

TEAM 23

Pitch Type Prediction

Baseball

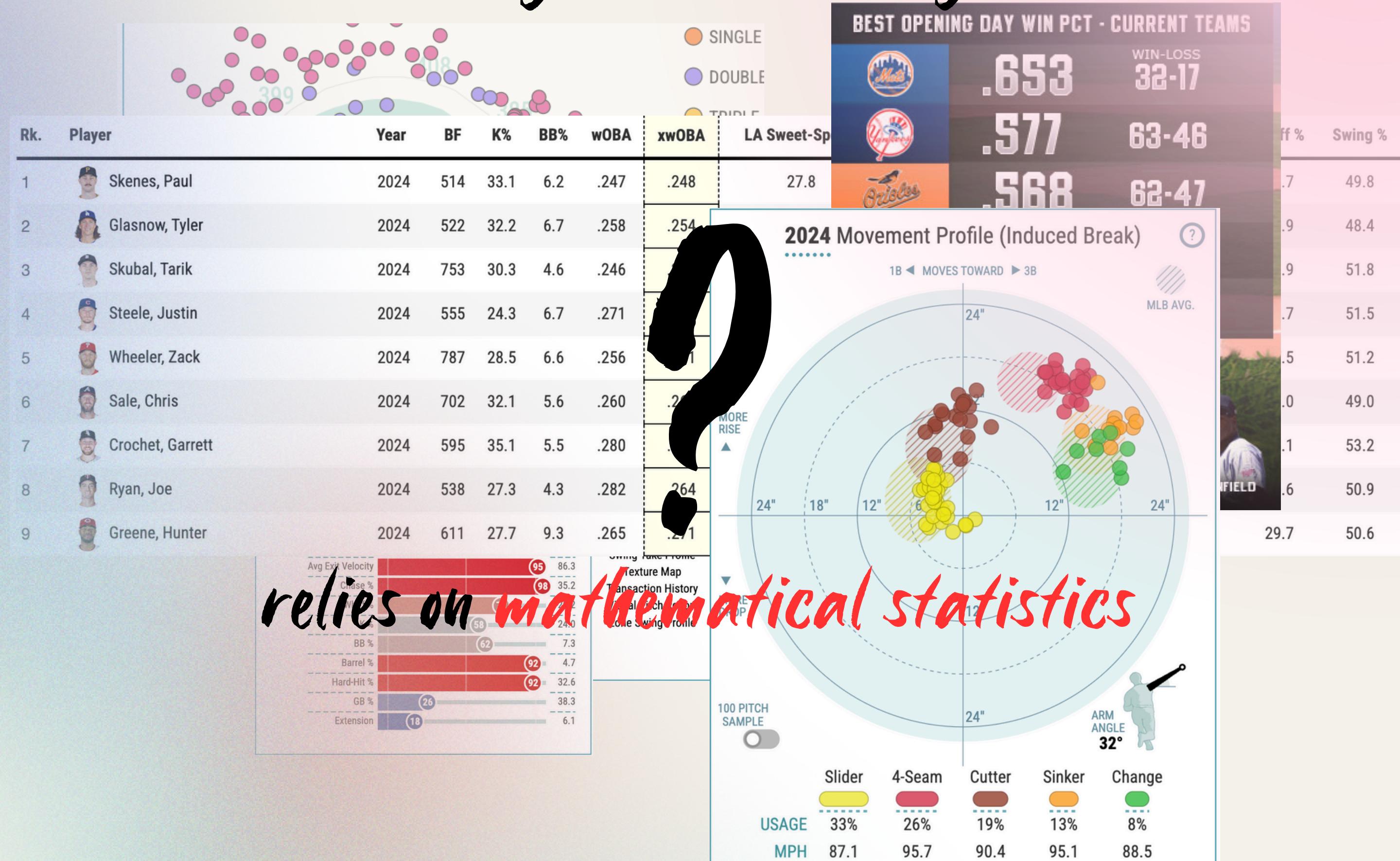


Using RNAOP model

# Outline

- Brief Introduction
- Question

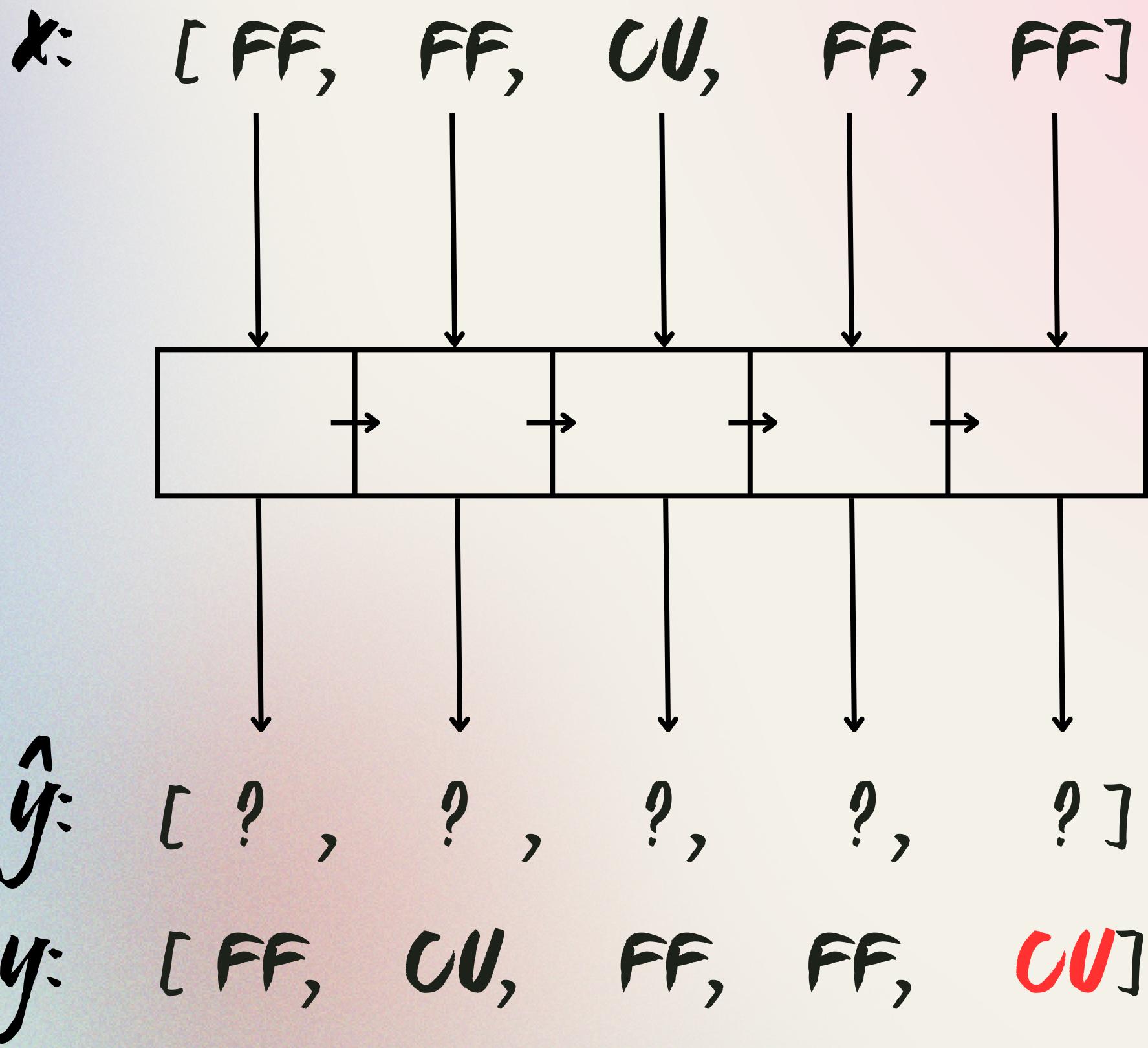
# Analytics Nowadays?



