

TEAM 23

Pitch Type Prediction

Baseball

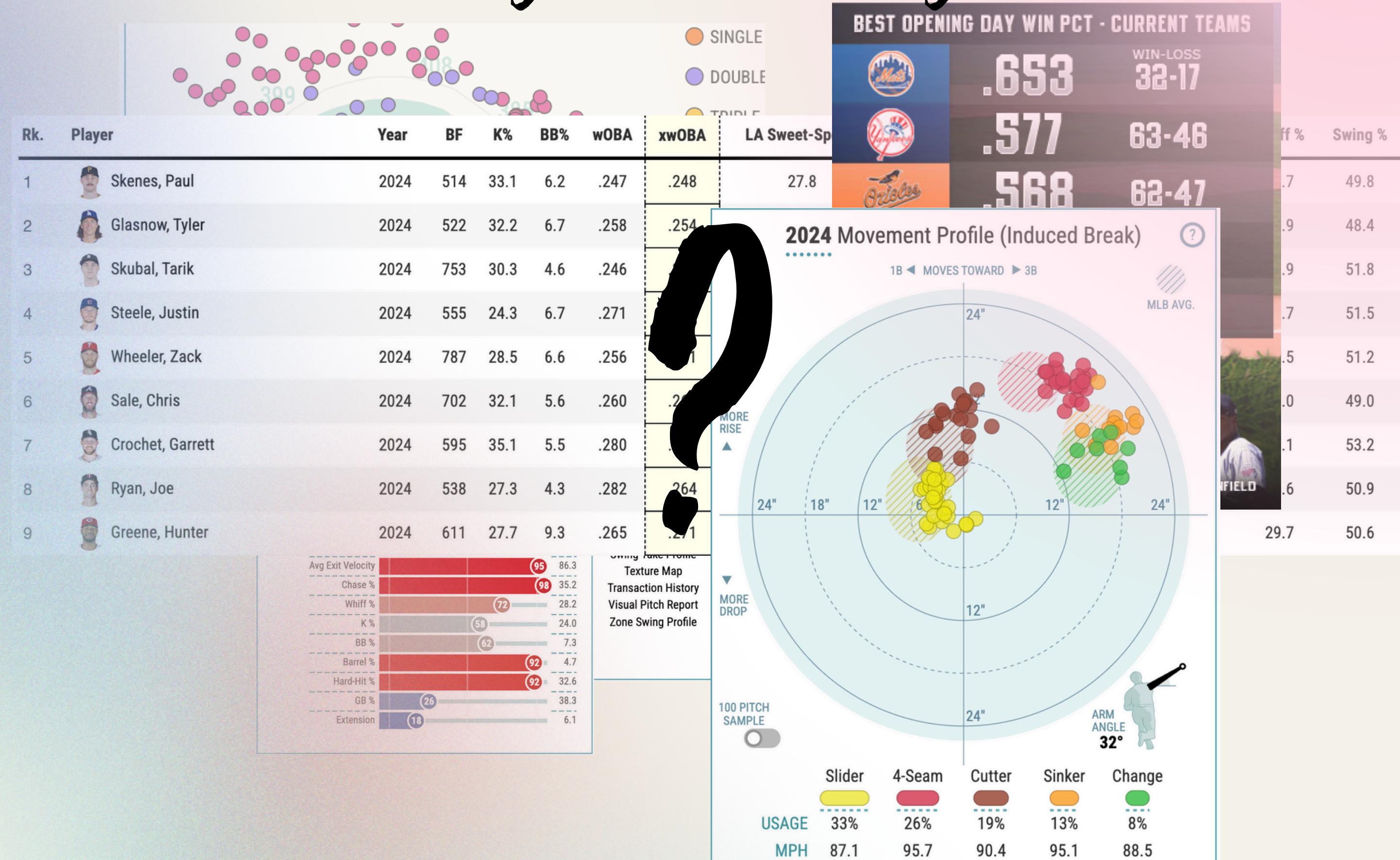


Using RNAOP model

Outline

- Introduction and Background
- Methodology
- Experiments
- Results
- Discussion and Conclusion

Analytics Nowadays ?



1. Data Processing

jldbc/pybaseball

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



44
Contributors

392
Used by

1k
Stars

339
Forks

Package Documentation

Player-Specific Queries

For a player-specific statcast query, pull pitching or batting data using the `statcast_pitcher` and `statcast_batter` functions. These take the same `start_dt` and `end_dt` arguments as the `statcast` function, as well as a `player_id` argument. This ID comes from MLB Advanced Media, and can be obtained using the function `playerid_lookup`. The returned columns match the set above, but filtered to rows for that specific pitcher or batter. A complete example:

```
# Find Clayton Kershaw's player id
from pybaseball import playerid_lookup
from pybaseball import statcast_pitcher
playerid_lookup('kershaw', 'clayton')
   name_last name_first key_mlbam key_retro key_bbref key_fangraphs mlb_played_first mlt
0    kershaw     clayton      477132  kersc001  kershcl01          2036        2008.0

# His MLBAM ID is 477132, so we feed that as the player_id argument to the following function
kershaw_stats = statcast_pitcher('2017-06-01', '2017-07-01', 477132)
kershaw_stats.groupby("pitch_type").release_speed.agg("mean")
pitch_type
CH      86.725000
CU      73.133333
FF      92.844622
SI      94.515385
SL      87.962381
Name: release_speed, dtype: float64
```

(Input: name, date)



```
1 input_last_name = "shota"
2 input_first_name = "imanaga"
3 input_start_date = '2024-04-01'
4 input_end_date = '2024-11-02'
5 # Look up player ID
6 pitcher_id = playerid_lookup(input_first_name, input_last_name).key_mlbam[0]
7
8 # Output path for the CSV
9 output_root = 'input_all_data_2024.csv'
10
11 # Fetch all data for the pitcher in the given date range
12 data = statcast_pitcher(input_start_date, input_end_date, player_id=pitcher_id)
13
14 # Save the complete dataset to a CSV file
15 data.to_csv(output_root, index=False, encoding='utf-8')
16
17 print(f"All data successfully saved to {output_root}!")
```

Raw data structure

```
→ pitch_type game_date release_speed release_pos_x release_pos_z \
0 FF 2024-04-01 92.4 2.36 5.68
1 FF 2024-04-01 93.5 2.19 5.69
2 FF 2024-04-01 92.7 2.35 5.59
3 FF 2024-04-01 93.3 2.38 5.67
4 FS 2024-04-01 84.1 2.50 5.43
...
2585 FF 2024-09-22 90.1 2.54 5.52
2586 FF 2024-09-22 91.3 2.51 5.49
2587 SI 2024-09-22 88.3 2.68 5.32
2588 FF 2024-09-22 90.9 2.60 5.44
2589 FS 2024-09-22 83.0 2.64 5.36

player_name batter pitcher events \
0 Imanaga, Shota 453568 684007 NaN
1 Imanaga, Shota 453568 684007 NaN
2 Imanaga, Shota 453568 684007 field_out
3 Imanaga, Shota 663898 684007 NaN
4 Imanaga, Shota 663898 684007 NaN
...
2585 Imanaga, Shota 702358 684007 ...
2586 Imanaga, Shota 702358 684007 NaN
2587 Imanaga, Shota 696285 684007 NaN
2588 Imanaga, Shota 696285 684007 NaN
2589 Imanaga, Shota 696285 684007 grounded_into_double_play

description ... n_thruorder_pitcher \
0 called_strike ... 1
1 ball ... 1
2 hit_into_play ... 1
3 foul ... 1
4 swinging_strike ... 1
...
2585 ball ... 3
2586 swinging_strike ... 3
```



```
1 data = pd.read_csv('input_all_data_2024.csv')
2 data = data.iloc[:-1].reset_index(drop=True)
3
4 print(data)
```

Methodology

Data Preprocessing

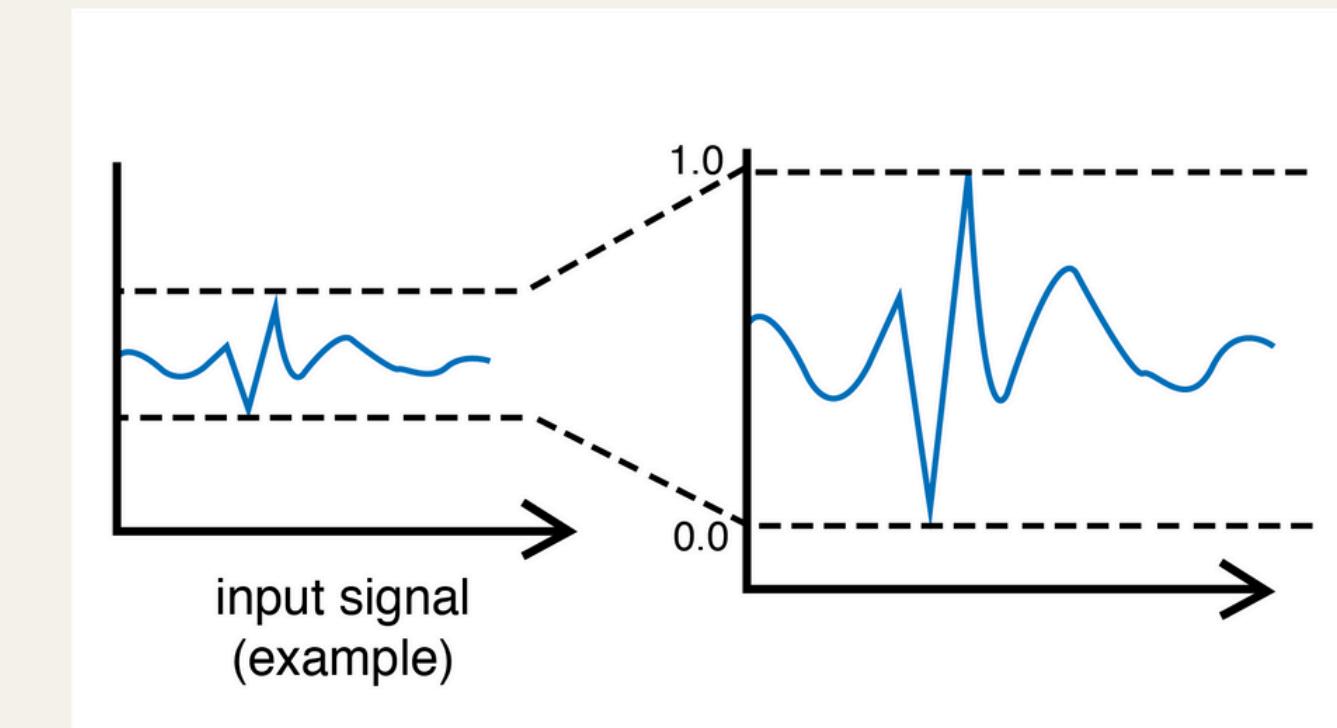
STEP 1: NORMALIZATION

SOME OF NUMERICAL FEATURES WERE SCALED TO A RANGE [-1, 1]
USING A MINMAXSCALER.

RELEASE_SPEED

LAUNCH_ANGLE

HIT_DISTANCE



Data Preprocessing

STEP 2: ONE-HOT ENCODING

CATEGORICAL FEATURES WERE CONVERTED INTO BINARY VECTORS
USING ONE-HOT ENCODING.

PITCH_TYPE

STRIKES

BALLS

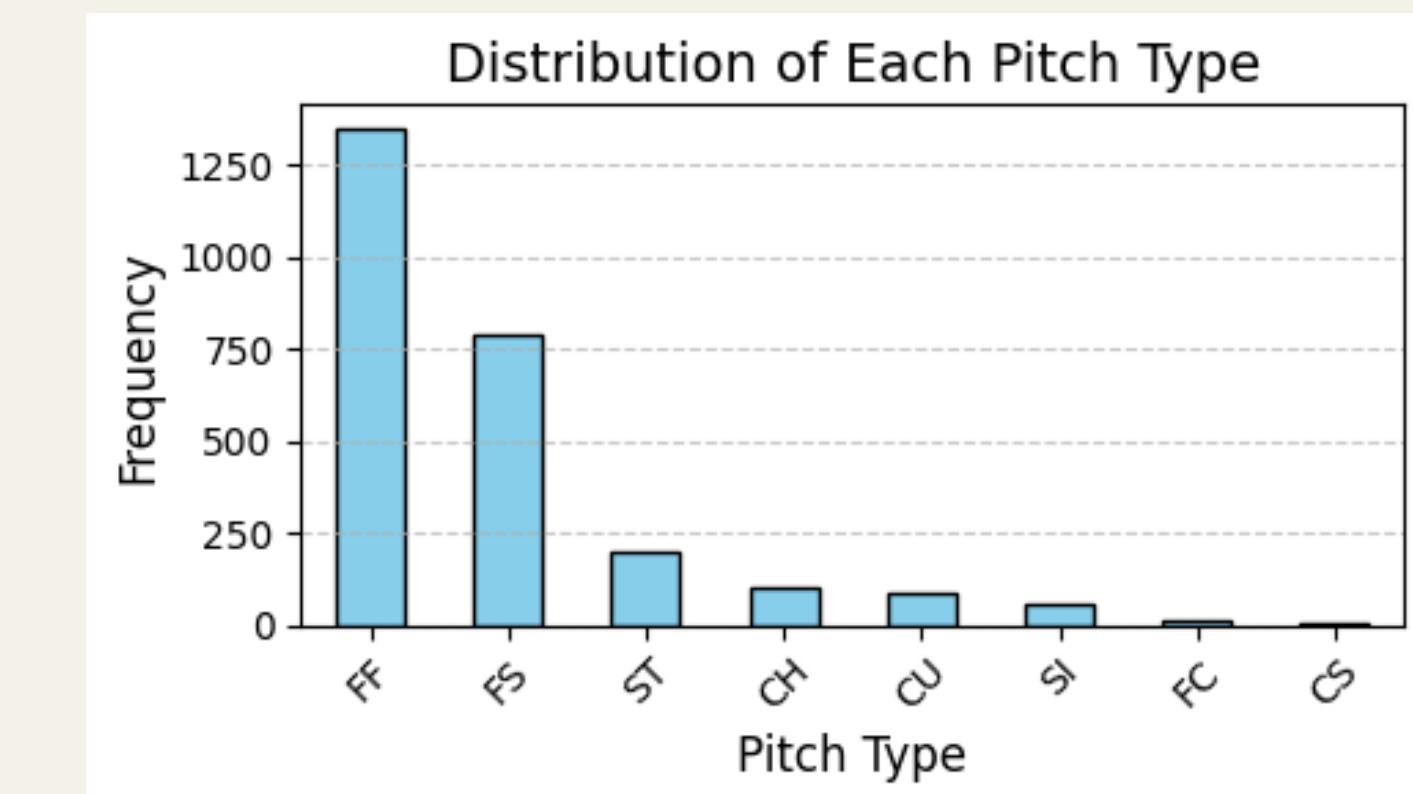
ZONE

Feature	Original Value	One-Hot Encoding
pitch_type	Fastball	[1, 0, 0, 0]
strikes	2	[0, 0, 1]
zone	3	[0, 0, 1, 0, 0]

