom middle):

<https://www.usatoday.com/gcdn/-mm-/cfddfdec04a3541d7c1e4182939fcd316de13cc/c=0-239-3933-2461/local/-/media/2016/11/03/USATODAY/USATODAY/636137488792285604-USP-MLB--World-Series-Chicago-Cubs-at-Cleveland-In.7.jpg>

Image 5 (right side):

<https://www.google.com/url?sa=i&url=https%3A%2F%2Fwww.chicagotribune.com%2F2016%2F11%2F03%2Fphotos-world-series-front-pages%2F&psig=AOvVaw2O4-H9GMPaEv-JriGYw3FB&ust=1763849512988000&source=images&cd=vfe&opi=89978449&ved=0CBYQjRxqFw6TCND3ic-hhJEDFQAAAAAdAAAAABAL>

Image 6 (right side)

[https://www.chicagotribune.com/wp-content/uploads/2024/04/CTC-L-Chicago-Cubs-03\\_187925179.jpg?w=1541](https://www.chicagotribune.com/wp-content/uploads/2024/04/CTC-L-Chicago-Cubs-03_187925179.jpg?w=1541)

WAR meaning Google search:

[https://www.google.com/search?q=what+does+wins+above+replacement+mean&oq=what+does+wins+above+re&gs\\_lcrp=EgZjaHJvbWUyBggAEAAyGAQyBwgAEAAyGAQyBggBEEUYOTIHCAlQABiABDIHCAMQABiABDIICAQQABgWGB4yCAGFEAAyFhgeMggIBhAAGBYyHjIICAcQABgWGB4yCAGlEAAyFhgeMggICRAAGBYyHtIICDQ0MzVqMGo3qAIAA&sourceid=chrome&ie=UTF-8](https://www.google.com/search?q=what+does+wins+above+replacement+mean&oq=what+does+wins+above+re&gs_lcrp=EgZjaHJvbWUyBggAEAAyGAQyBwgAEAAyGAQyBggBEEUYOTIHCAlQABiABDIHCAMQABiABDIICAQQABgWGB4yCAGFEAAyFhgeMggIBhAAGBYyHjIICAcQABgWGB4yCAGlEAAyFhgeMggICRAAGBYyHtIICDQ0MzVqMGo3qAIAA&sourceid=chrome&ie=UTF-8)

Replacement-level player meaning Google search:

[https://www.google.com/search?q=replacement-level+player+meaning&oq=replacement-level+player+meaning&gs\\_lcrp=EgZjaHJvbWUyBggAEUUYOTIHCAlQABgWGB4yDQgCEAAyhgMYgAQYigUyDQgDEAAyhgMYgAQYigUyCggFEAAyGAQYogQyBwgGEAAy7wUyBwgHEAAy7wXSAQgxOTMajBqN6gCALACAA&sourceid=chrome&ie=UTF-8](https://www.google.com/search?q=replacement-level+player+meaning&oq=replacement-level+player+meaning&gs_lcrp=EgZjaHJvbWUyBggAEUUYOTIHCAlQABgWGB4yDQgCEAAyhgMYgAQYigUyDQgDEAAyhgMYgAQYigUyCggFEAAyGAQYogQyBwgGEAAy7wUyBwgHEAAy7wXSAQgxOTMajBqN6gCALACAA&sourceid=chrome&ie=UTF-8)

All-time WAR ranking website: [https://www.baseball-reference.com/leaders/WAR\\_career.shtml](https://www.baseball-reference.com/leaders/WAR_career.shtml)

Baseball Data + screenshot: <https://www.baseball-reference.com>

Wii sports image: <https://i.ytimg.com/vi/hp3De8lflWQ/hqdefault.jpg>