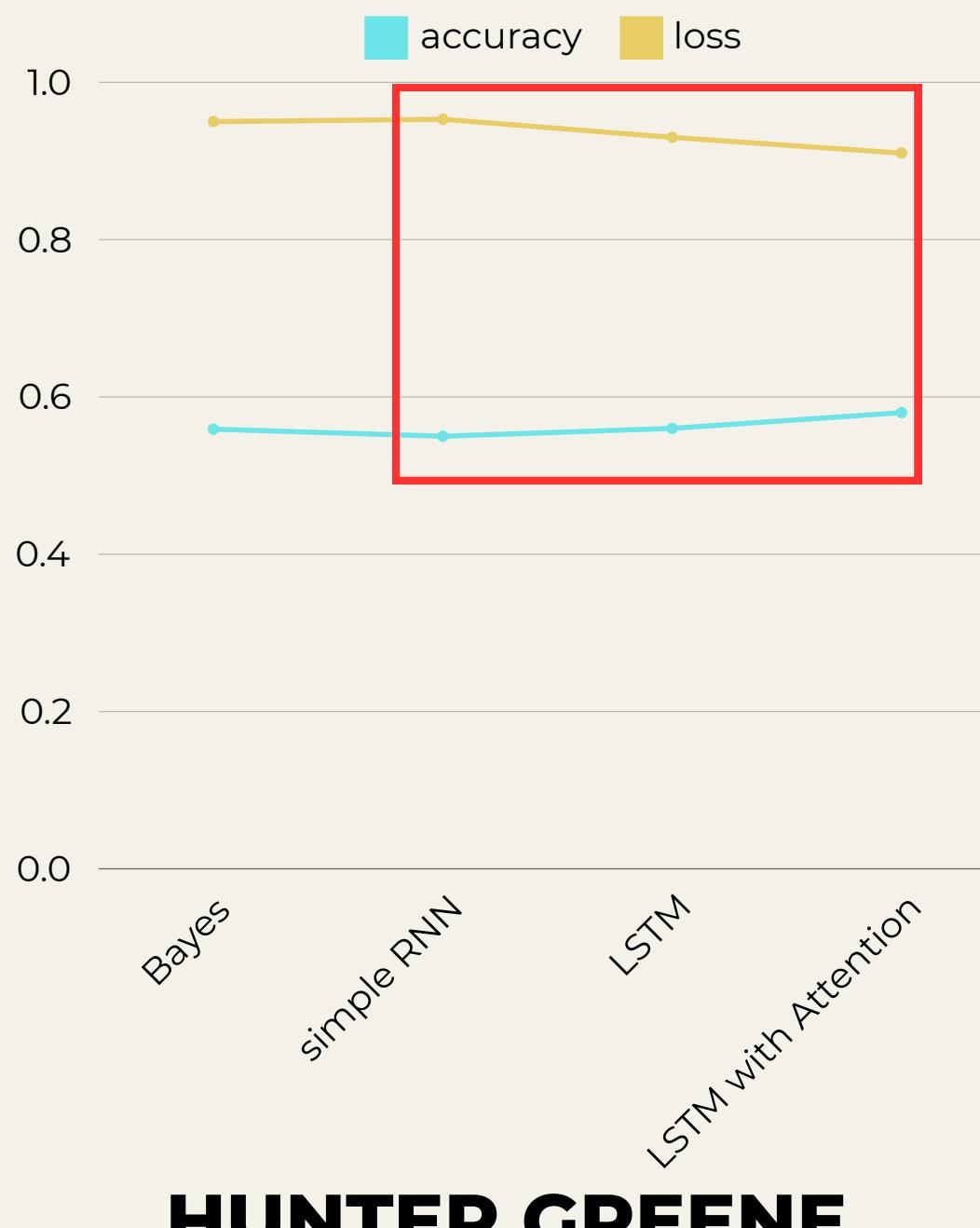
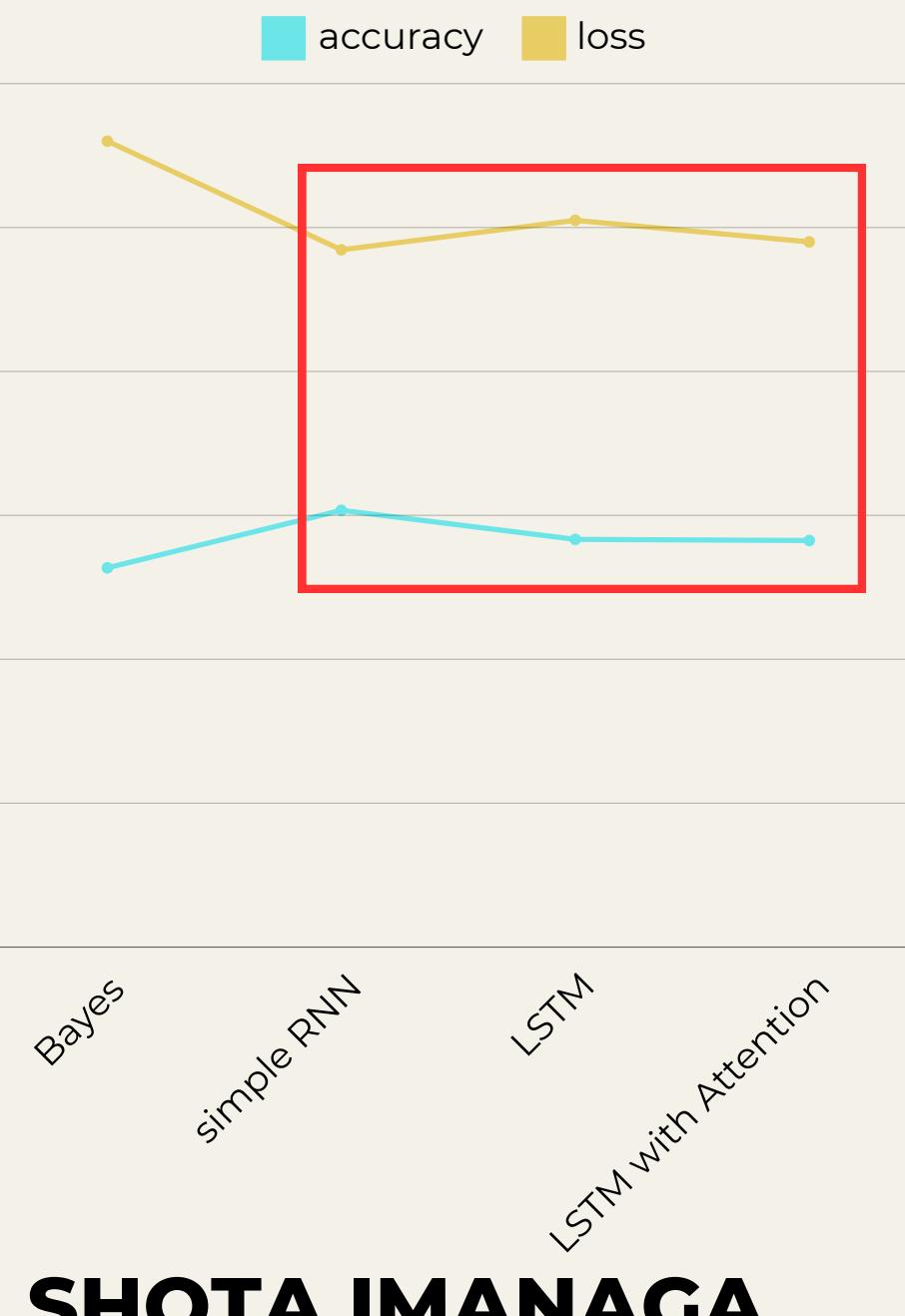
*“We all think that it may overlook game context and abstract features like pitcher-catcher synergy or strategy shifts based on the count.”*

# USE RNN MODEL TO LEARN THE RELATIONSHIP BETWEEN PITCHES

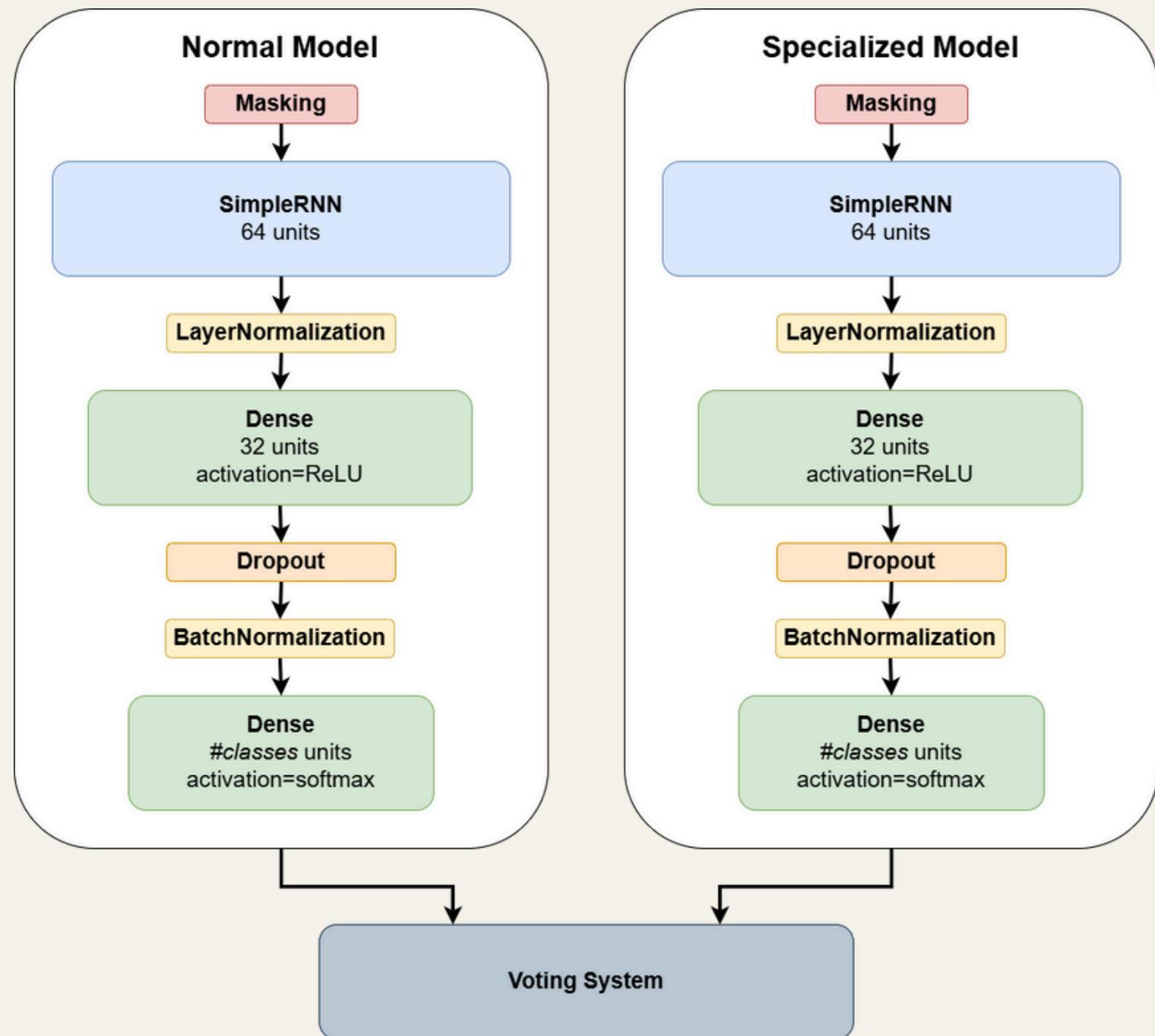
(and lots of other features...)



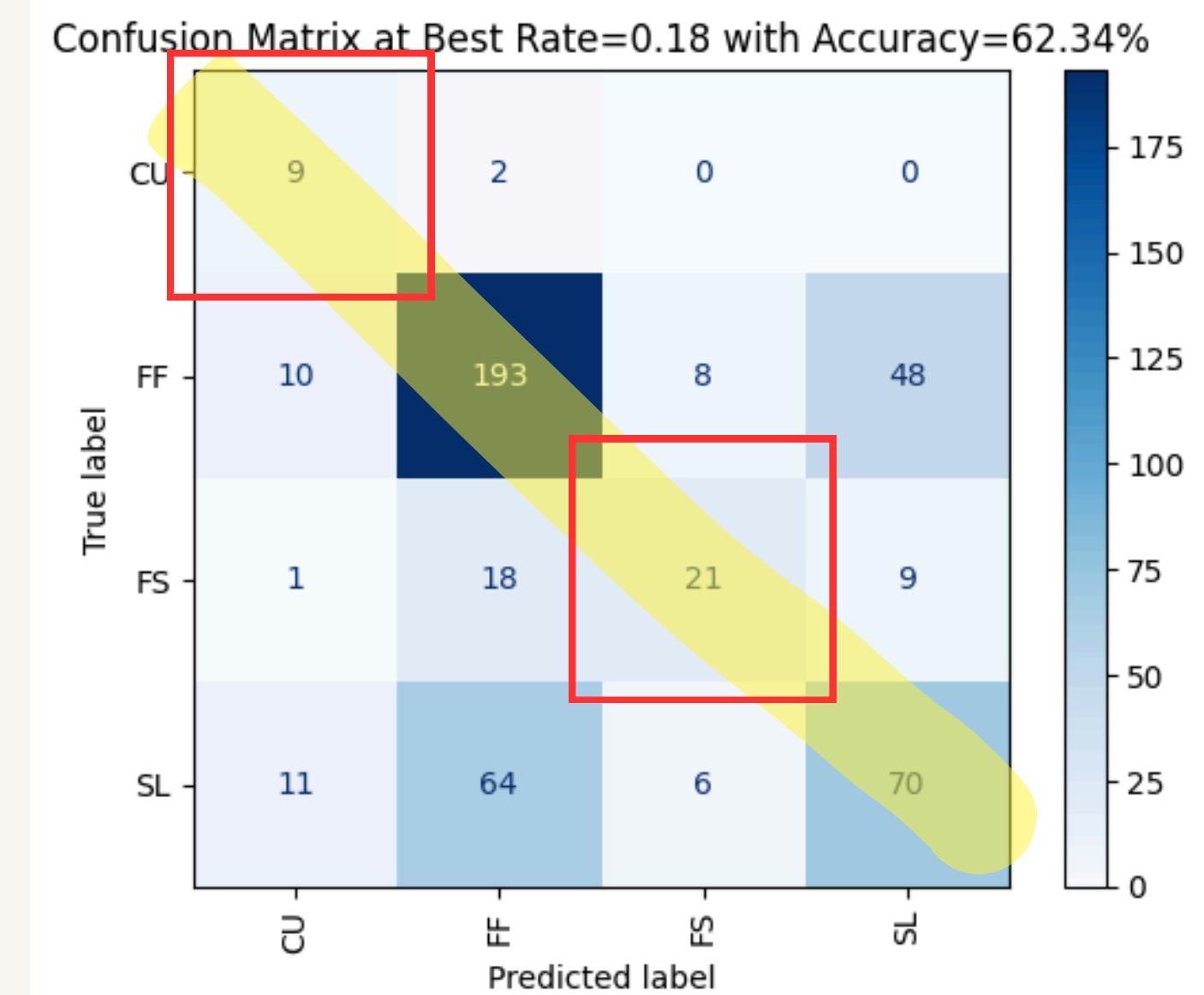
# After comparing four versions of model...