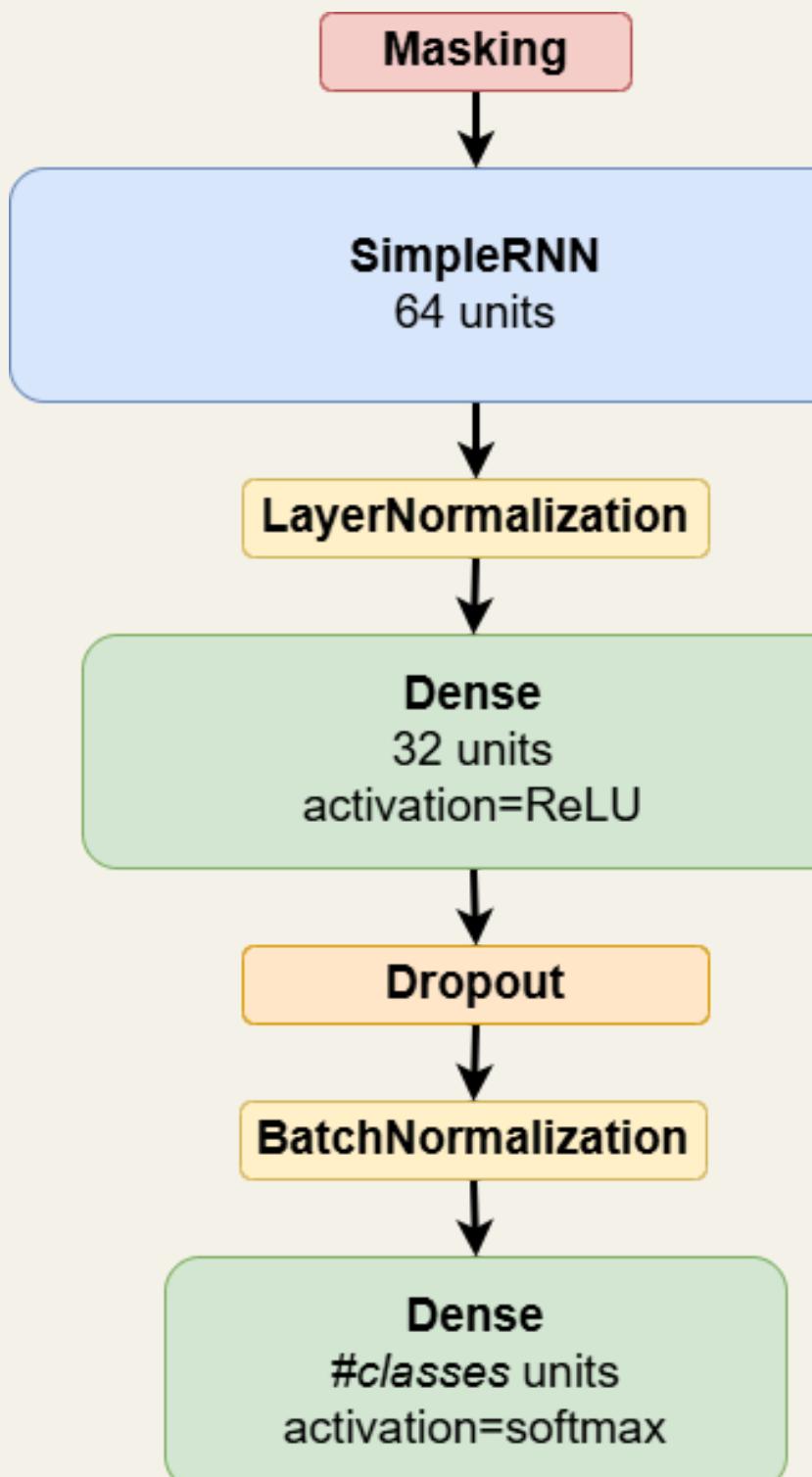
Data Preprocessing

STEP 3: ADDING FREQUENCY FEATURES

ADDED ADDITIONAL FEATURES BASED ON THE FREQUENCY OF
DIFFERENT PITCH TYPES.

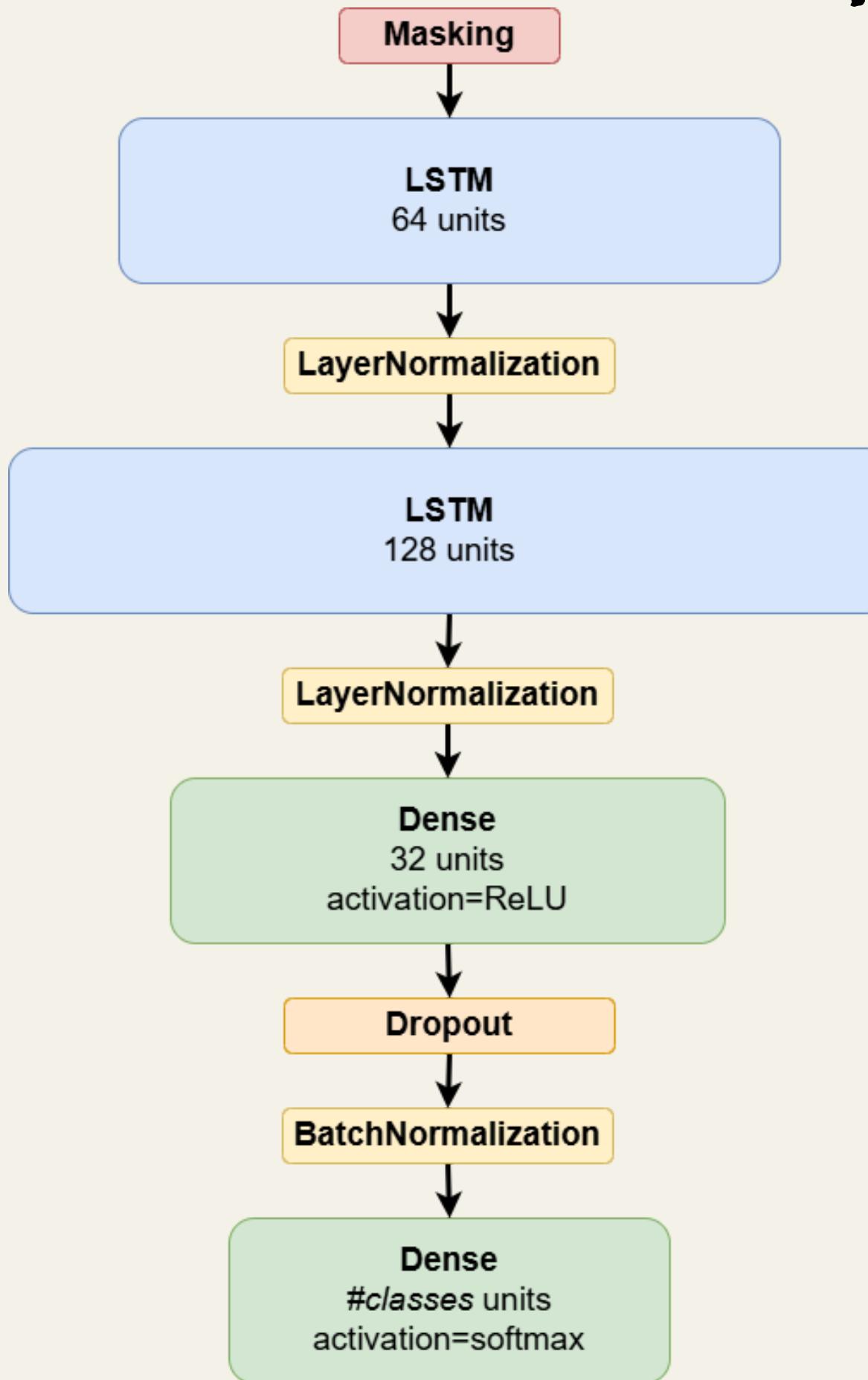


Model Design: RNN



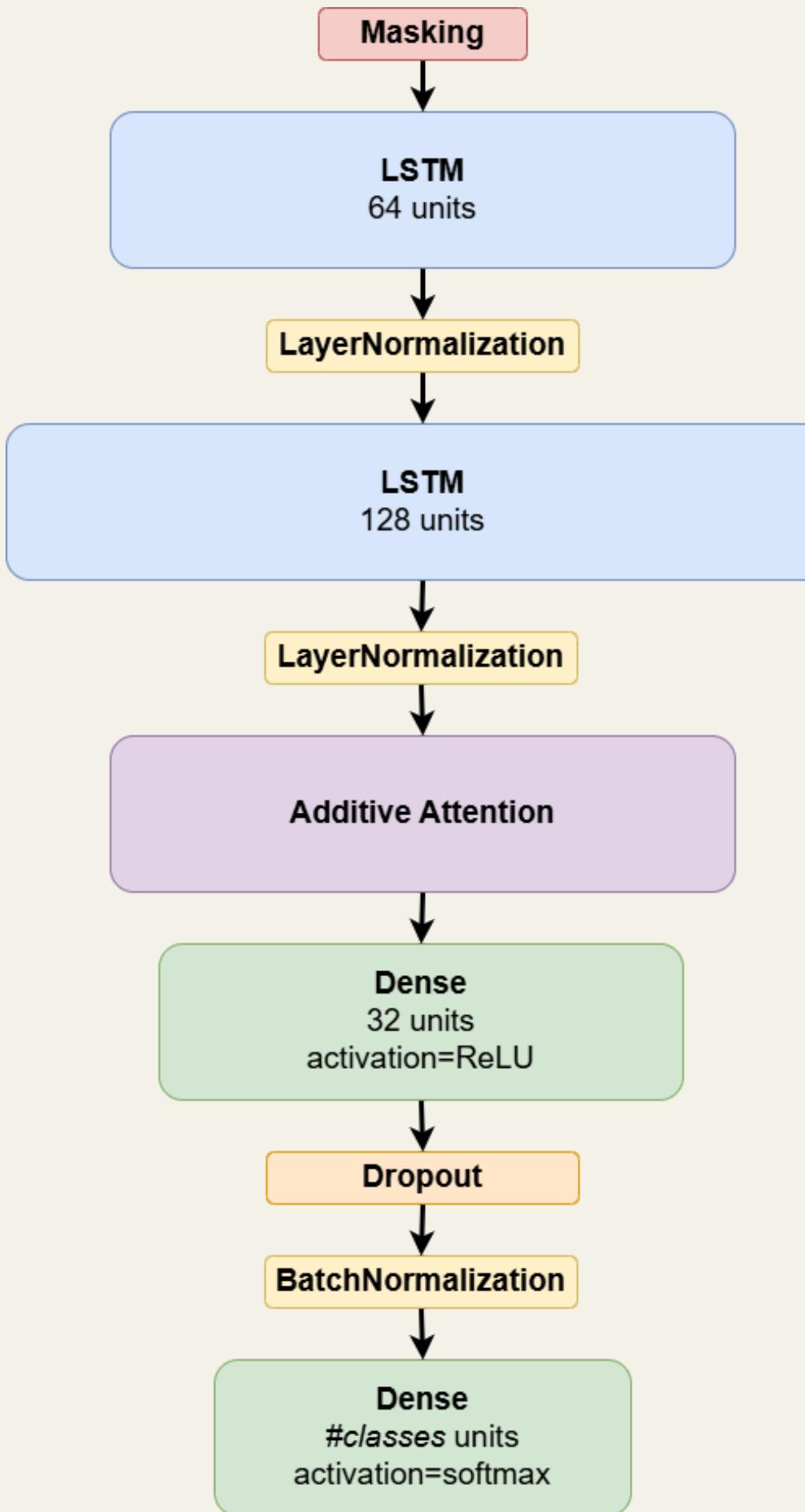
- MASKING LAYER ENSURES PADDED VALUES IN EACH SEQUENCE ARE IGNORED
- DROPOUT LAYER RANDOMLY DROPS UNITS TO PREVENT OVERFITTING
- LAYER NORMALIZATION LAYER NORMALIZES ACROSS FEATURES FOR EACH TIMESTEP
- BATCH NORMALIZATION LAYER NORMALIZES ACROSS THE BATCH FOR EACH FEATURE

Model Design: LSTM



- STACKED LSTM ENABLE THE MODEL TO LEARN COMPLEX, HIERARCHICAL SEQUENTIAL PATTERNS
- IT CAPTURES LOW-LEVEL TEMPORAL DEPENDENCIES IN THE 1ST LAYER AND REFINES THEM INTO HIGHER-LEVEL ABSTRACTIONS IN THE 2ND LAYER.

Model Design: LSTM with Attention

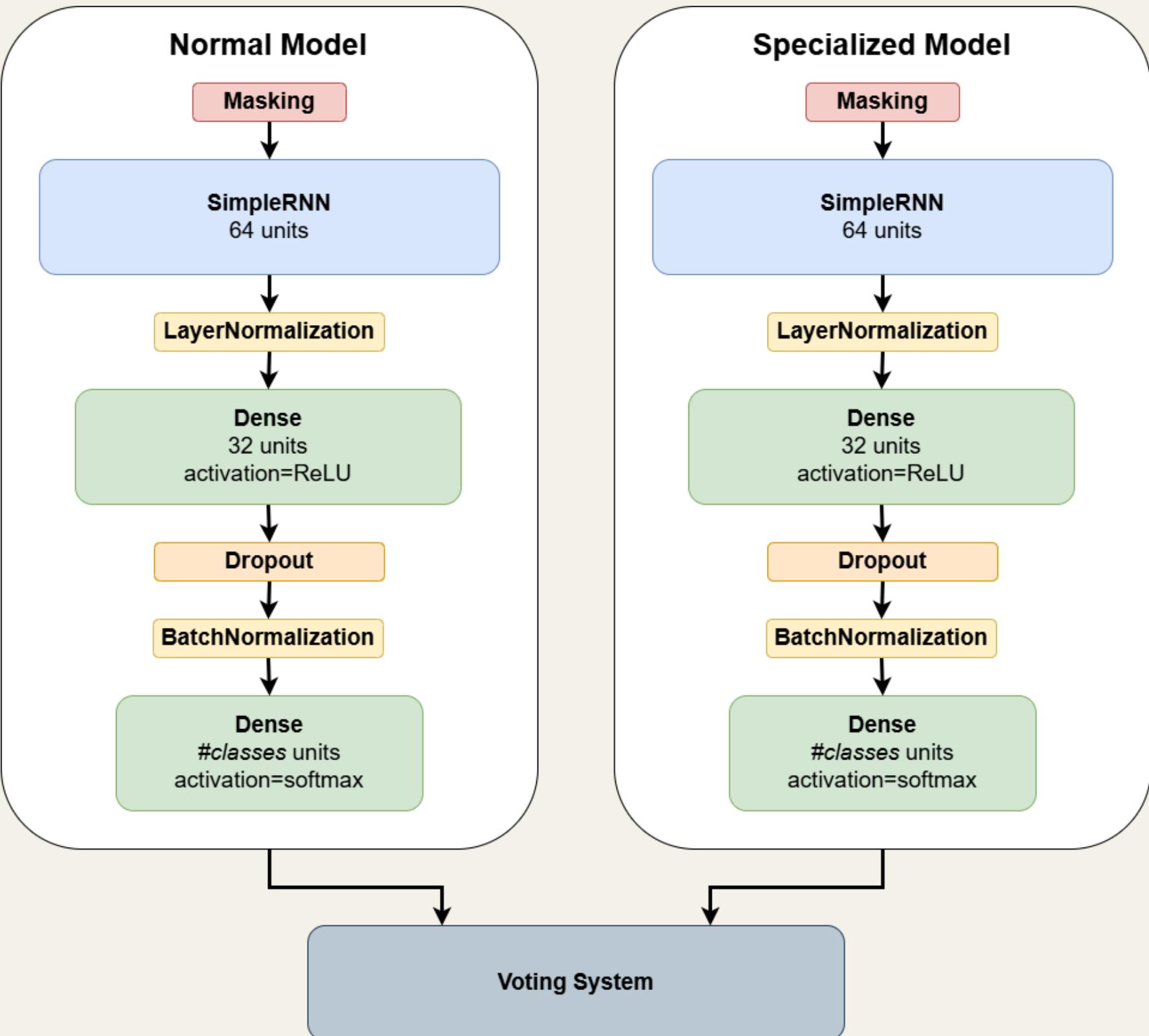


- ATTENTION LAYER HELPS THE MODEL TO FOCUS ON RELEVANT TIMESTEPS IN THE INPUT SEQUENCE

$$C = \sum_{t=1}^T a_t \cdot o_t$$

where a_t is the attention weights,
 o_t is the output of the LSTM layer at timestep t ,
 C is the output of the Attention layer