# FINAL VERSION MODEL HAS THE HIGHEST ACCURACY



## SHOTA IMANAGA



# **RELIEF PITCHER**

**RELIEF PITCHER**

**RELIEF PITCHER**



**MASON MILLER**

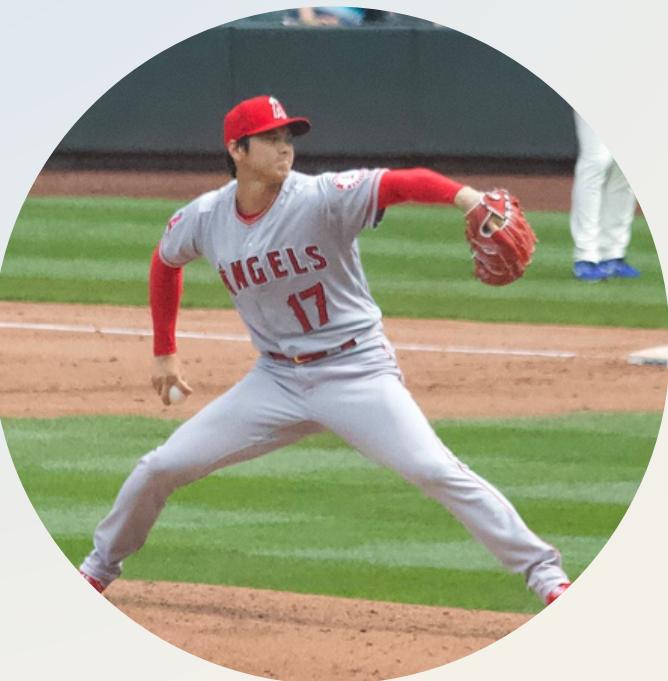
*Our accuracy of predicting pitch type in relief pitcher is about...*

**69%** , **55%**

ACCURACY WHICH WAY HIGHER THAN BASELINE PREDICTING BY BAYES CLASSIFIER

# Practical Value of our Model:

PROVIDES A POWERFUL TOOL TO ANALYZE THE PITCHING TENDENCIES



PITCHERS



COACHES



CATCHERS

# 1. Data Processing

Q1. I DIDN'T UNDERSTAND THE DATA SLICING PART TO ADD MORE DATA TO THE LIMITED ONES. CAN YOU EXPLAIN IT IN A MORE SPECIFIC WAY TO EXPLAIN THE TECHNIQUE YOU USED AND HOW IT INCREASES THE DATA SIZE?

# Data Slicing: Take an example

input: min = 2, max = 3

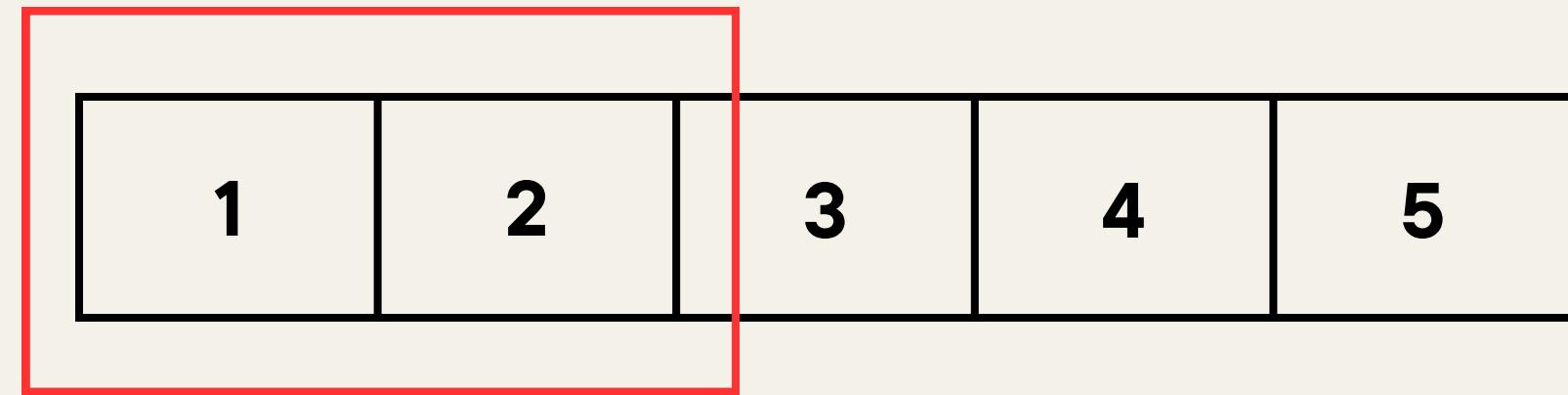
original sequence:

1	2	3	4	5
---	---	---	---	---

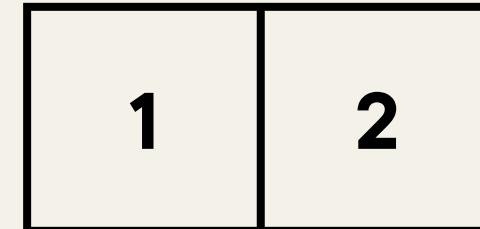
# Data Slicing

*length = 2*

original sequence:



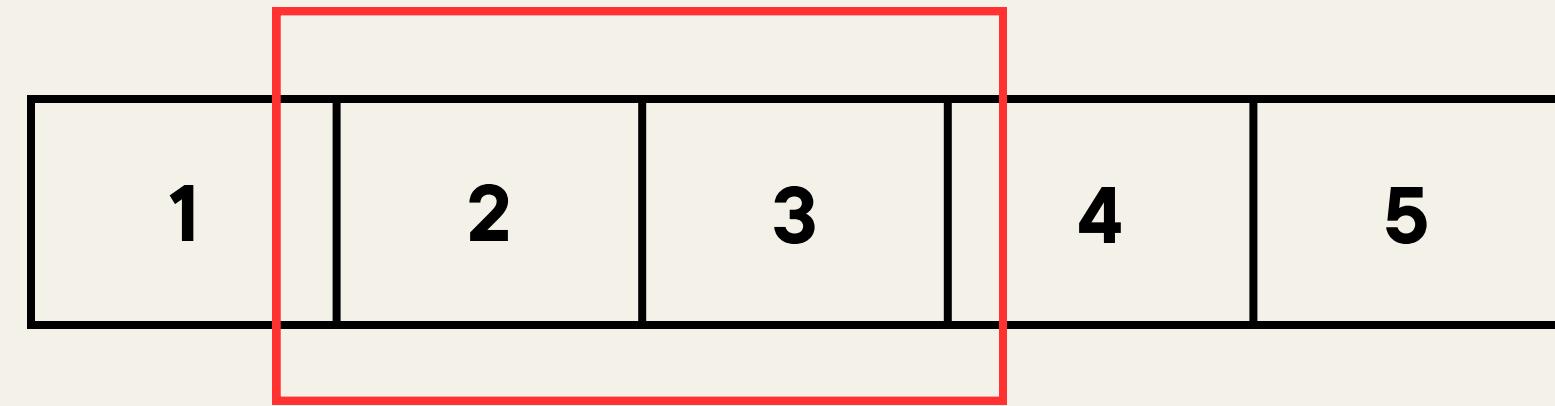
subsequences:



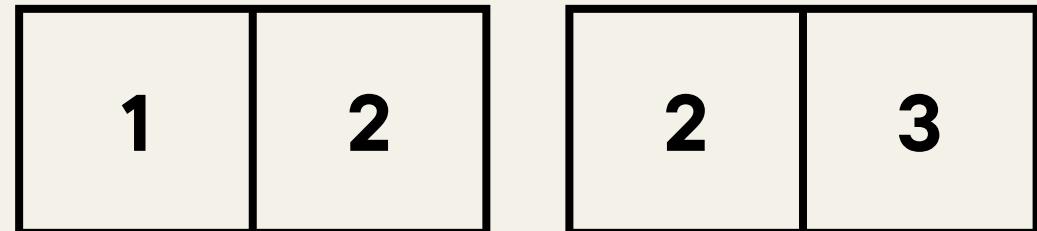
# Data Slicing

*length = 2*

original sequence:



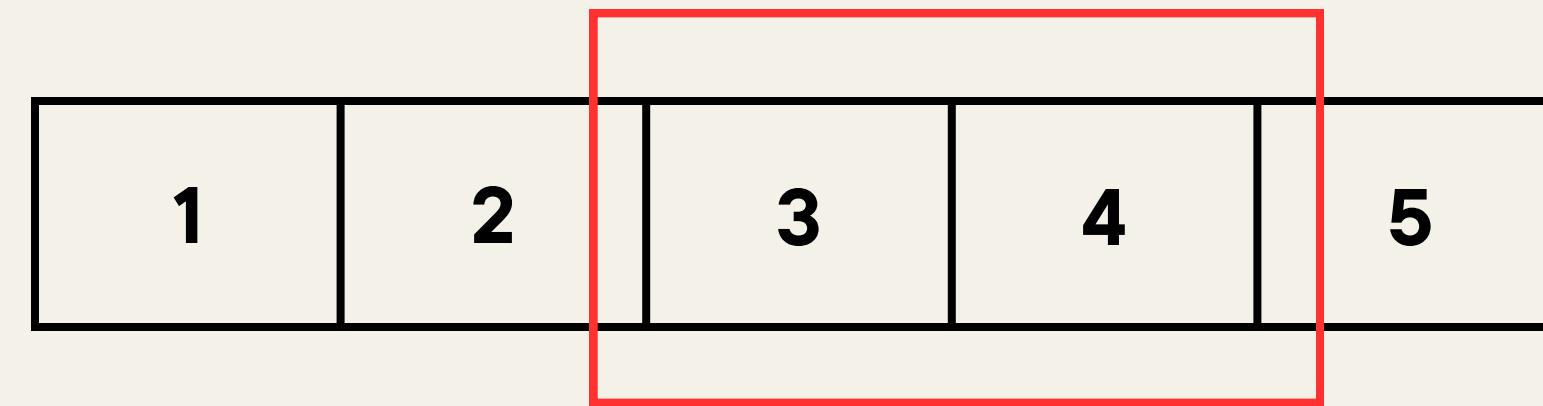
subsequences:



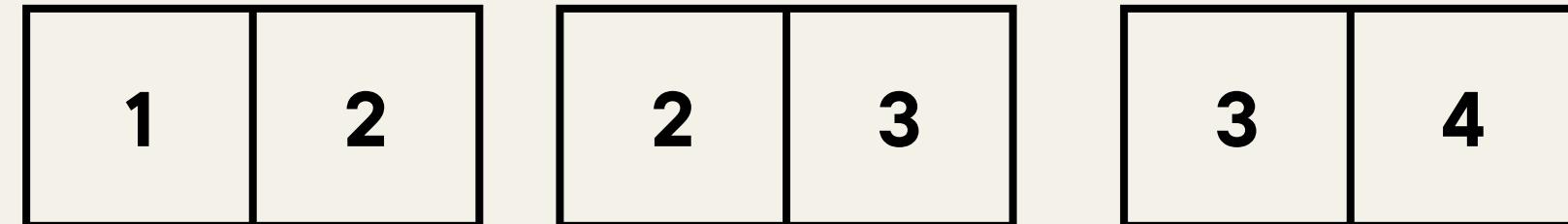
# Data Slicing

*length = 2*

original sequence:



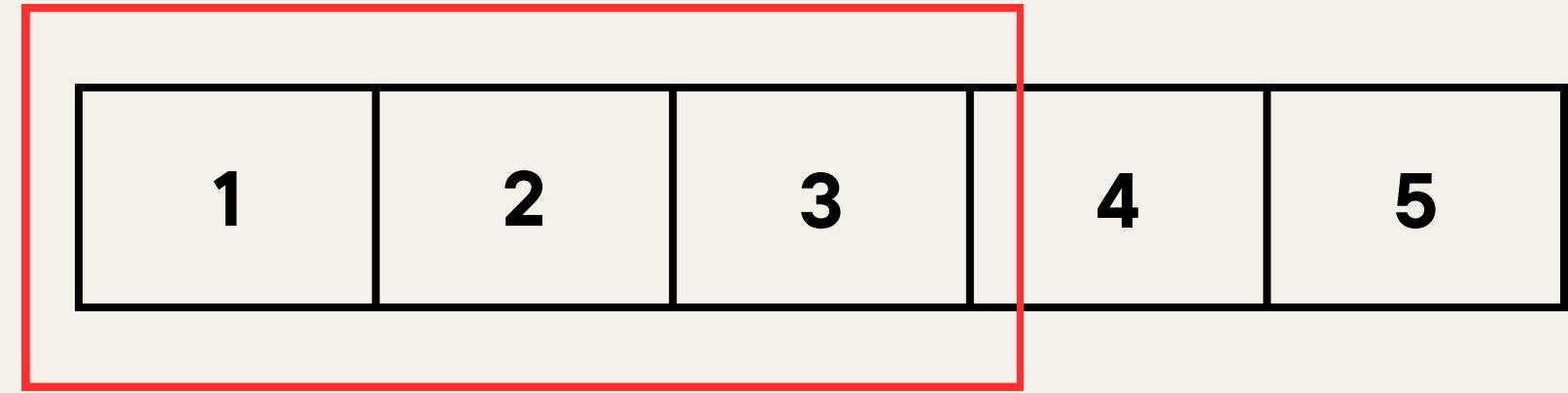
subsequences:



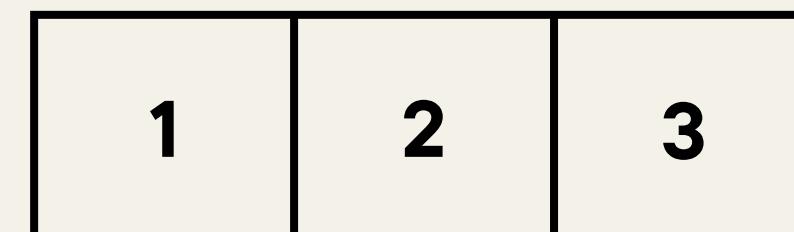
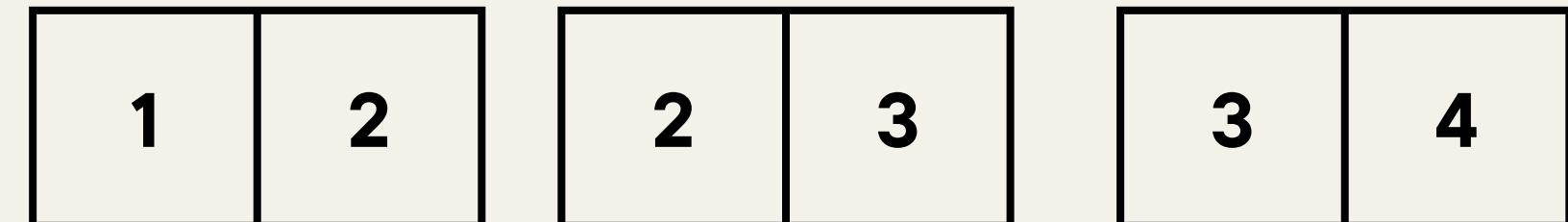
# Data Slicing

*length = 3*

original sequence:



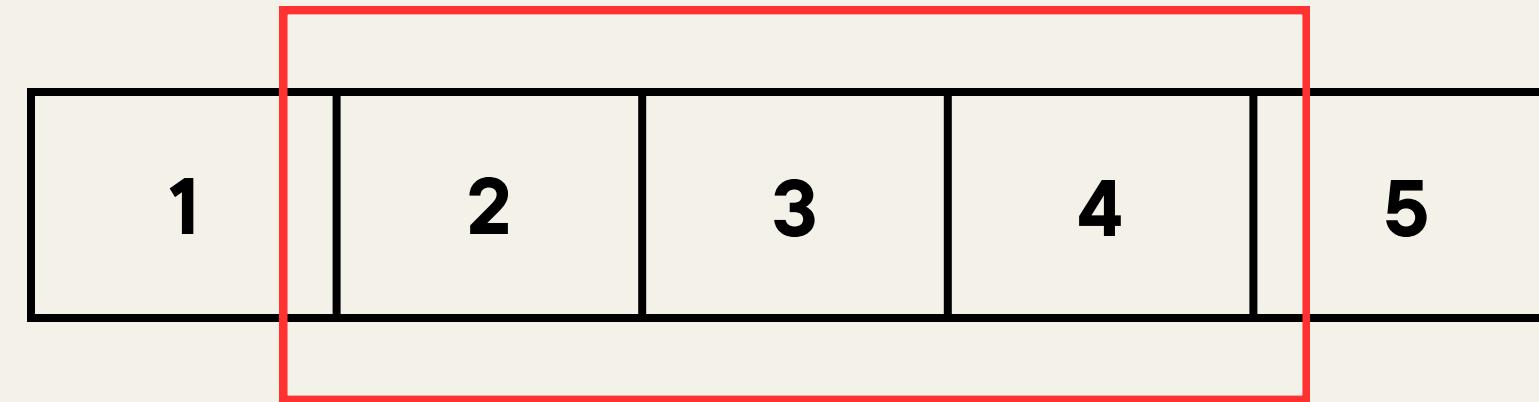
subsequences:



# Data Slicing

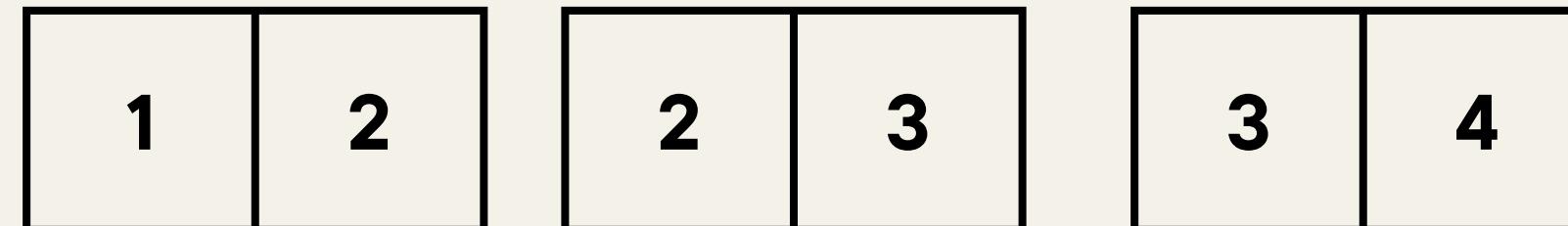
*length = 3*

original sequence:

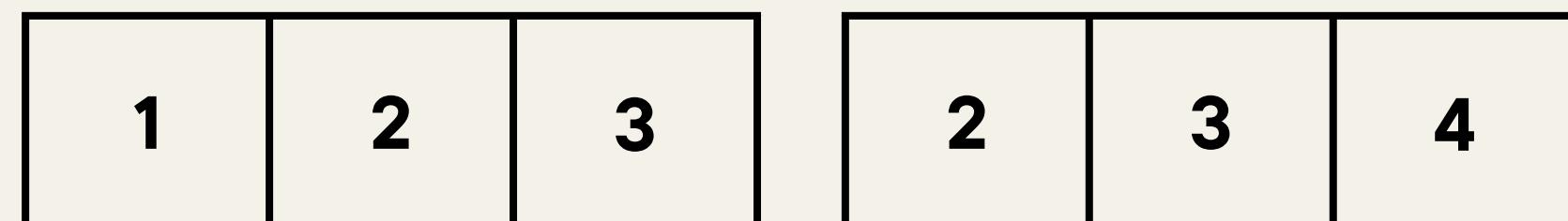


1

subsequences:



↓



S

*By doing this,*

TAKE THE CASE OF SHOTA IMANAGA AS EXAMPLE.

original sequences:

2590

->

19273

subsequences:

**Q2. WHY DON'T YOU USE DATA AUGMENTATION TO GENERATE ADDITIONAL DATA TO BALANCE THE DATA DISTRIBUTION?**

# Why not data augmentation

1. TRIED SMOTE-TS (SYNTHETIC MINORITY OVERSAMPLING  
TECHNIQUE FOR TIME SERIES).

- NO SIGNIFICANT IMPROVEMENTS.