Model Design: Specialist Model



- **NORMAL MODEL FOCUSES ON PREDICTING FASTBALL PITCH TYPES**
- **SPECIALIZED MODEL FOCUSES ON PREDICTING OFF-SPEED PITCH TYPES**
- **VOTING SYSTEM WEIGHTS THE PREDICTIONS FROM BOTH MODELS AND OUTPUTS THE FINAL PREDICTION**



Experiments

Data Preprocessing

TO PREPARE THE DATA FOR TIME-SERIES MODELING, WE
EXPERIMENTED WITH TWO METHODS FOR GENERATING
SUBSEQUENCES:

- BATTER-BASED
- DATE-BASED

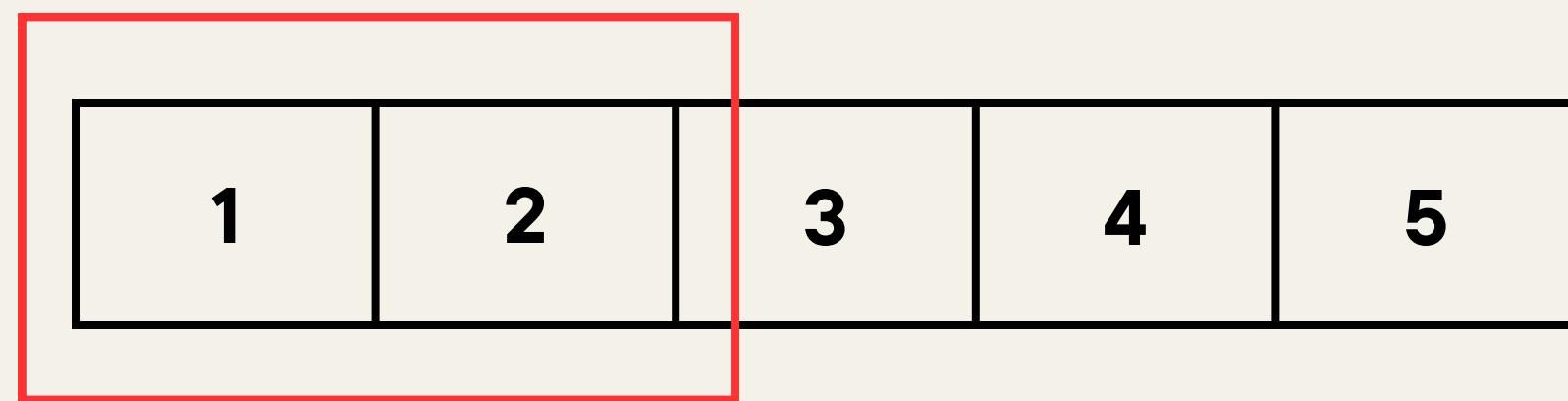
Data Slicing: Take an example

input: min = 2, max = 3

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

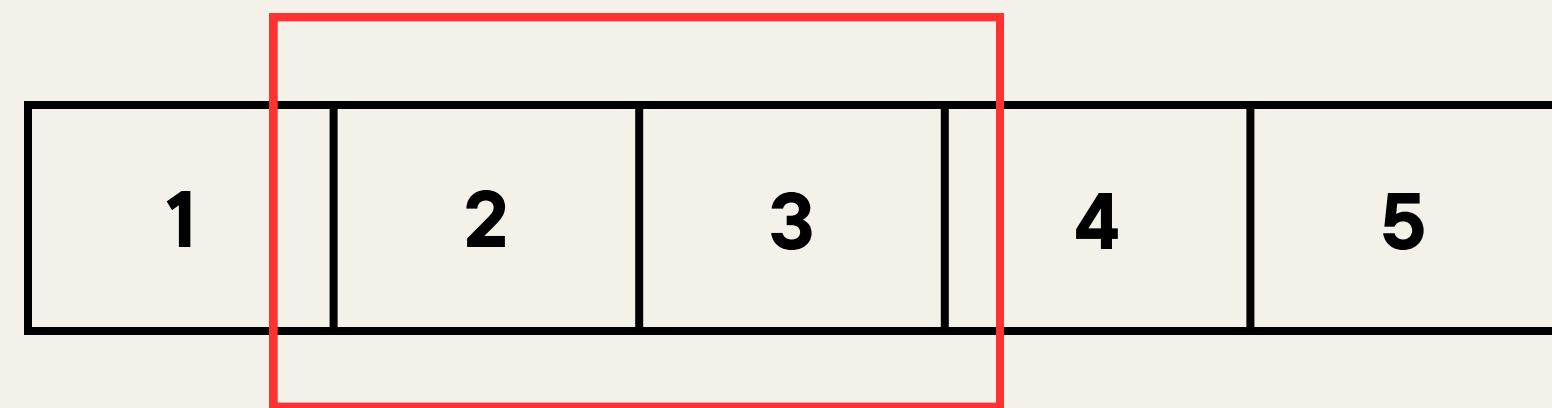
Data Slicing

length = 2



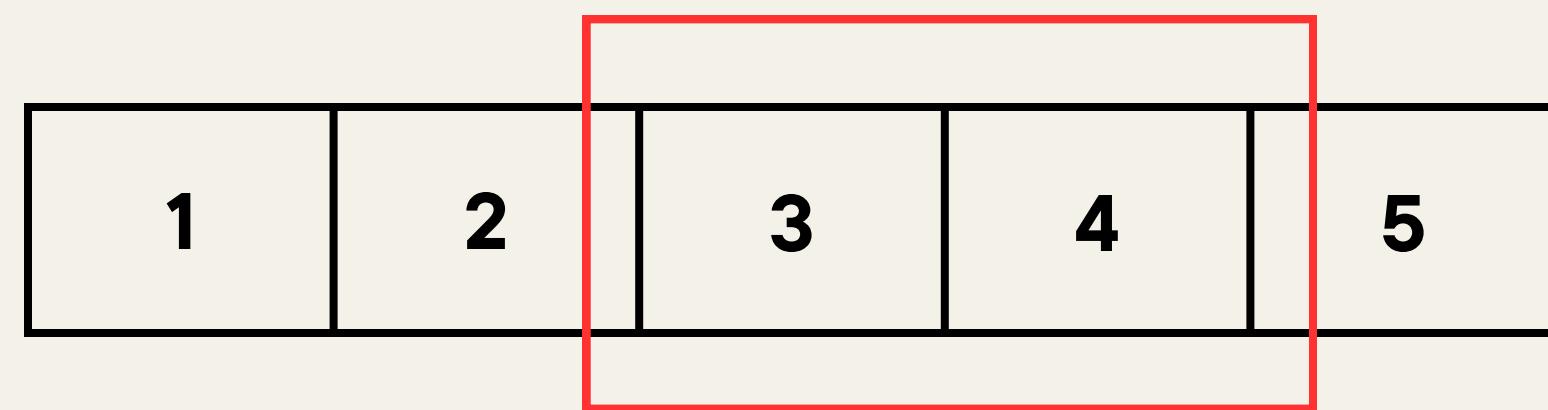
Data Slicing

length = 2



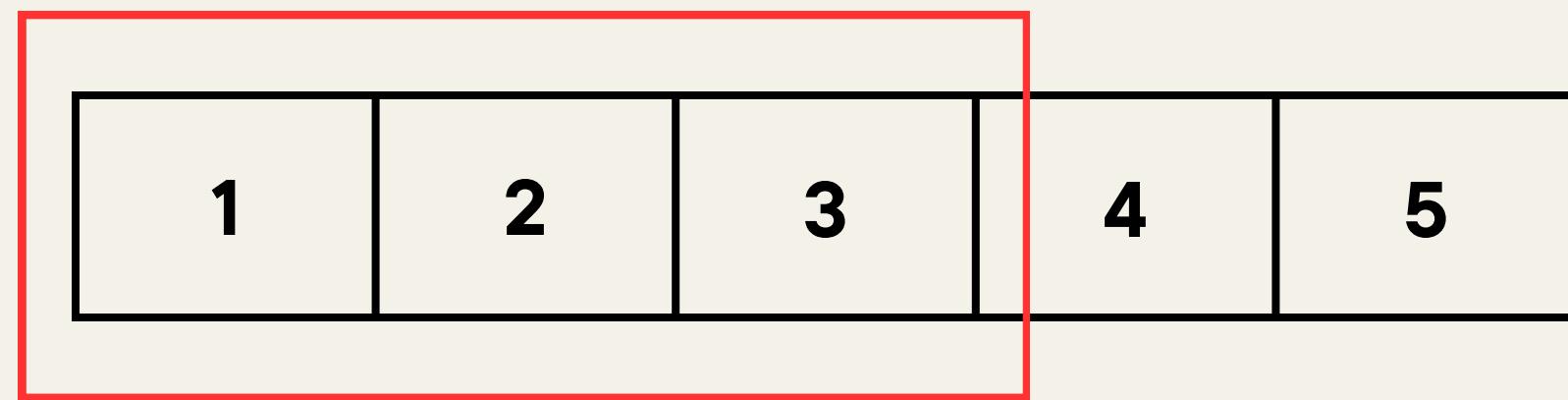
Data Slicing

length = 2



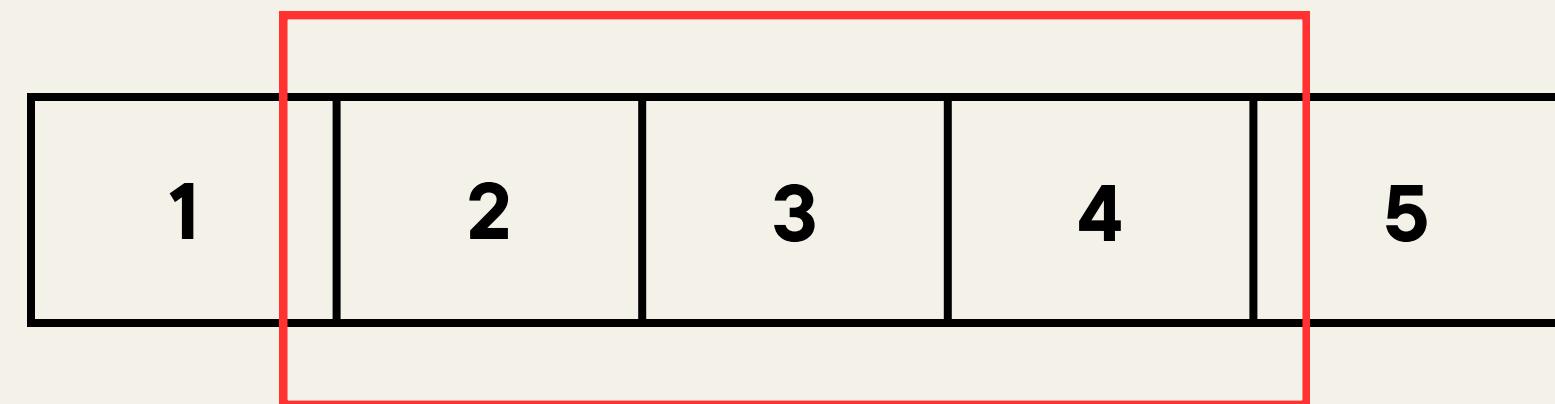
Data Slicing

length = 3



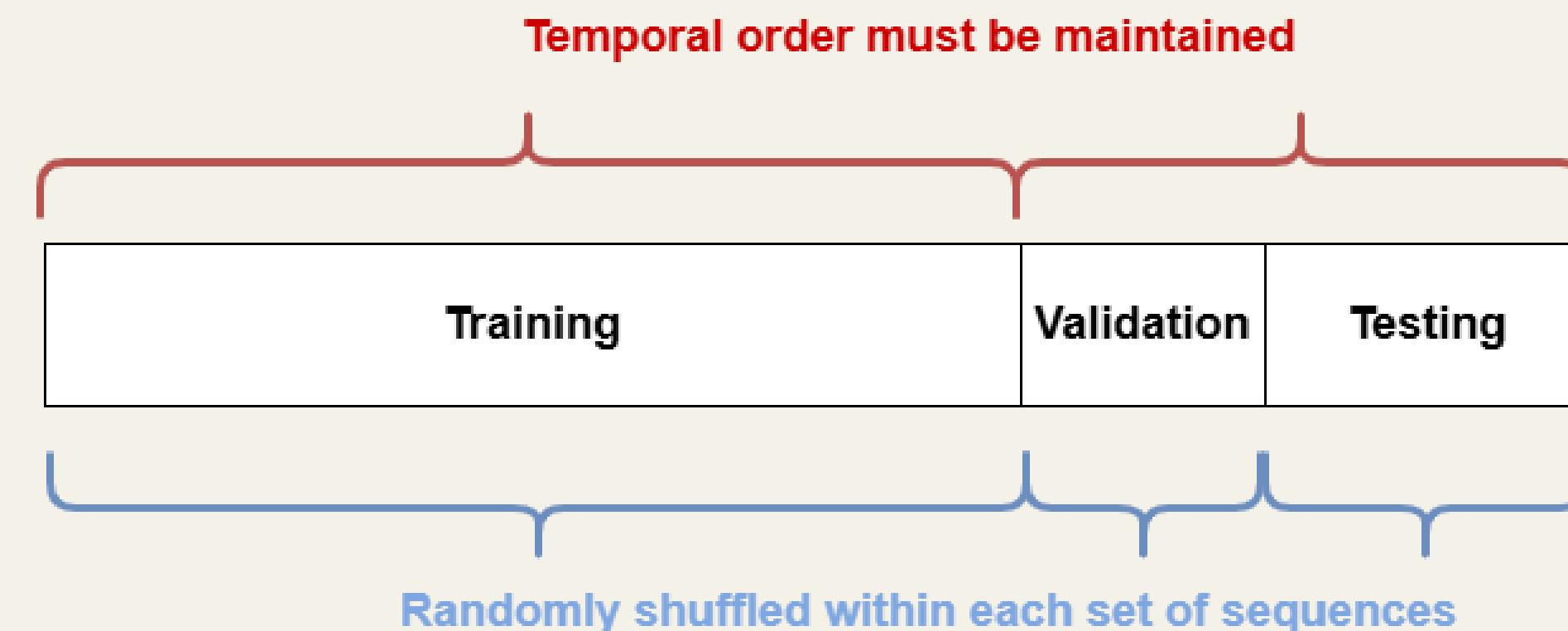
Data Slicing

length = 3



Data Splitting

- DATA IS SPLIT WITH THE RATIO, TRAINING:VALIDATION:TESTING = 0.64:0.16:0.2
- THE DATASETS {TRAINING, VALIDATION} AND {TESTING} ARE SPLIT WITHOUT BREAKING THEIR TEMPORAL ORDER
- SUBSEQUENCES IN EACH OF THE ABOVE DATASET ARE RANDOMLY SHUFFLED TO PREVENT THE MODEL FROM MISTAKENLY LEARNING TEMPORAL RELATIONS BETWEEN THESE SUBSEQUENCES



yes/no

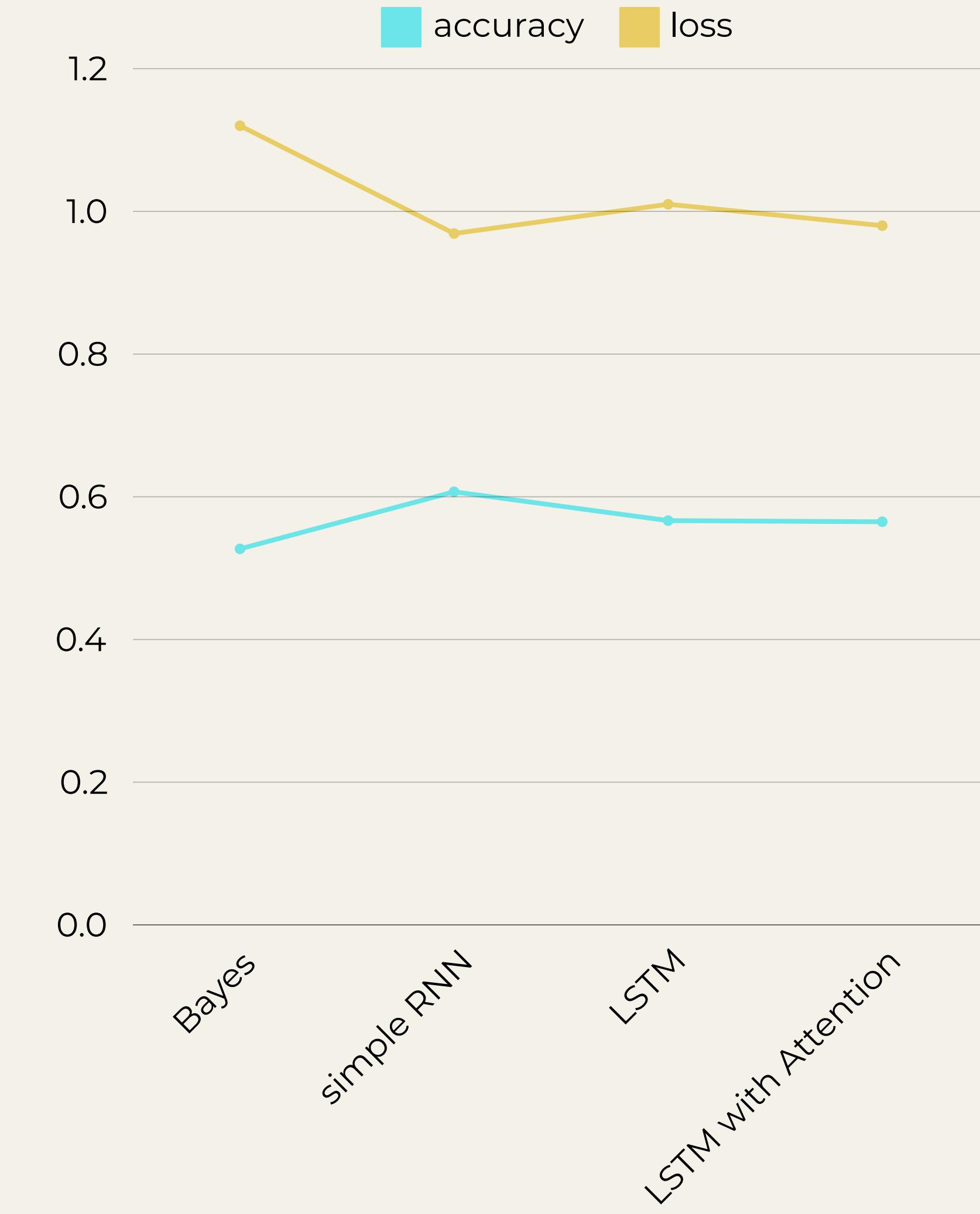
STARTING PITCHER

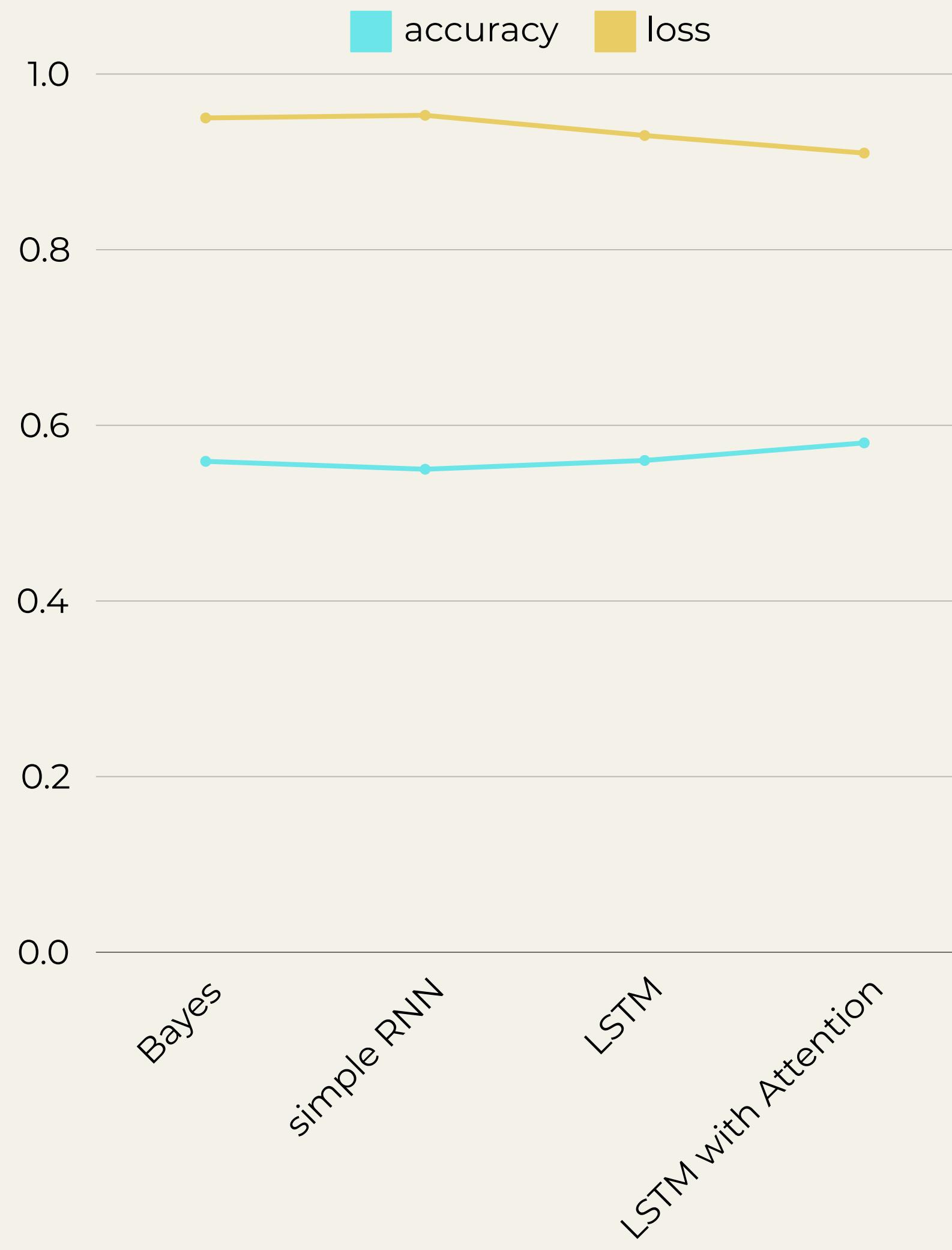
STARTING PITCHER

STARTING PITCHER



SHOTA IMANAGA





STARTING PITCHER

STARTING PITCHER

STARTING PITCHER



HUNTER GREENE

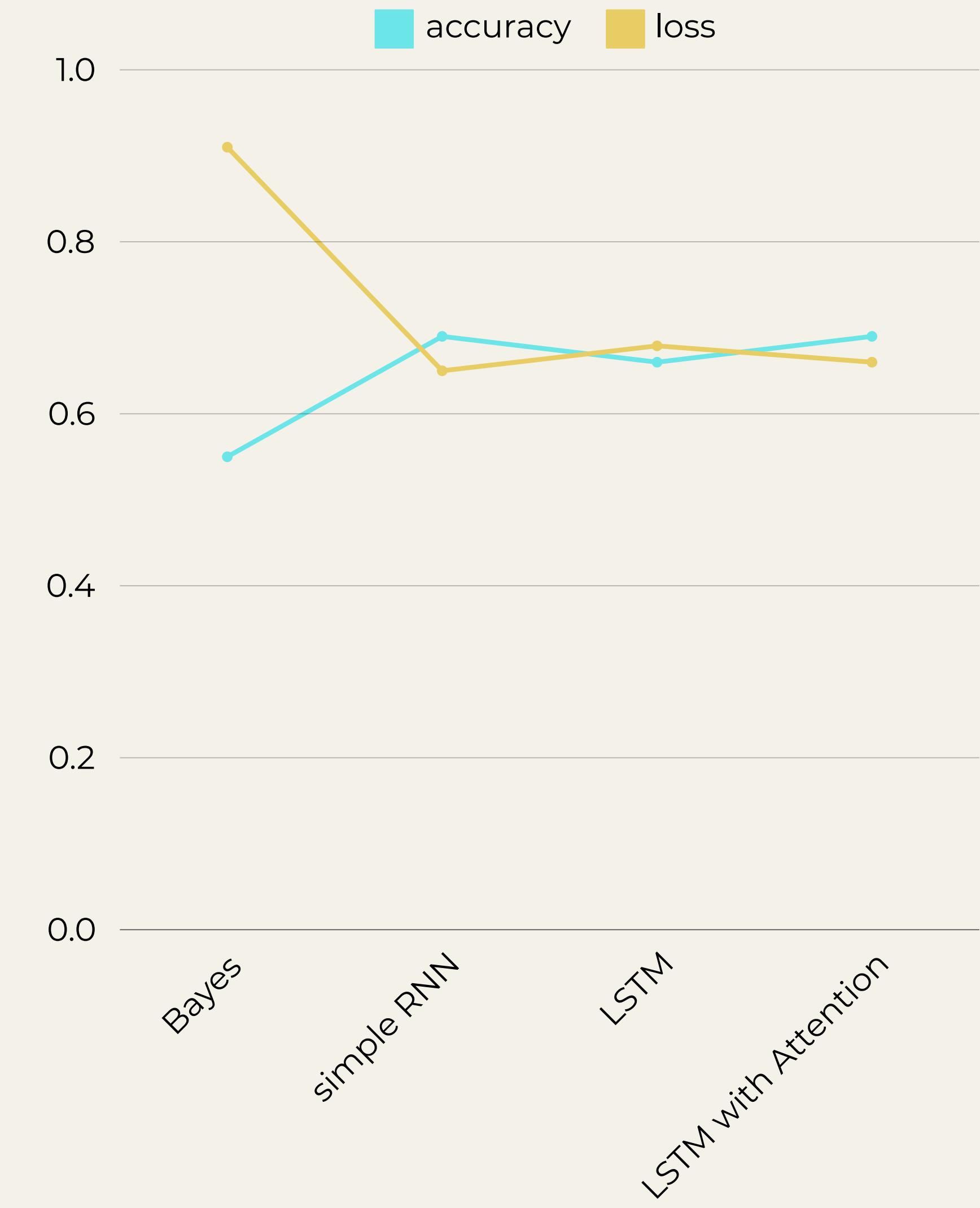
RELIEF PITCHER

RELIEF PITCHER

RELIEF PITCHER



MASON MILLER



Our accuracy of predicting pitch type in relief pitcher



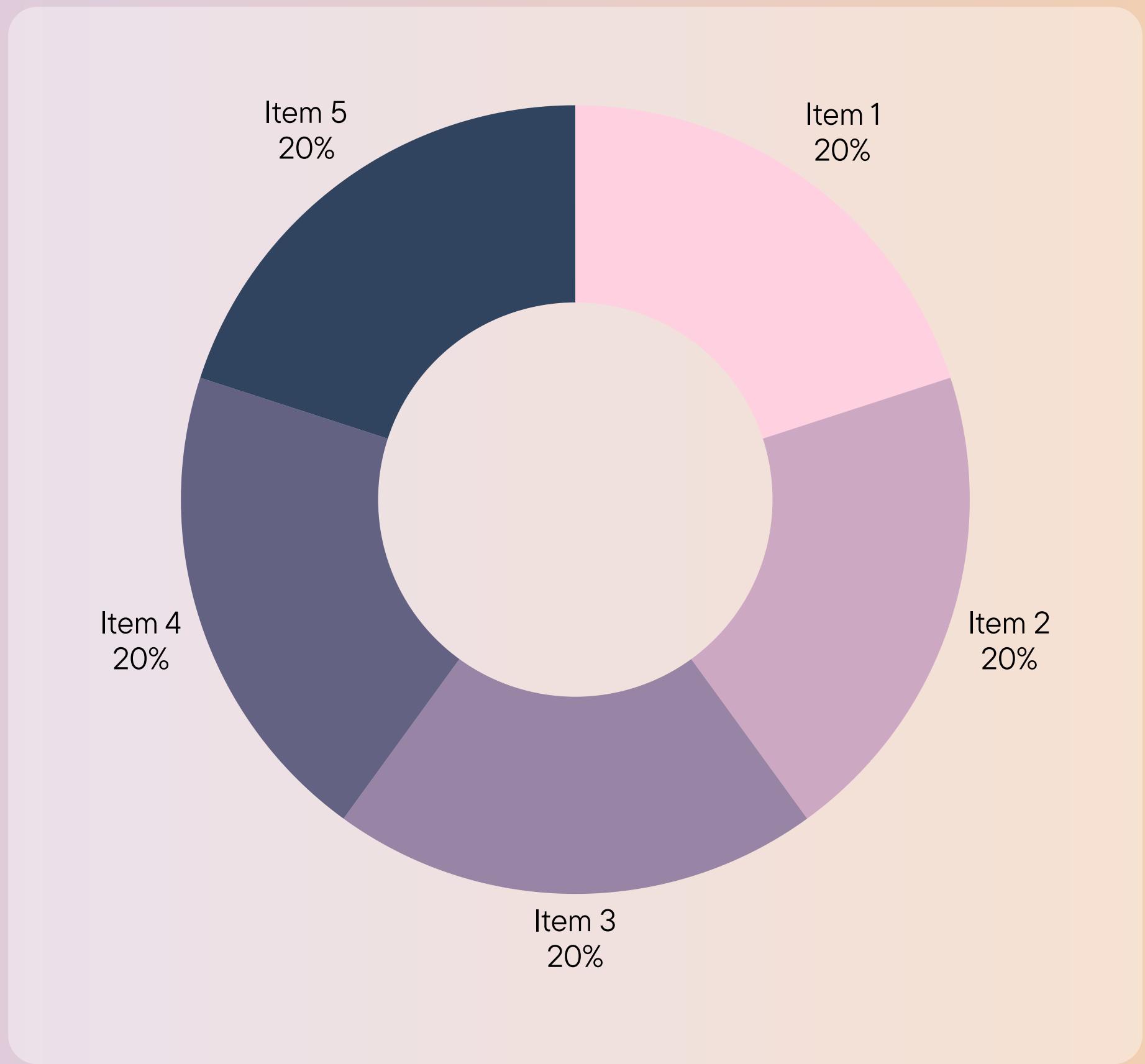
**FAIL TO SEE THE BENEFITS OF DIVERSITY BECAUSE OF A LACK OF
INCLUSIVE PRACTICES**

CONFRONTING UNCONSCIOUS BIAS

Awareness is the first step to addressing unconscious bias.