2. DATA IS BASED ON MANY REAL-WORLD BASEBALL DYNAMIC  
ASPECTS.

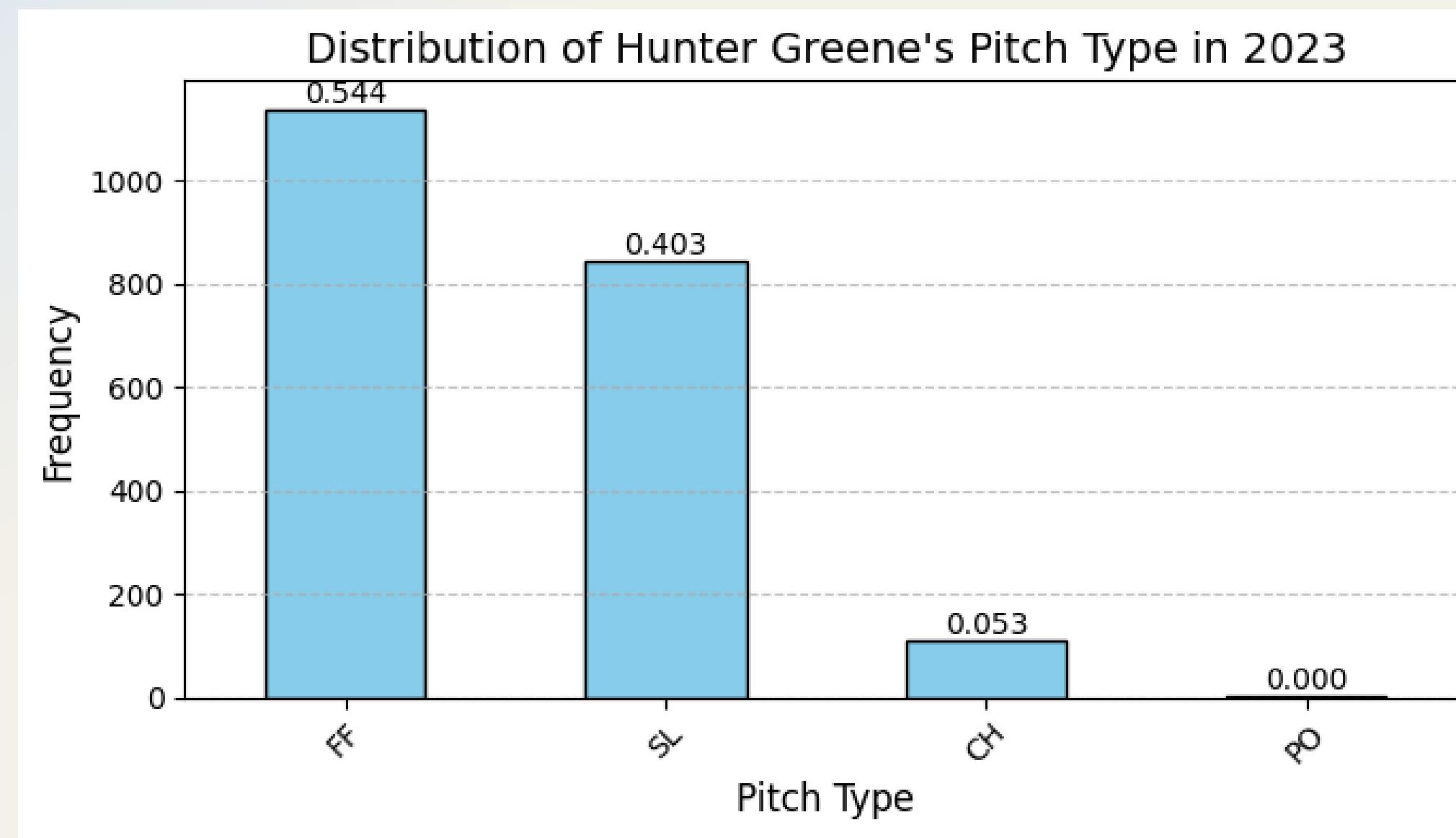
- NAIVE AUGMENTATION MIGHT LOSE THE RELATIONSHIP  
BETWEEN PITCHES.

Q3. YOU MENTIONED ADDING THE FREQUENCY FEATURE. I AM  
CURIOUS ABOUT WHEN YOU CALCULATE THE FREQUENCY—DO  
YOU CALCULATE IT AFTER SPLITTING THE DATA INTO TRAIN,  
VALIDATION, AND TEST SETS?

HOW DO YOU PREVENT A SITUATION WHERE THE DISTRIBUTIONS  
OF THE TRAIN AND VALIDATION SETS ARE SIGNIFICANTLY  
DIFFERENT?

# Frequency Feature?

- BASIC KNOWLEDGE OF THE PITCHER
- THE FREQUENCY DATA ARE FROM PREVIOUS YEAR



Q4. I NOTICED THAT THE LENGTH OF THE SUBSEQUENCES YOU USE FOR DATA SLICING RANGES FROM TWO TO SIX.

WOULD IT BE BETTER TO USE LONGER SEQUENCES FOR RNN TRAINING?

- WE HAVE TRIED TO CUT THE SEQUENCE BY BATTER AND BY DATE.
- THE SIZE OF THE SUBSEQUENCES: BATTER:2~6, DATE: 20 ~ 30.
- WE FIND OUT BATTER CUTTING SCHEME WITH SMALLER PERFORMS BETTER.

Q5. IN THE RAW DATA STRUCTURE, IS THERE ANY DATA COLUMN(FEATURE) ABOUT THE BATTER OR THE CATCHER, APART FROM THE BATTER'S ID? SINCE WE THINK THE BATTER AND THE CATCHER MIGHT INFLUENCE THE PITCH TYPE THROWN BY PITCHER.

Q6. IN THE CURRENT BASEBALL MATCHUP, THE PITCHER'S STRATEGY IS CLOSELY TIED TO THE CATCHER'S DECISIONS. HAVE YOU INCORPORATED THIS RELATIONSHIP INTO YOUR MODEL DURING TRAINING?

# jldbc/pybaseball

Pull current and historical baseball statistics using Python (Statcast, Baseball Reference, FanGraphs)



44 Contributors    392 Used by    1k Stars    339 Forks

*fielder\_2.1: catcher id*

```
Index(['pitch_type', 'game_date', 'release_speed', 'release_pos_x',
       'release_pos_z', 'player_name', 'batter', 'pitcher', 'events',
       'description', 'spin_dir', 'spin_rate_deprecated',
       'break_angle_deprecated', 'break_length_deprecated', 'zone', 'des',
       'game_type', 'stand', 'p_throws', 'home_team', 'away_team', 'type',
       'hit_location', 'bb_type', 'balls', 'strikes', 'game_year', 'px_x',
       'px_z', 'plate_x', 'plate_z', 'on_3b', 'on_2b', 'on_1b',
       'outs_when_up', 'inning', 'inning_topbot', 'hc_x', 'hc_y',
       'tfs_deprecated', 'tfs_zulu_deprecated', 'fielder_2', 'umpire', 'sv_id',
       'vx0', 'vy0', 'vz0', 'ax', 'ay', 'az', 'sz_top', 'sz_bot',
       'hit_distance_sc', 'launch_speed', 'launch_angle', 'effective_speed',
       'release_spin_rate', 'release_extension', 'game_pk', 'pitcher.1',
       'fielder_2.1', 'fielder_3', 'fielder_4', 'fielder_5', 'fielder_6',
       'fielder_7', 'fielder_8', 'fielder_9', 'release_pos_y',
       'estimated_ba_using_speedangle', 'estimated_woba_using_speedangle',
       'woba_value', 'woba_denom', 'babip_value', 'iso_value',
       'launch_speed_angle', 'at_bat_number', 'pitch_number', 'pitch_name',
       'home_score', 'away_score', 'bat_score', 'fld_score', 'post_away_score',
       'post_home_score', 'post_bat_score', 'post_fld_score',
       'if_fielding_alignment', 'of_fielding_alignment', 'spin_axis',
       'delta_home_win_exp', 'delta_run_exp'],
      dtype='object')
```

# jldbc/pybaseball

Pull current and historical baseball statistics using Python (Statcast, Baseball Reference, FanGraphs)



44 Contributors    392 Used by    1k Stars    339 Forks

## woba:

Weighted On Base Average  
(加權上壘率)

## babip:

Batting Average on Balls put  
Into Play (場內被安打率)

## iso:

Isolated Power (純長打率)

*batter id, hit\_distance...*

```
Index(['pitch_type', 'game_date', 'release speed', 'release_pos_x',
       'release_pos_z', 'player_name', 'batter', 'pitcher', 'events',
       'description', 'spin_dir', 'spin_rate_deprecated',
       'break_angle_deprecated', 'break_length_deprecated', 'zone', 'des',
       'game_type', 'stand', 'p_throws', 'home_team', 'away_team', 'type',
       'hit_location', 'bb_type', 'balls', 'strikes', 'game_year', 'px',
       'pxz', 'plate_x', 'plate_z', 'on_3b', 'on_2b', 'on_1b',
       'outs_when_up', 'inning', 'inning_topbot', 'hc_x', 'hc_y',
       'tfs_deprecated', 'tfs_zulu_deprecated', 'fielder_2', 'umpire', 'sv_id',
       'vx0', 'vy0', 'vz0', 'ax', 'ay', 'az', 'sz_top', 'sz_bot',
       'hit_distance_sc' 'launch_speed', 'launch_angle', 'effective_speed',
       'release_spin_rate', 'release_extension', 'game_pk', 'pitcher.1',
       'fielder_2.1', 'fielder_3', 'fielder_4', 'fielder_5', 'fielder_6',
       'fielder_7', 'fielder_8', 'fielder_9', 'release_pos_y',
       'estimated ba using speedangle', 'estimated woba using speedangle',
       'woba_value', 'woba_denom', 'babip_value', 'iso_value',
       'launch_speed_angle', 'at_bat_number', 'pitch_number', 'pitch_name',
       'home_score', 'away_score', 'bat_score', 'fld_score', 'post_away_score',
       'post_home_score', 'post_bat_score', 'post_fld_score',
       'if_fielding_alignment', 'of_fielding_alignment', 'spin_axis',
       'delta_home_win_exp', 'delta_run_exp'],
      dtype='object')
```

# 2 Model Design

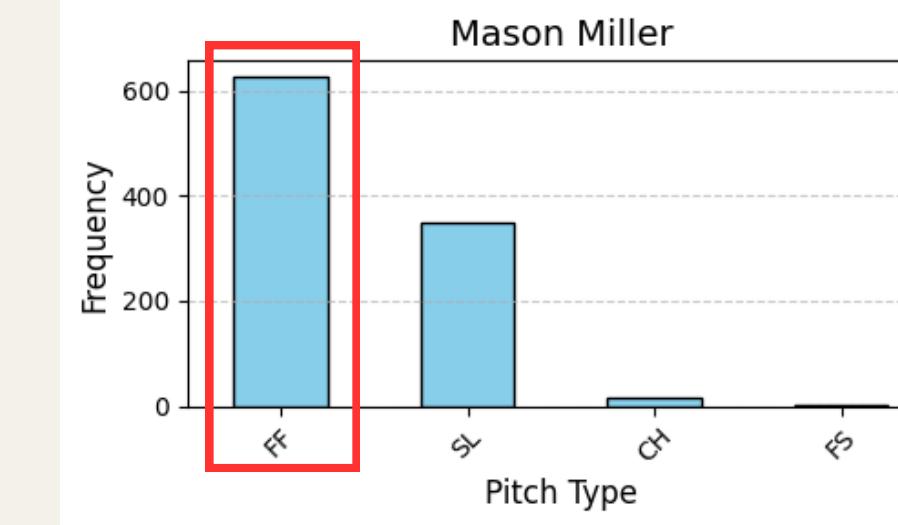
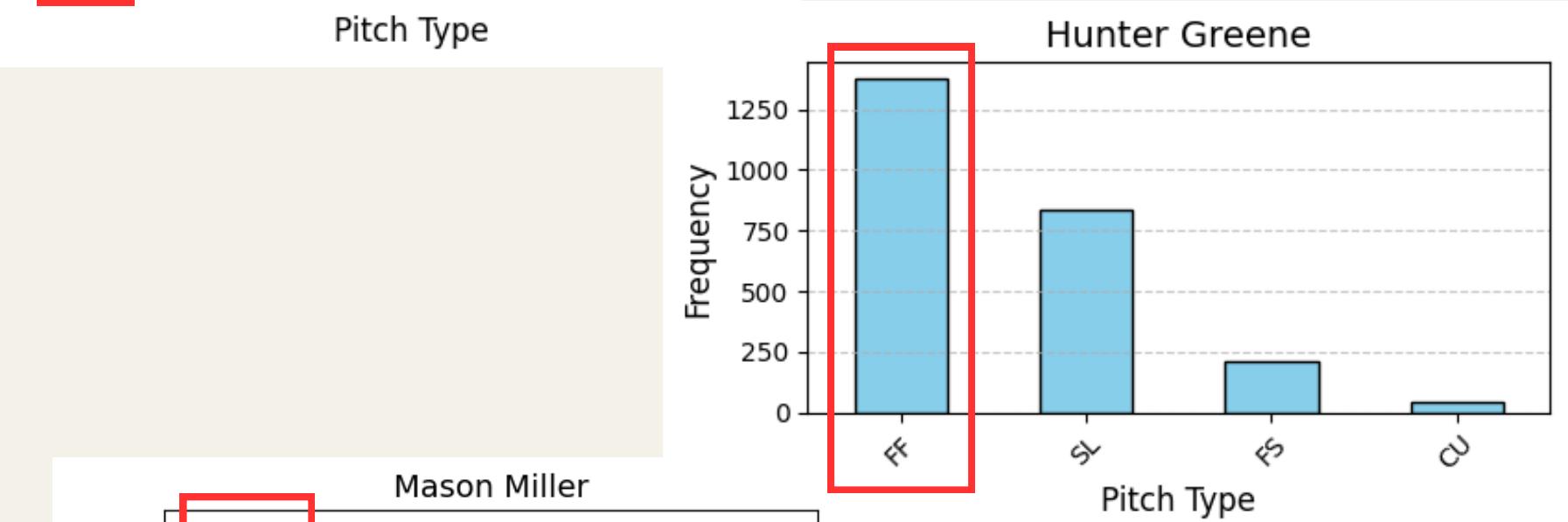
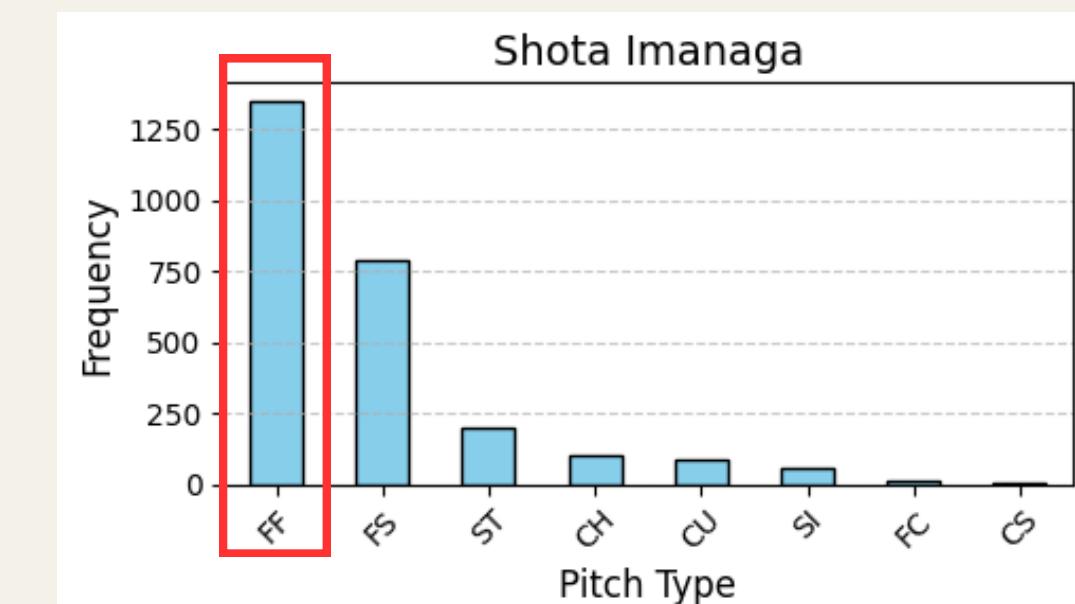
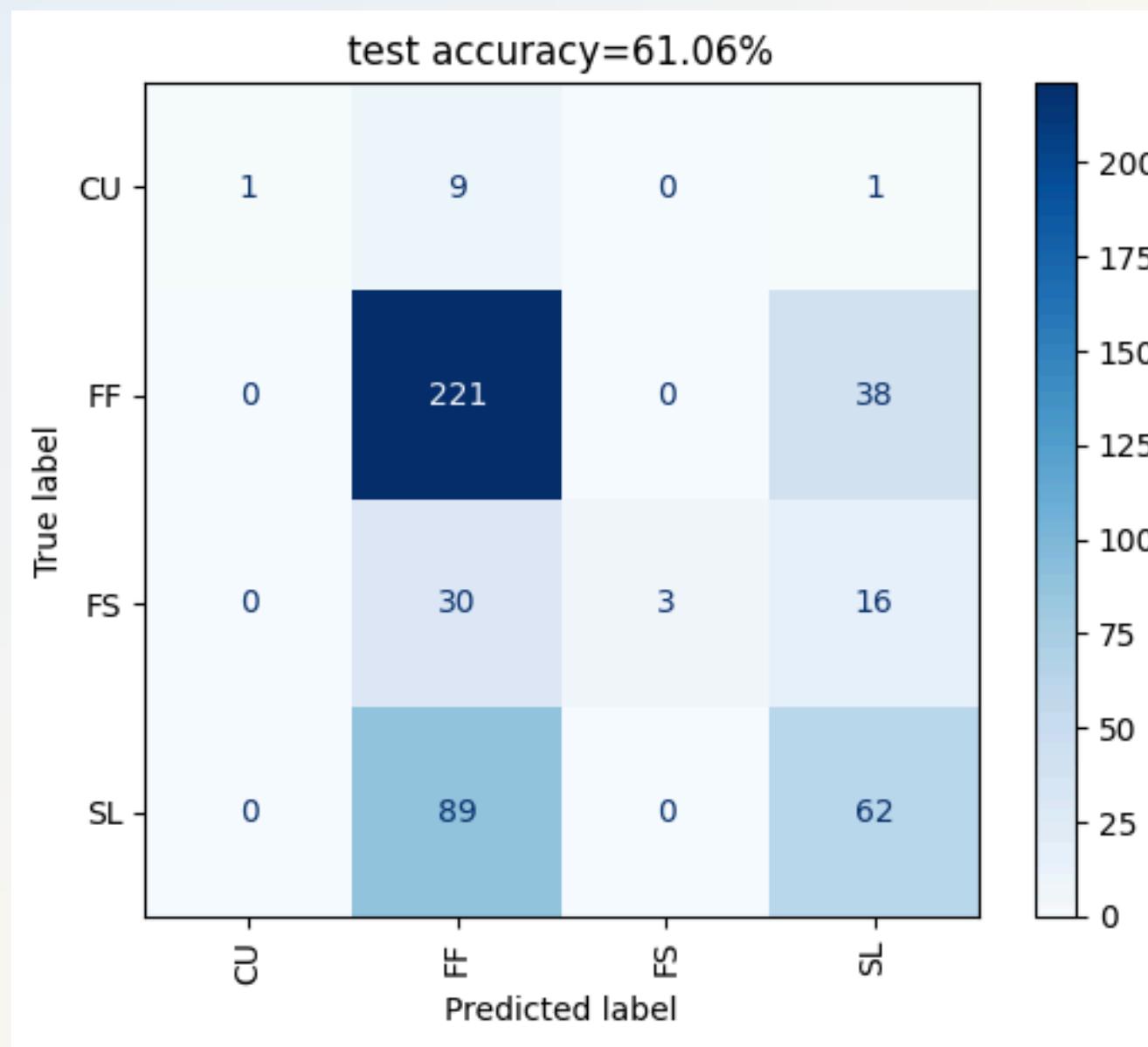
Q7. IN THE DISCUSSION PART, YOU NOTICE THAT THE PREDICTION OF THE NORMAL MODEL FOCUSES ONLY ON FASTBALLS, AND YOU DECIDE TO ADD THE SPECIALIST MODEL TO NORMAL RNN ONES AND CREATE A VOTING SYSTEM THAT SHOWS A SLIGHTLY BETTER PERFORMANCE.