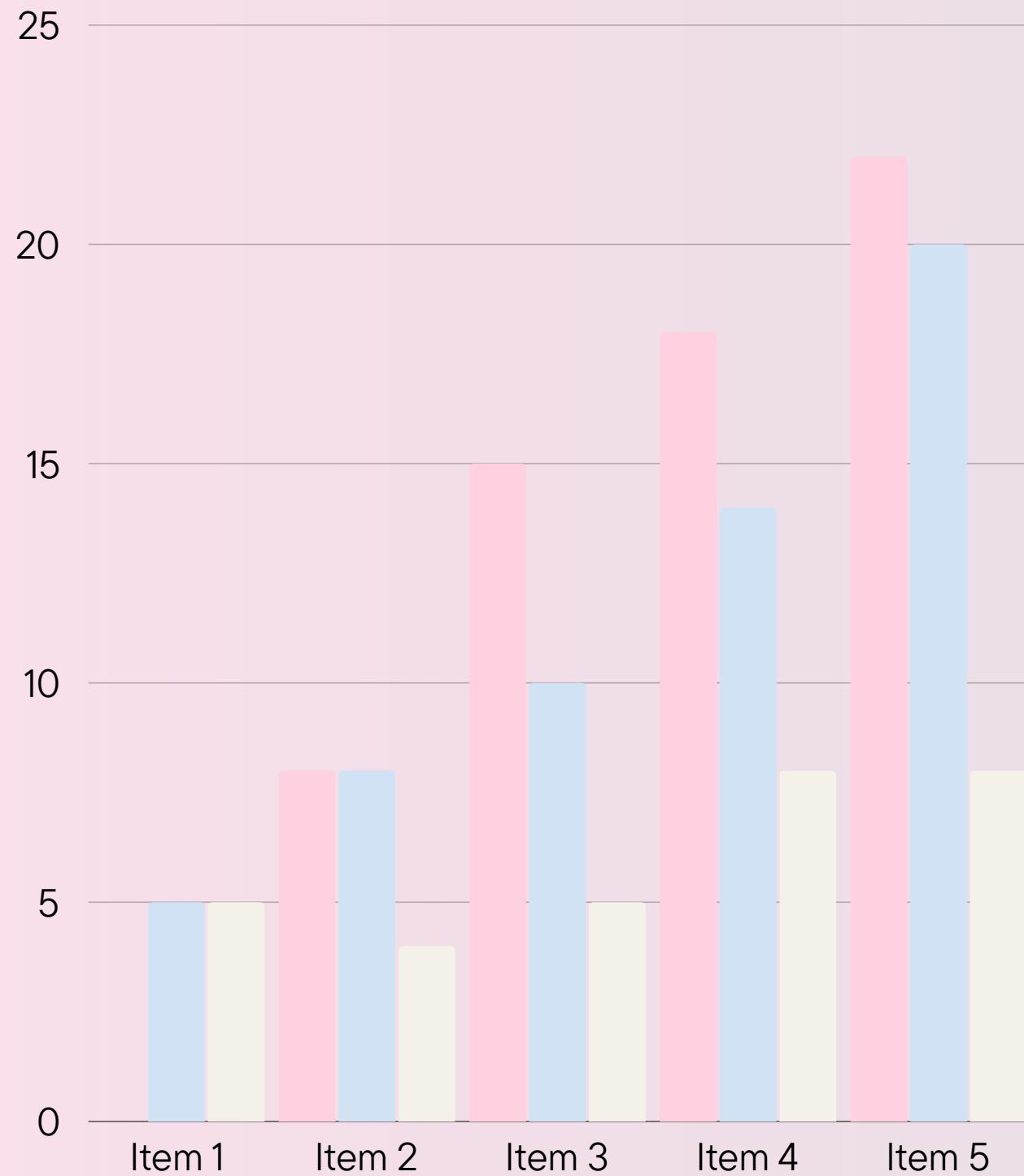
Here are some behaviors employees identified as showing biases and microaggressions.

Visualize complicated and dense information with graphs and charts. These are visual aids that help add more context to the topic you are discussing.



DRIVING GROWTH THROUGH D&I

Practicing diversity and inclusion increases talent retention, employee satisfaction, and productivity.



See how businesses have grown based on their commitment to diversity.

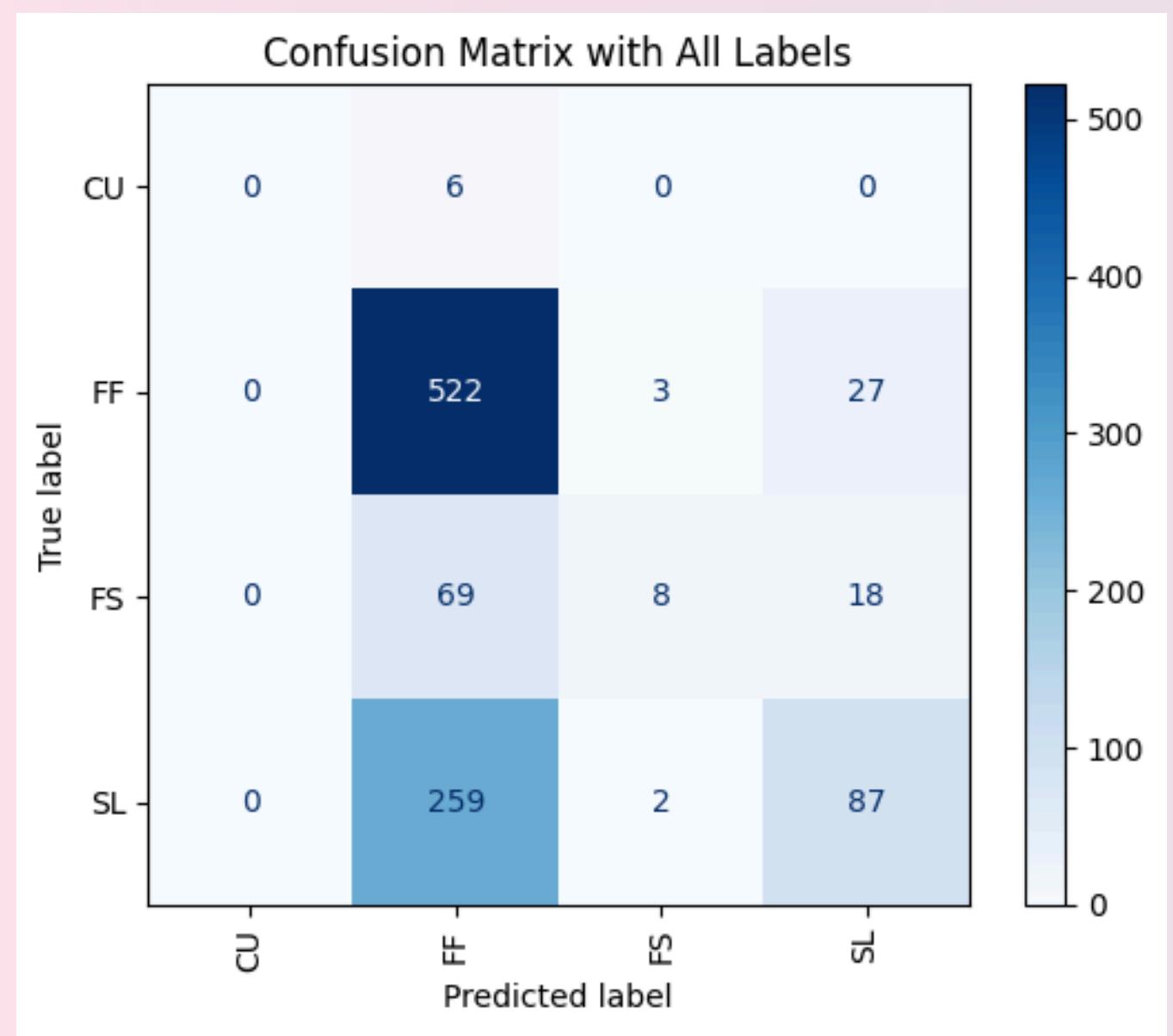
Visualize complicated and dense information with graphs and charts. These are visual aids that help add more context to the topic you are discussing.

Discussion & Conclusion

1. Accuracy and Loss Are Not Enough to Evaluate Models.

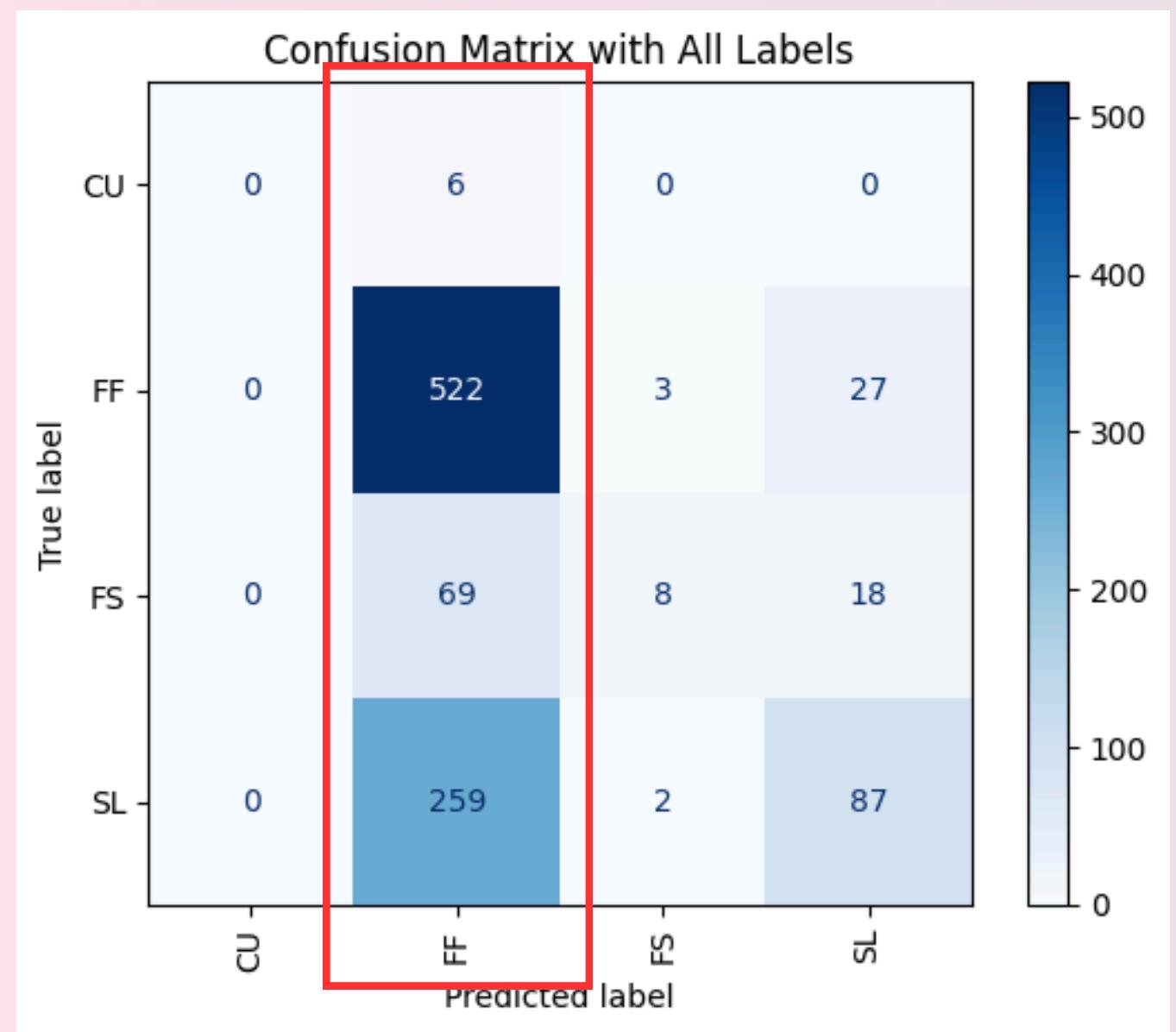
1. Accuracy and Loss Are Not Enough: Confusion Matrix

HUNTER GREENE: 61.68%



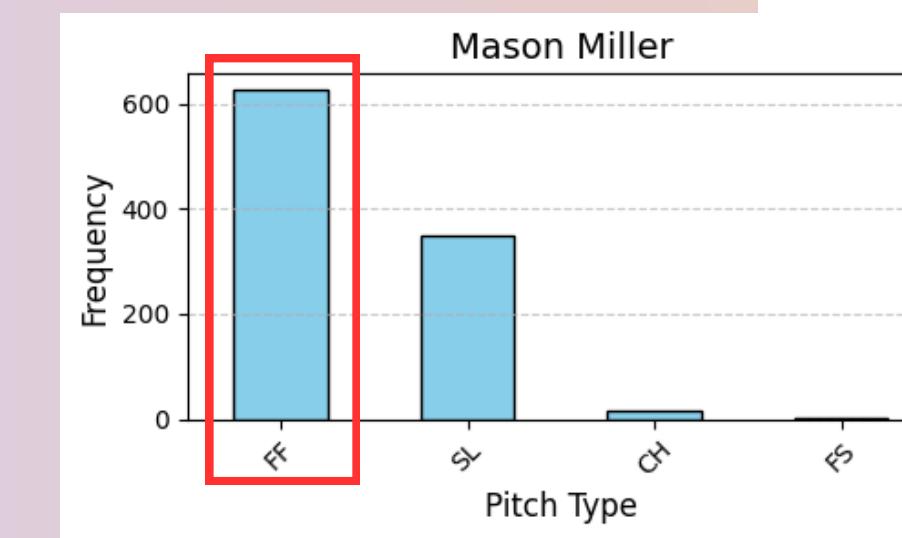
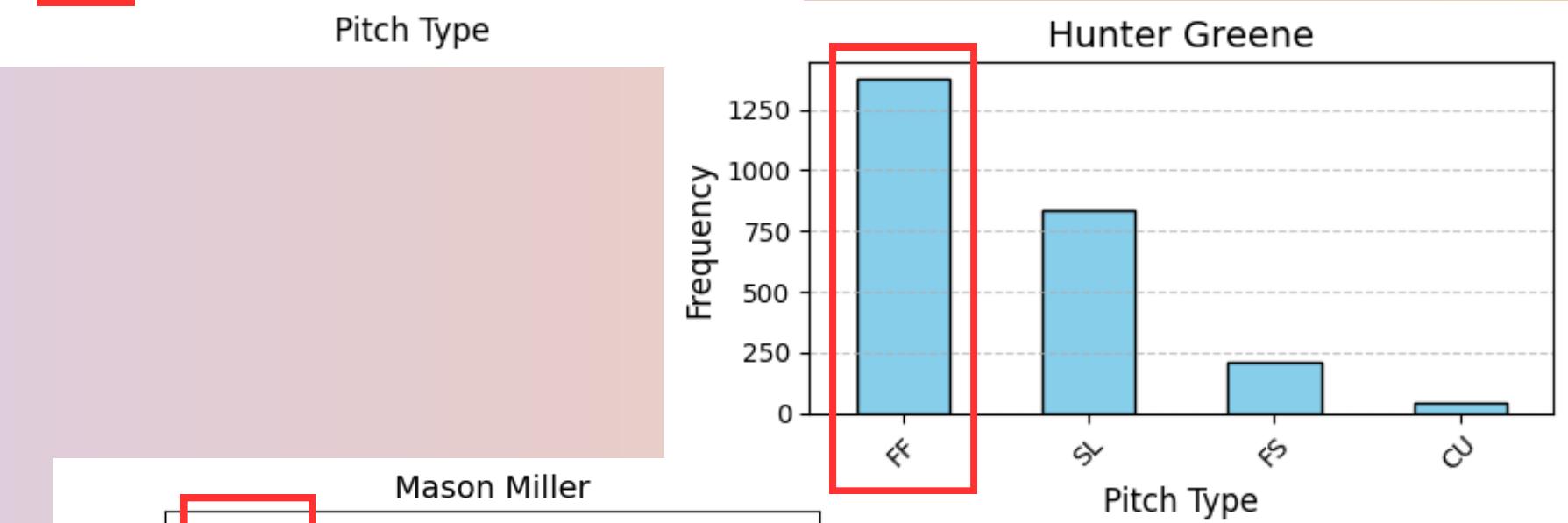
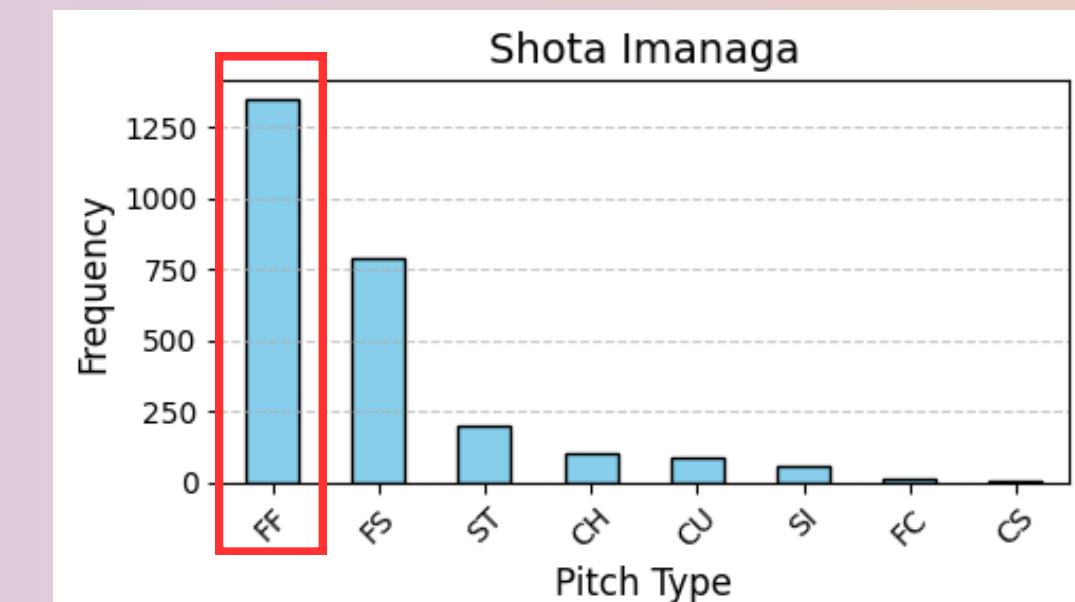
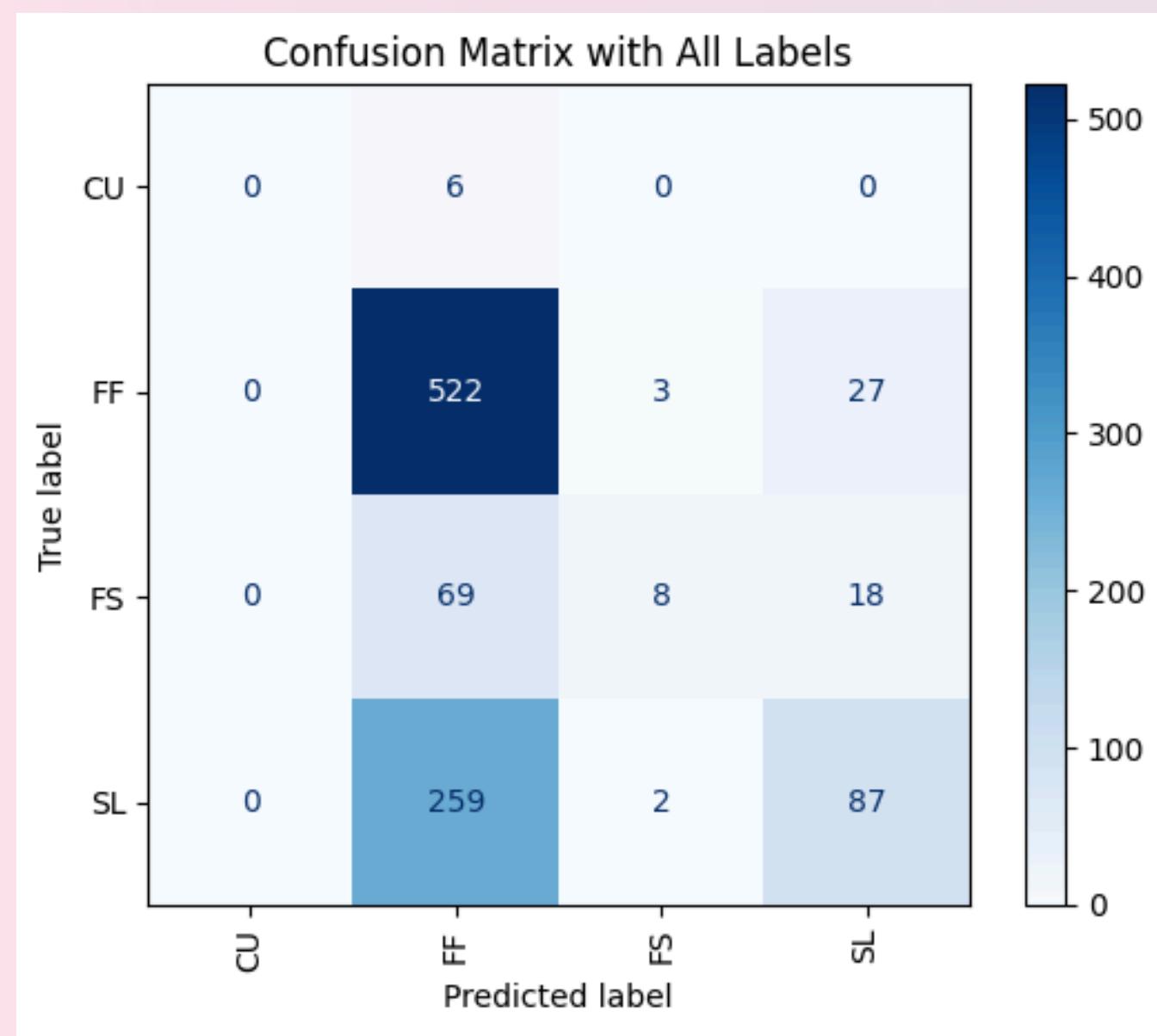
1. Accuracy and Loss Are Not Enough: Confusion Matrix

HUNTER GREENE: 61.68%



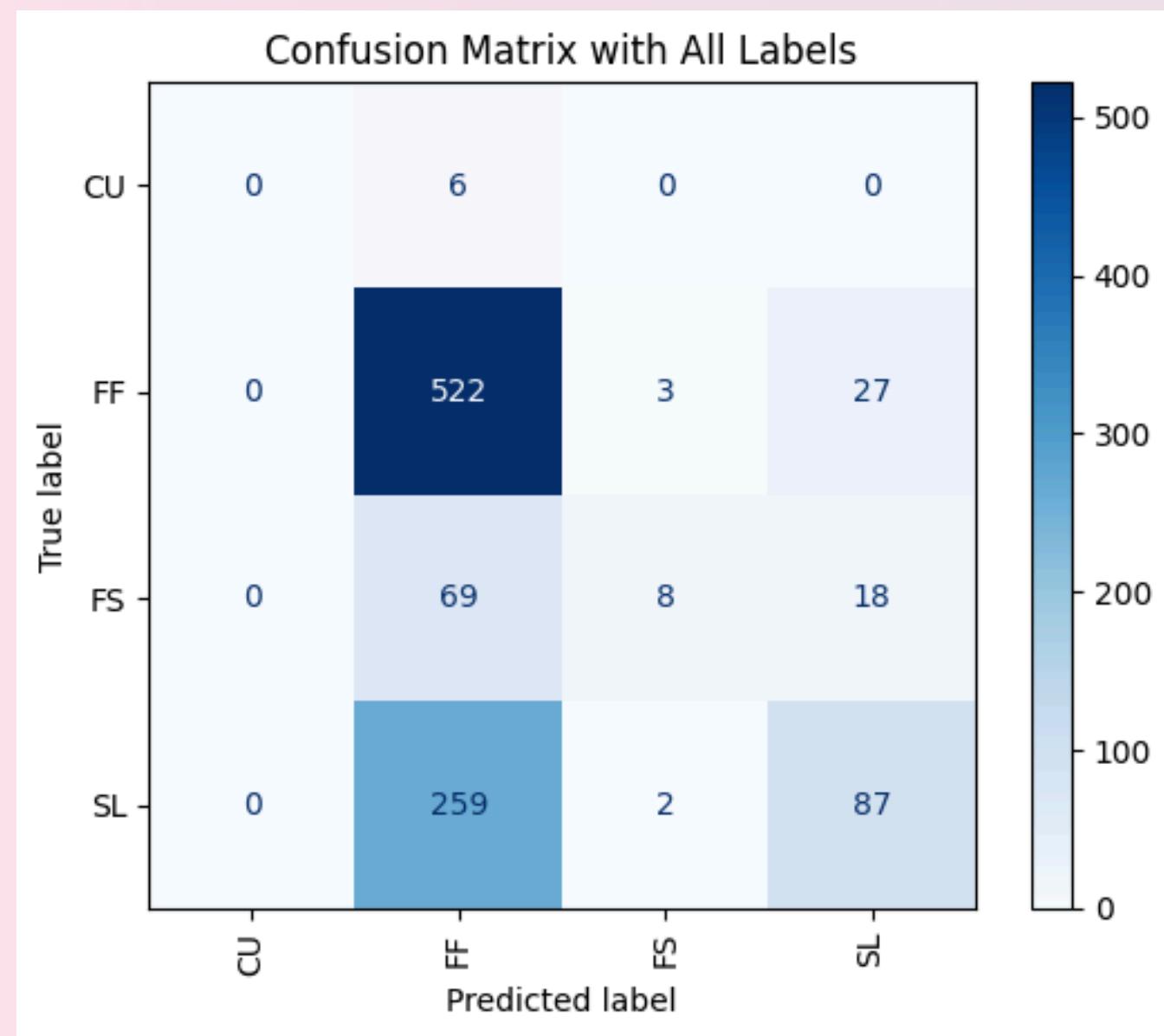
1. Accuracy and Loss Are Not Enough: Confusion Matrix