Q7. (CONT.) I ALSO NOTICED IN THAT PART, IN WHICH YOU COMPARE THE ORIGINAL ONES WITH THE NEW ONES (WITH THE SPECIALIST MODEL), THAT THE TWO CONFUSION MATRICES DIDN'T HAVE THE SAME AMOUNT OF GROUND TRUTH, I BELIEVE THAT ADDING THE SPECIALIST MODEL DOES GAIN SOME IMPROVEMENT BUT MAYBE COMPARE IT WITH THE SAME AMOUNT OF DATA THAT CAN GET A CLEARER VISION OF THE IMPROVEMENT?

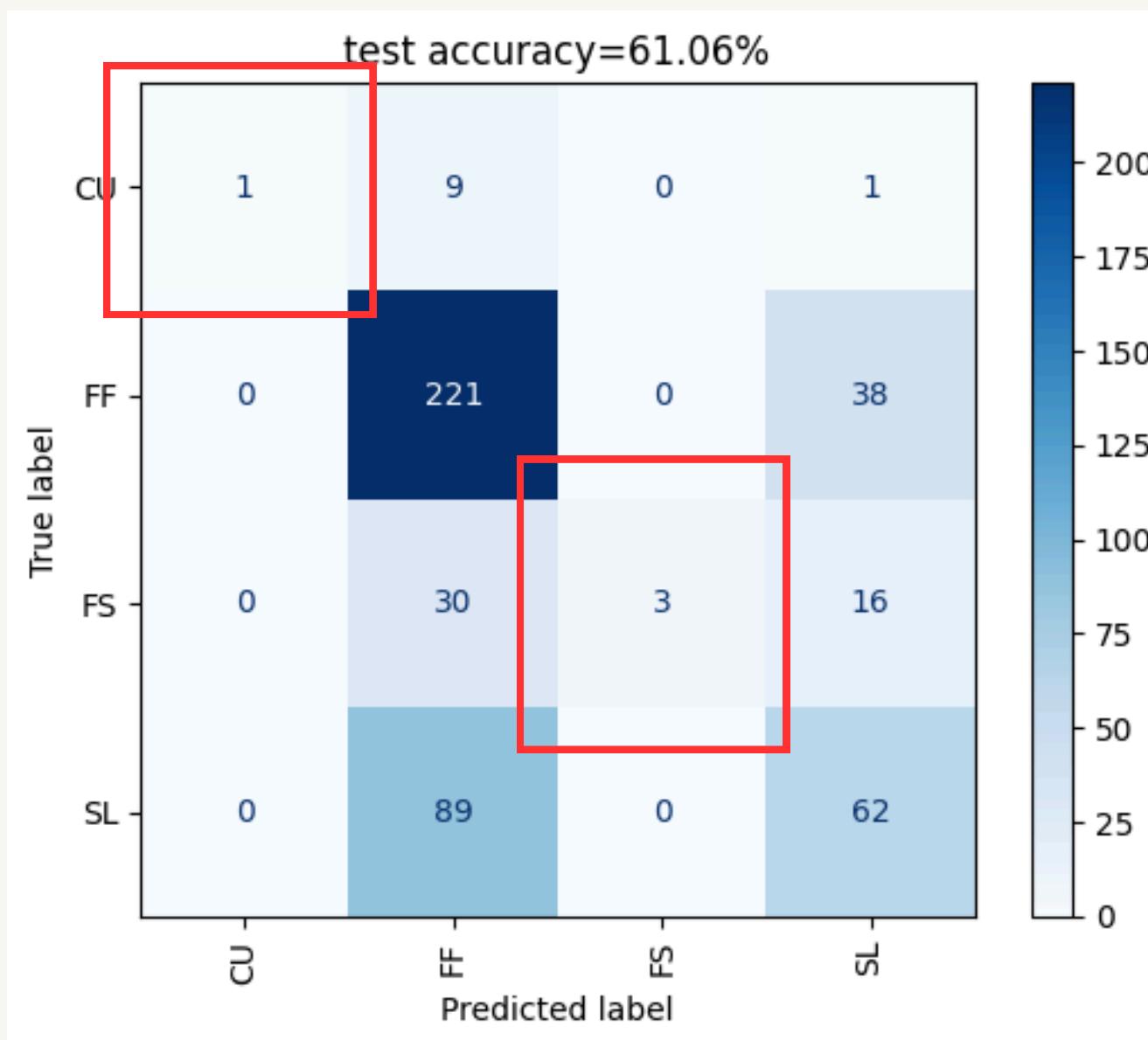
# Accuracy and Loss Are Not Enough:

## Confusion Matrix

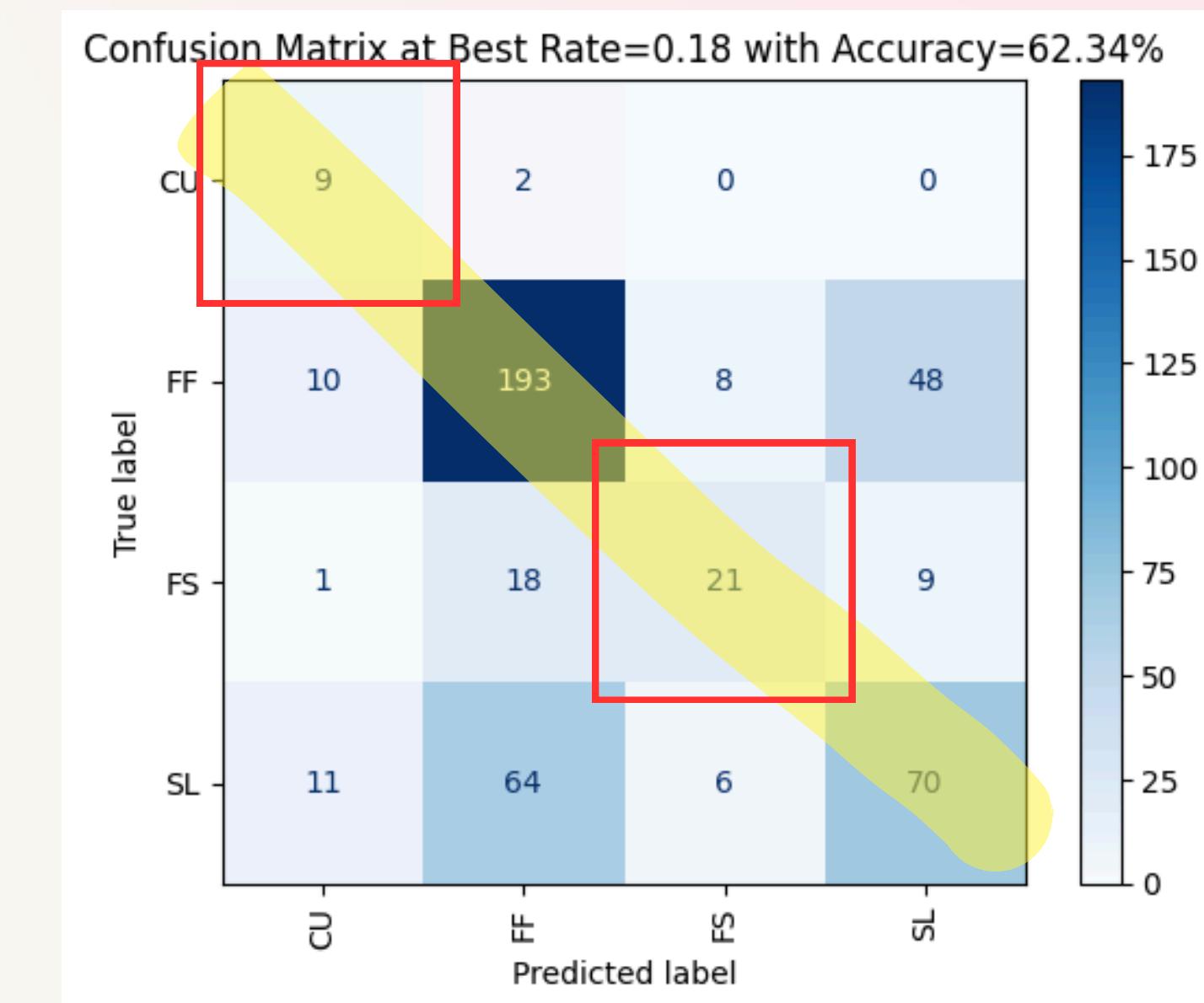
HUNTER GREENE: 61.06%



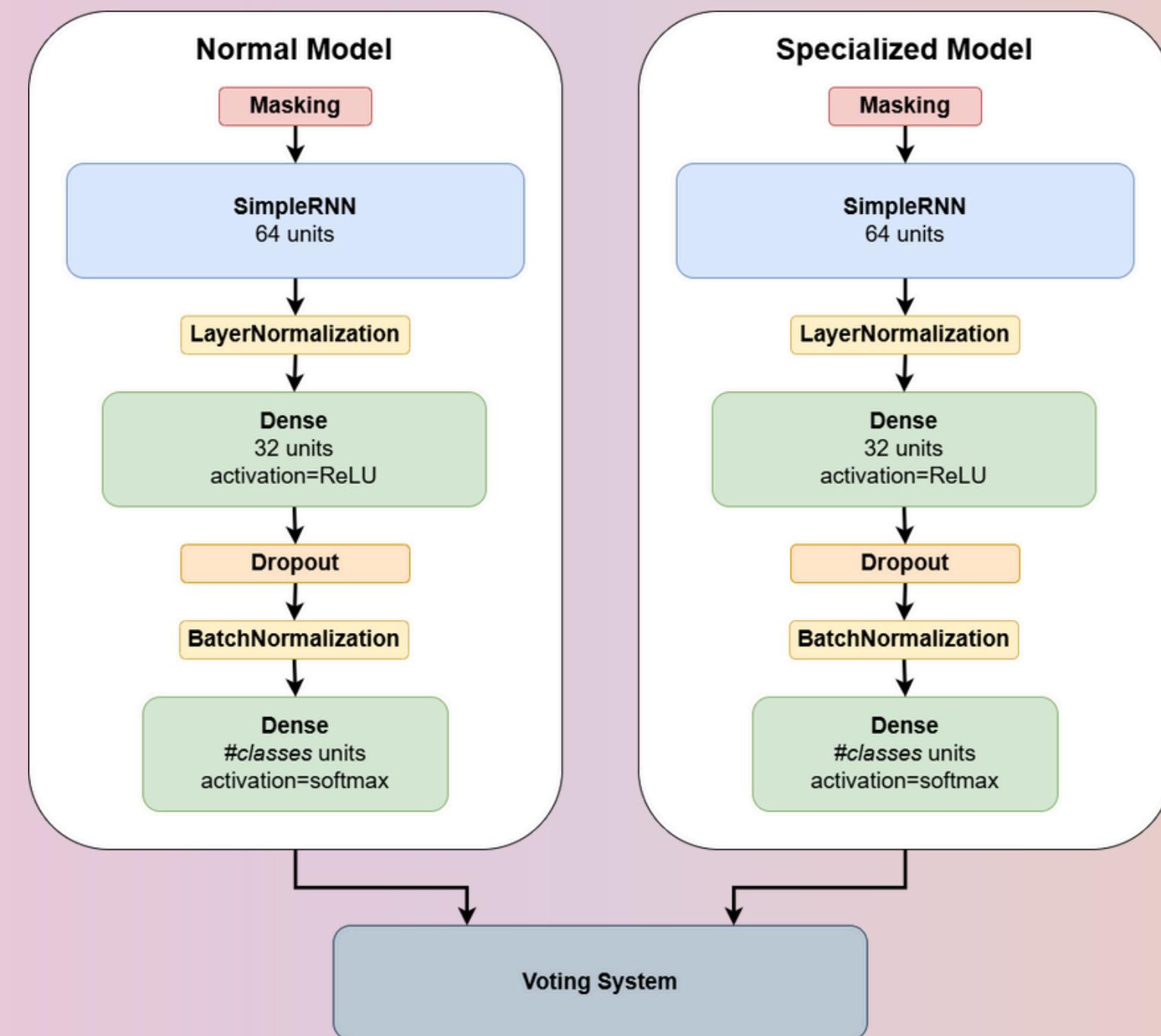
# SIMPLE RNN: 61.06%



# ADDING SPECIALIST MODEL: 62.34%



# QB. CAN YOU EXPLAIN WHY YOU CHOOSE ONLY TWO BASE LEARNERS IN A VOTING SYSTEM WITH WEIGHTED AVERAGING ?



# Why 2 Base Learners in Voting System

1. THE WORD "VOTING" HERE IS A BIT INFORMAL. A MORE PRECISE TERM WOULD BE "ENSEMBLE".
2. WE CHOSE TWO BASE LEARNERS FOR GENERALIZABILITY
  - ONE FOR ALL PITCH TYPES
  - ONE FOR LESS COMMON PITCH TYPES (EXCLUDING THE 2 MOST FREQUENT FS AND CU)

# 3. Result Analysis

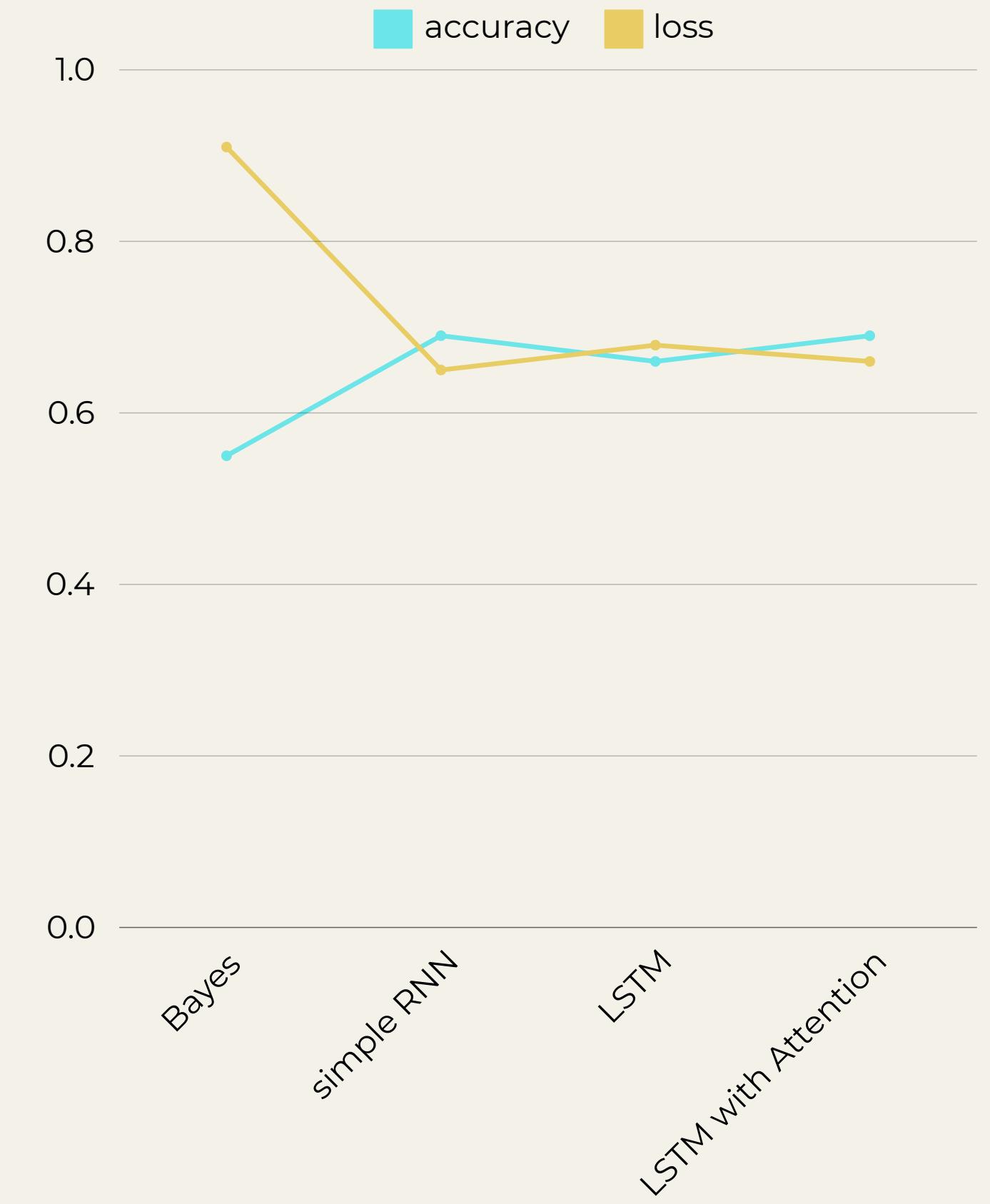
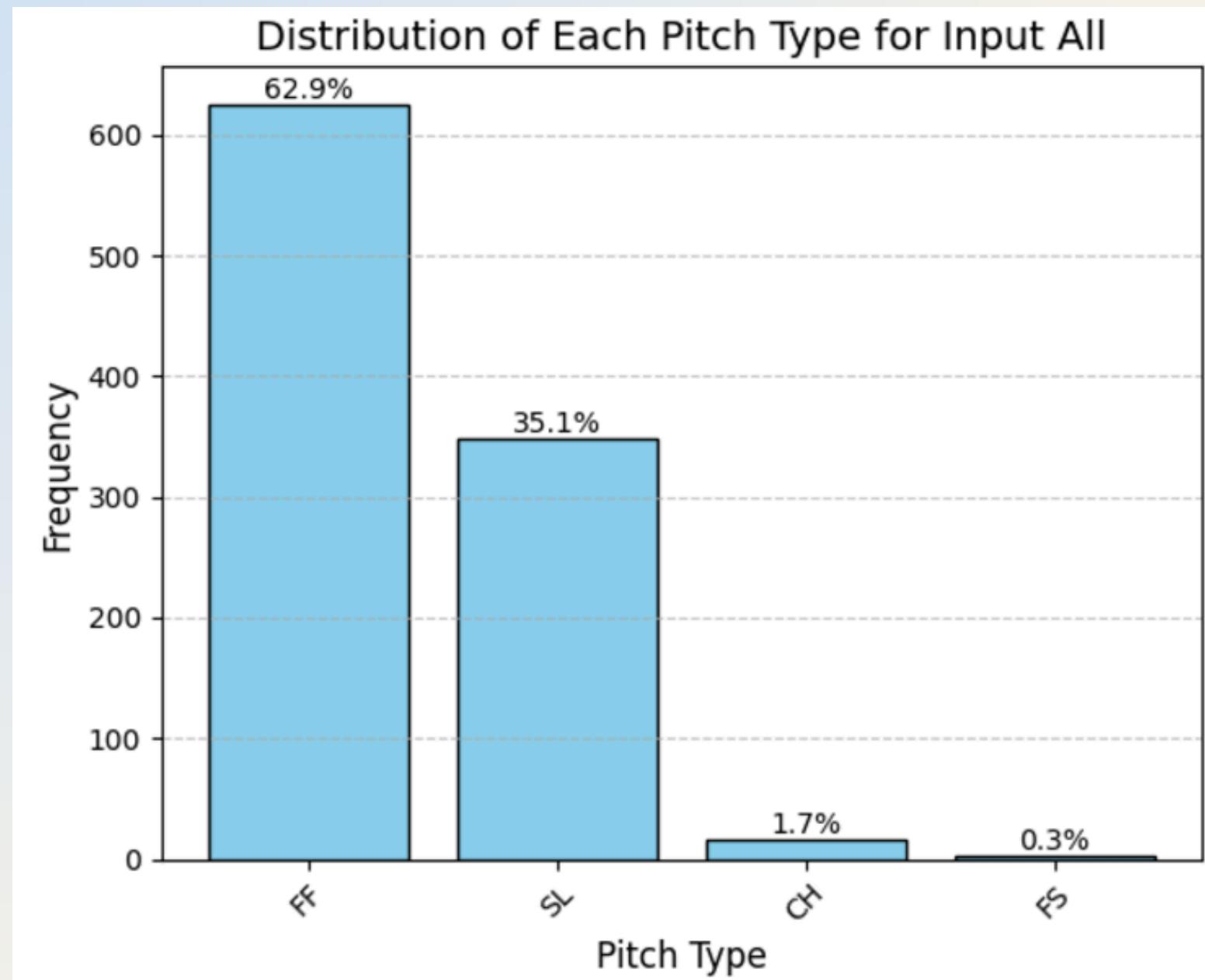
Q9. CAN YOU EXPLAIN AGAIN WHY YOU THINK THE ONES WITH FEWER PITCH TYPES HAVE BETTER RESULTS, AND ALSO IN THE VIDEO, I NOTICED THAT THE MILLER ONE HAS BETTER RESULTS THAN THE OTHER TWO, COULD YOU ALSO DISCUSS THIS PART TO SHOW THE POSSIBLE REASON THAT MILLER'S STAND OUT?

WHEN A PITCHER HAS A LIMITED VARIETY OF PITCH TYPES, THE MODEL HAS FEWER CATEGORIES TO DISTINGUISH, LEADING TO HIGHER ACCURACY.

MILLER'S BETTER RESULTS CAN BE ATTRIBUTED TO HIS FEWER PITCH TYPES, WHICH GENERALLY IMPROVES ACCURACY.

**Q10. IS 60% A GOOD ACCURACY COMPARED TO THE ACCURACY  
OF EASIER GUESSING METHODS?**

# MASON MILLER



THANKS FOR LISTENING

TEAM23