HUNTER GREENE: 61.68%



1. Accuracy and Loss Are Not Enough: Confusion Matrix

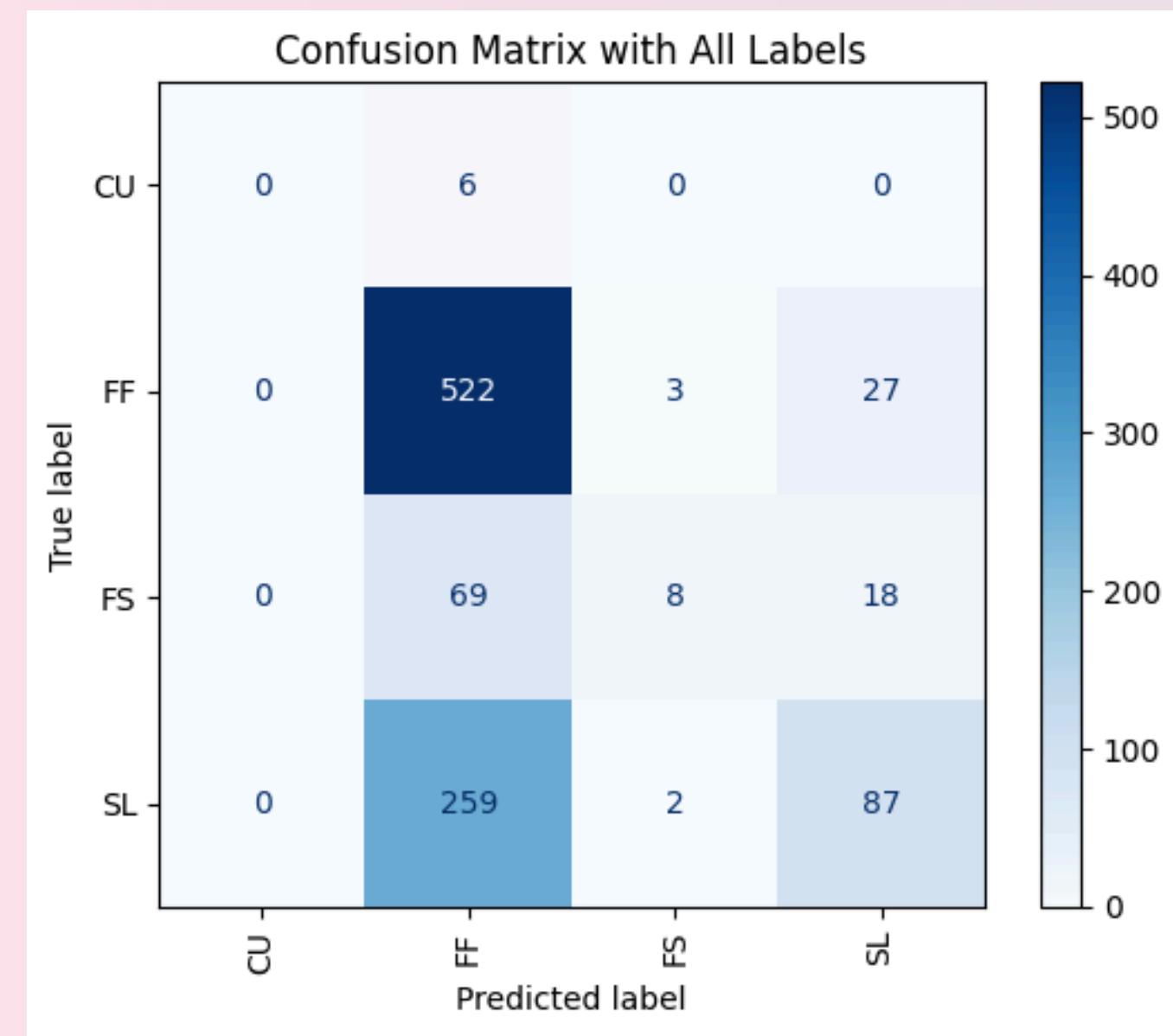
HUNTER GREENE: 61.68%



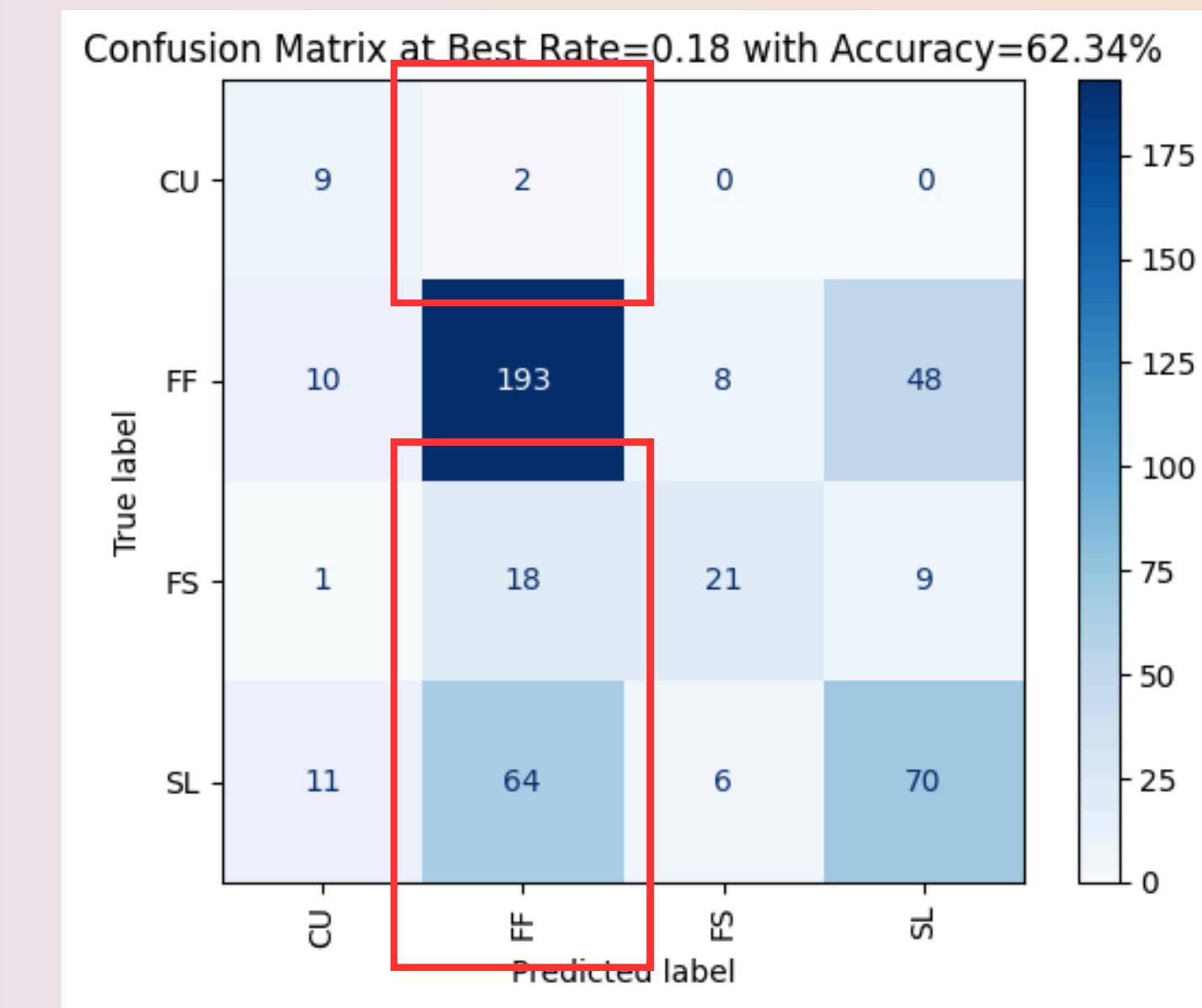
Adding Specialist Model?

1. Accuracy and Loss Are Not Enough: Confusion Matrix

SIMPLE RNN: 61.68%

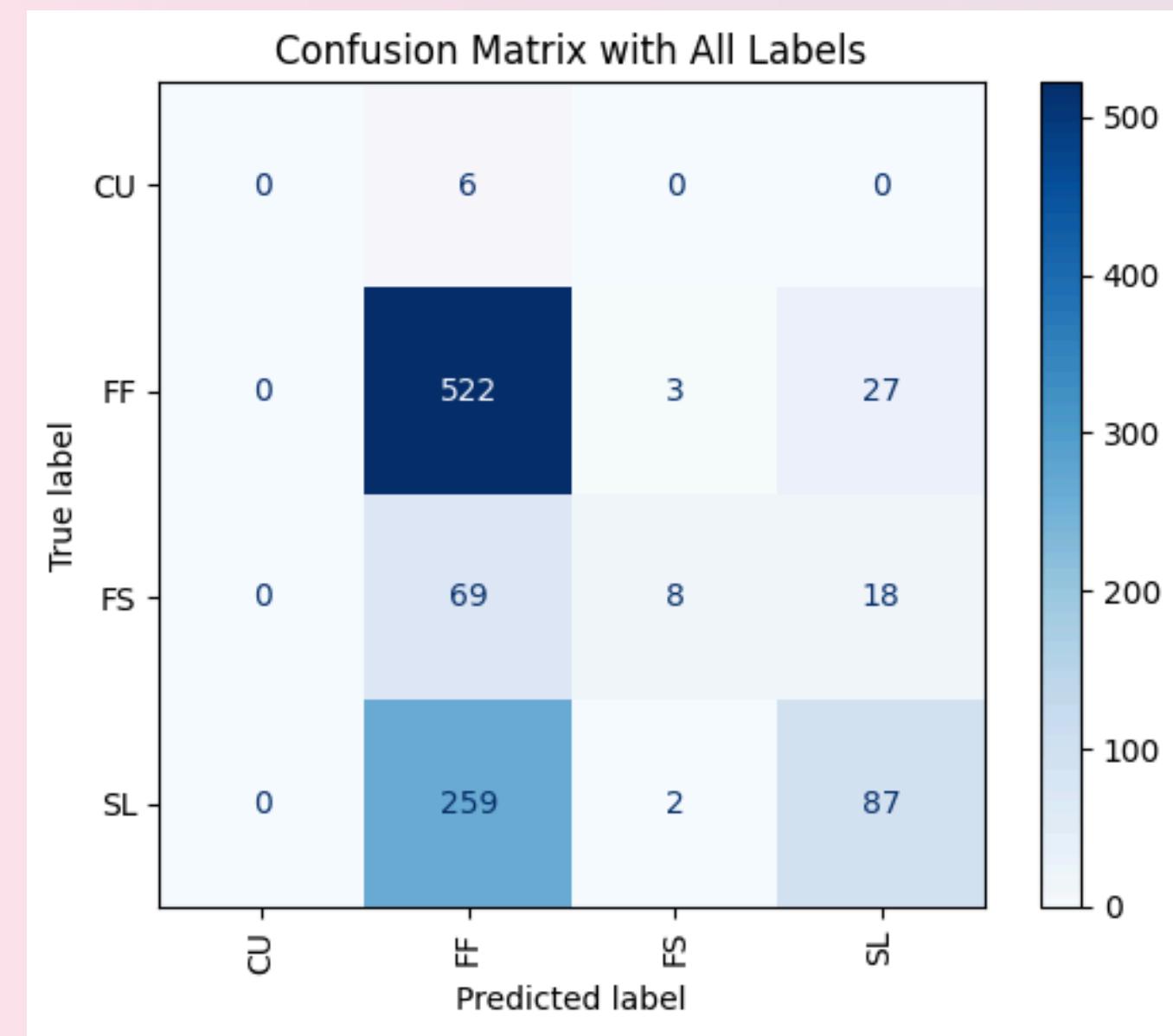


ADDING SPECIALIST MODEL: 62.34%

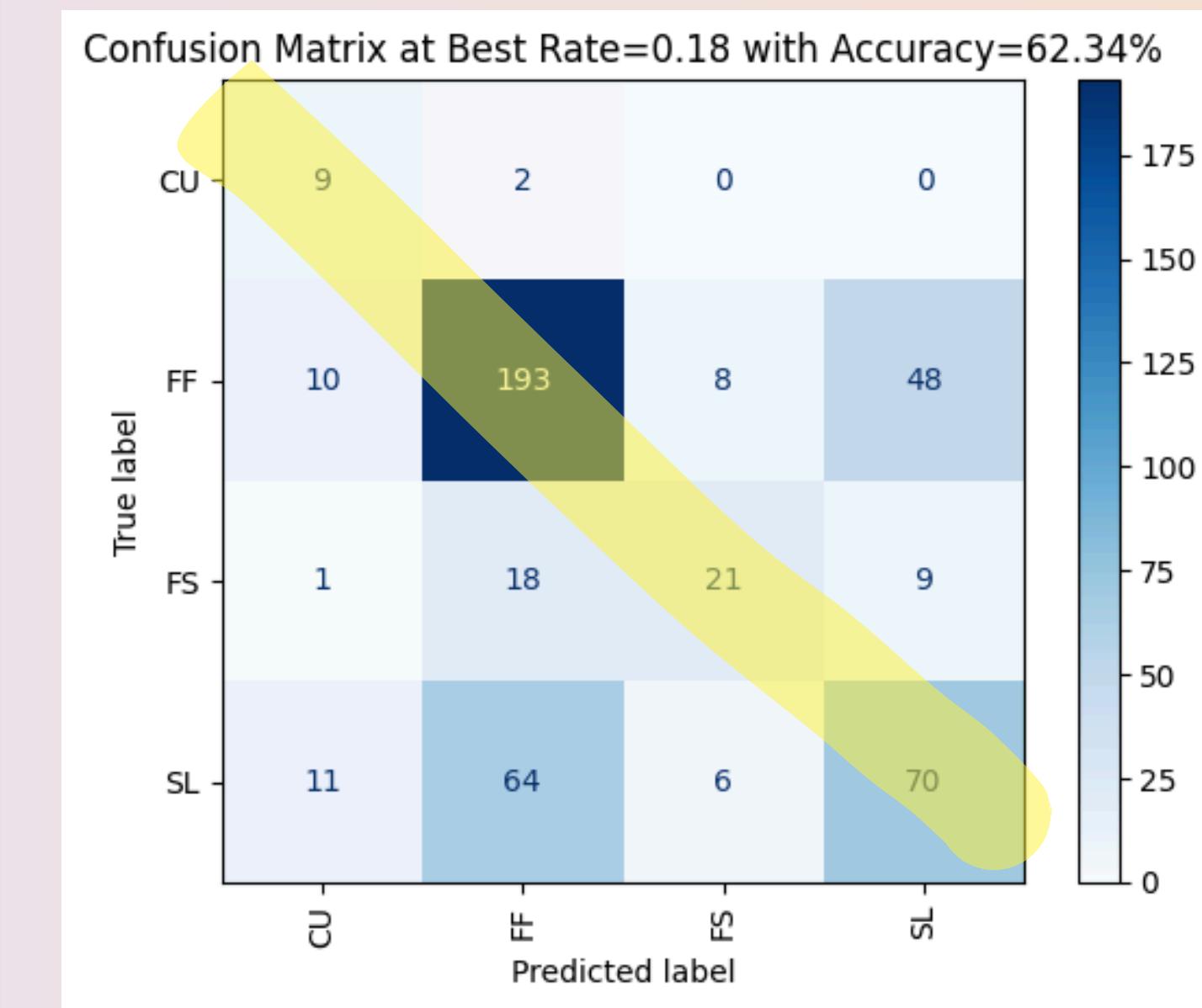


1. Accuracy and Loss Are Not Enough: Confusion Matrix

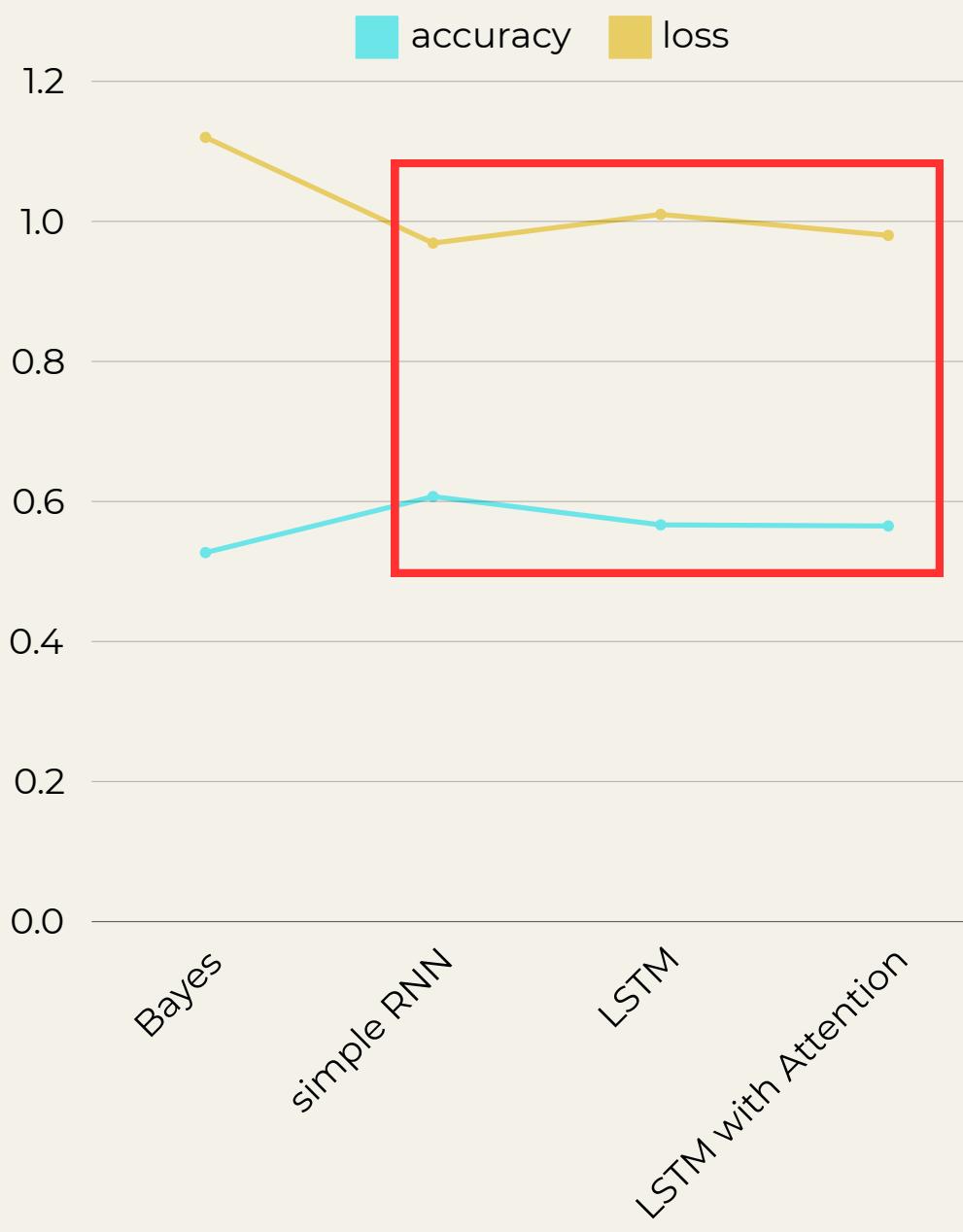
SIMPLE RNN: 61.68%



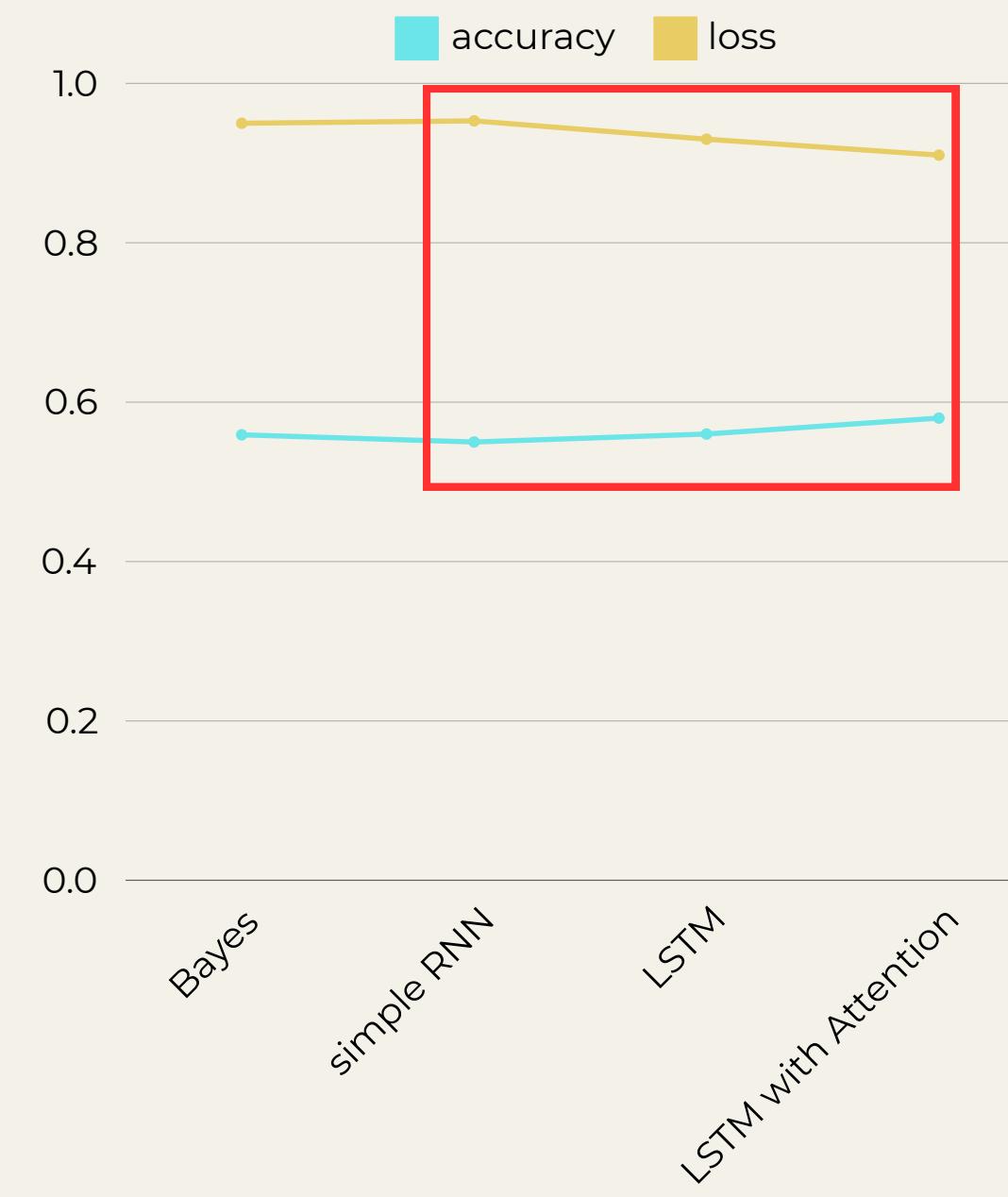
ADDING SPECIALIST MODEL: 62.34%



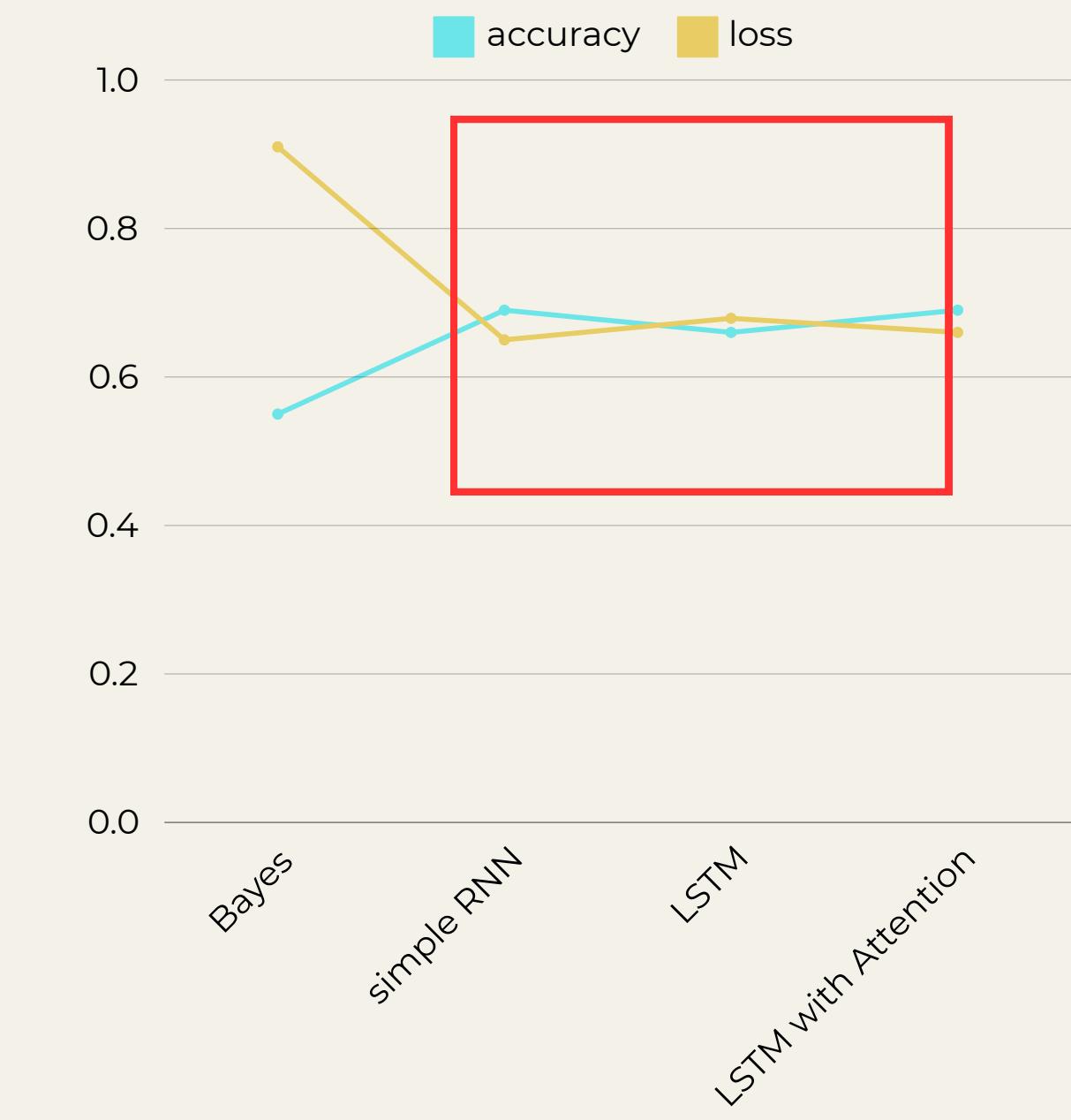
2. LSTM and LSTM+Attention Do Not Outperform on Short Sequences



SHOTA IMANAGA



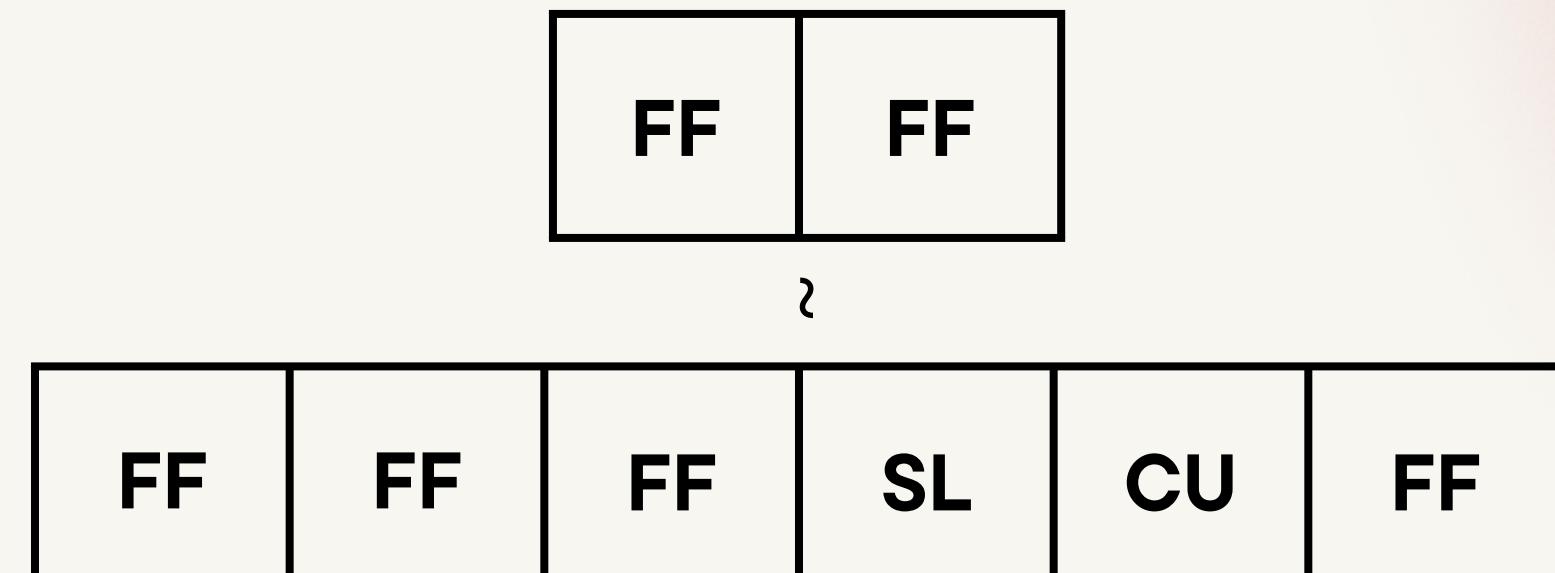
HUNTER GREENE



MASON MILLER

2. LSTM and LSTM+Attention Do Not Outperform on *Short* Sequences

(1) LONG-TERM MEMORY ADVANTAGE A NOT FULLY UTILIZED



2. LSTM and LSTM+Attention Do Not Outperform on *Short Sequences*

(1) MEMORY ADVANTAGE NOT FULLY UTILIZED

FF	FF
----	----

?

FF	FF	FF	SL	CU	FF
----	----	----	----	----	----

FF	FF	FF	SL	CU
----	----	----	----	----

(2) BIAS ACCUMULATION OF ATTENTION LAYER

OUR MODEL OUTPERFORMS BAYESIAN
BASELINES

Conclusion:

Conclusion:

OUR MODEL OUTPERFORMS BAYESIAN
BASELINES

OUR MODEL PERFORMS BETTER FOR
PITCHERS WITH FEWER PITCH TYPES

Conclusion:

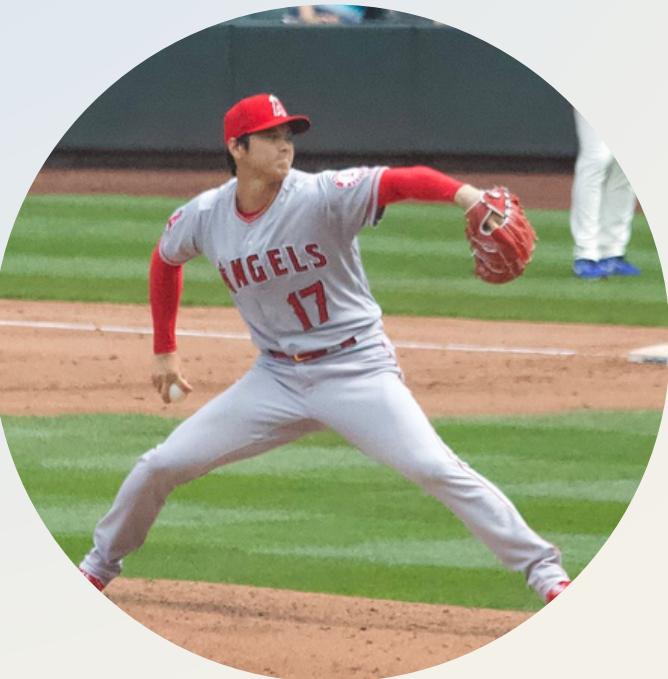
OUR MODEL OUTPERFORMS BAYESIAN
BASELINES

OUR MODEL PERFORMS BETTER FOR
PITCHERS WITH FEWER PITCH TYPES

BUT FOR PITCHERS WITH MORE PITCH
TYPES, OUR MODEL STILL REMAINS SOLID

Practical Value of our Model:

PROVIDES A POWERFUL TOOL TO ANALYZE THE PITCHING TENDENCIES



PITCHERS



COACHES



CATCHERS



THANKS FOR LISTENING

TEAM23