



24/02/2021

claire SUN

Sentiment Analysis for Stock Price Prediction

Claire Sun

1

Introduction

2

Design

3

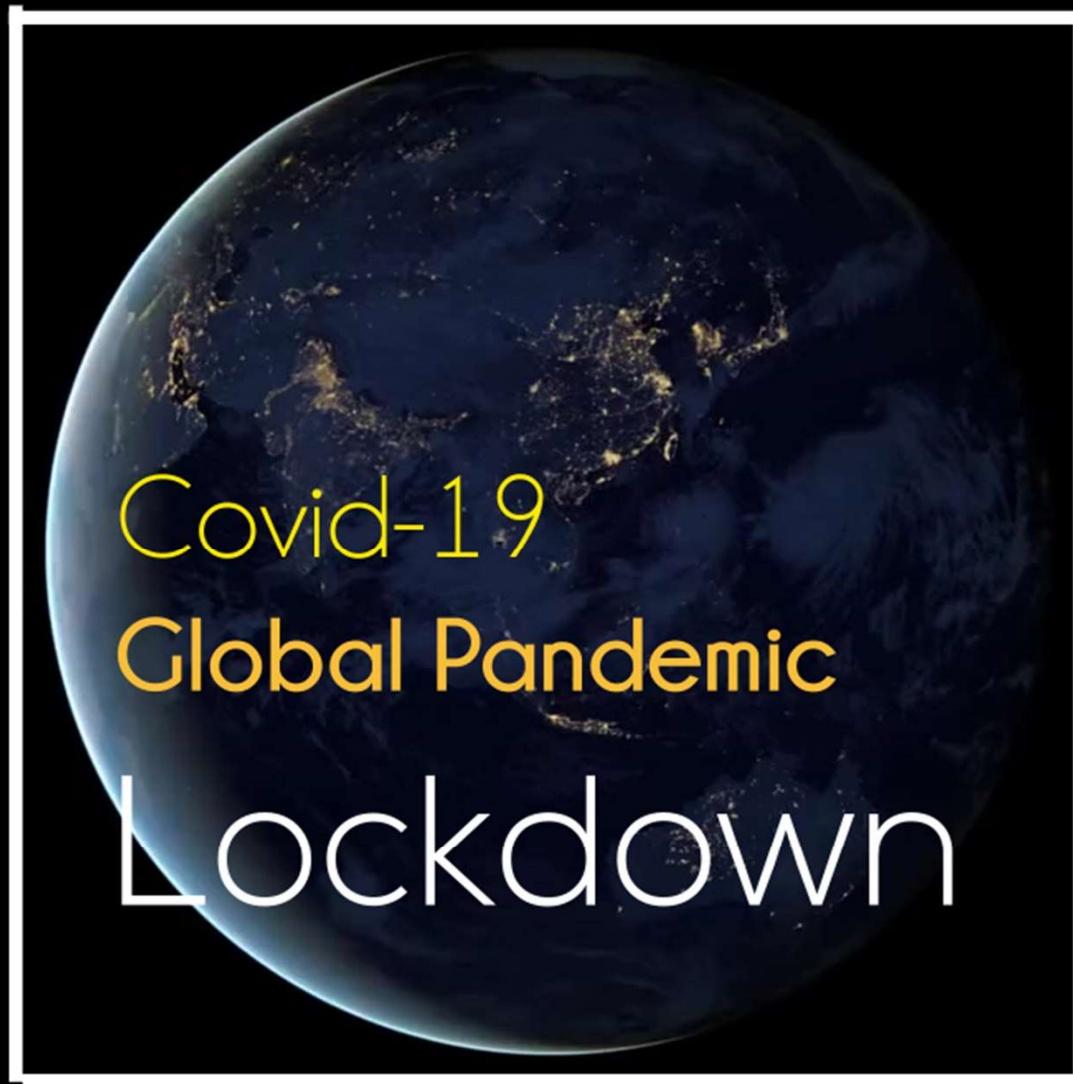
Demo

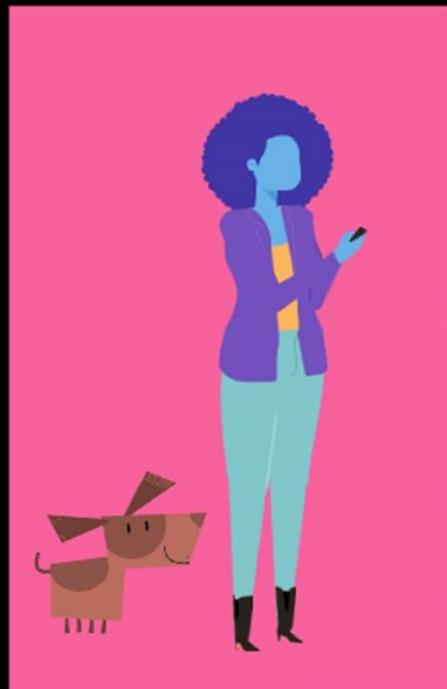
4

Conclusion

**The main topics we are going to
cover today**

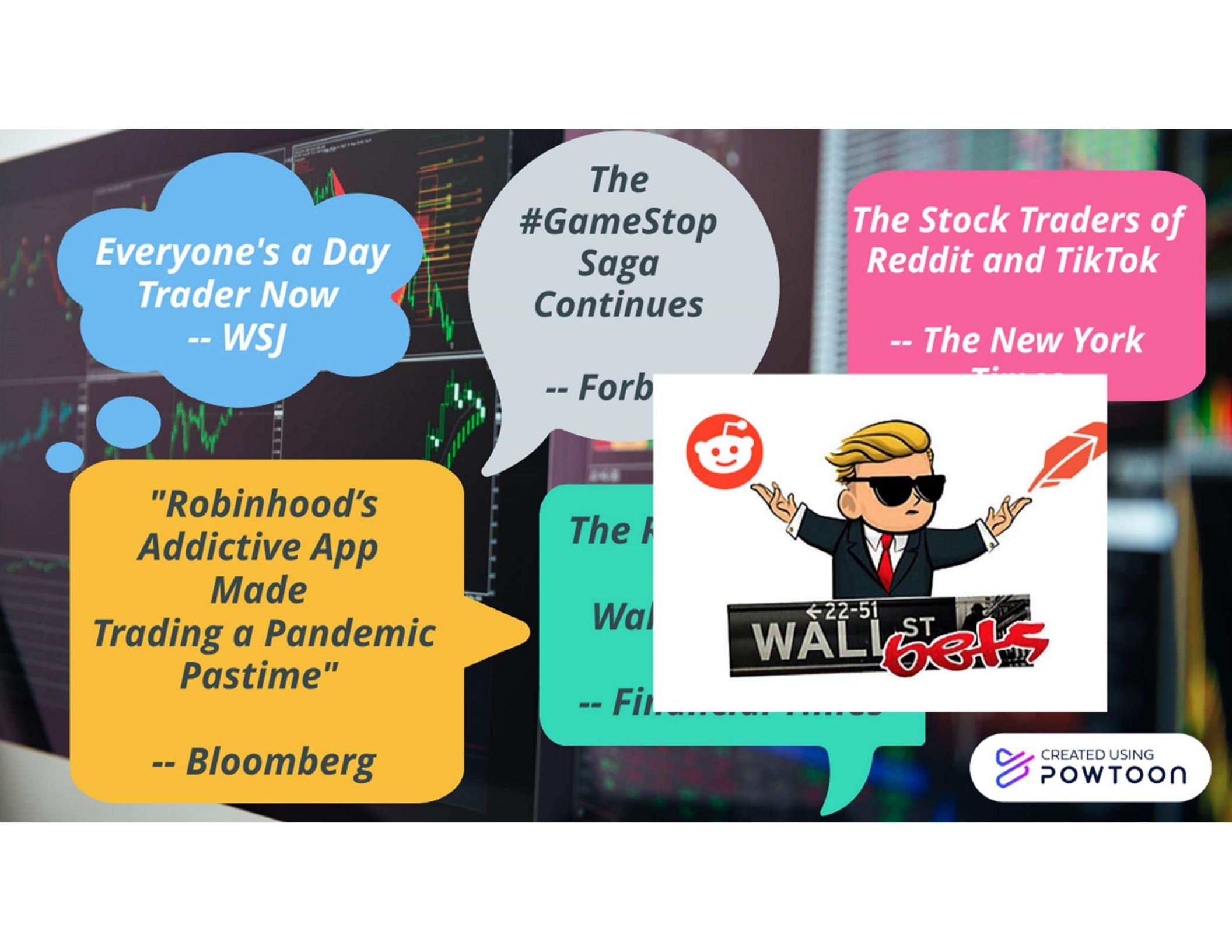
2020





U5ErkJgg=="})base64,iVBORw0KG
0C1V9QD4NEOxs9xBQHQVCwSJF
hwVYBIdLn9vkLp79QcBCTDMiy3w2
QtfMBSO **TRADING** Uz4BDMle
MledhqOu/AzVSmzIJKUz4BDMled
/AzVSmzZ49CUjCC0yvim98iqtJT2L2
IQx8Q7hQYFek4AkixXFe1rsFR4I/RT





*Everyone's a Day
Trader Now*
-- WSJ

*The
#GameStop
Saga
Continues*

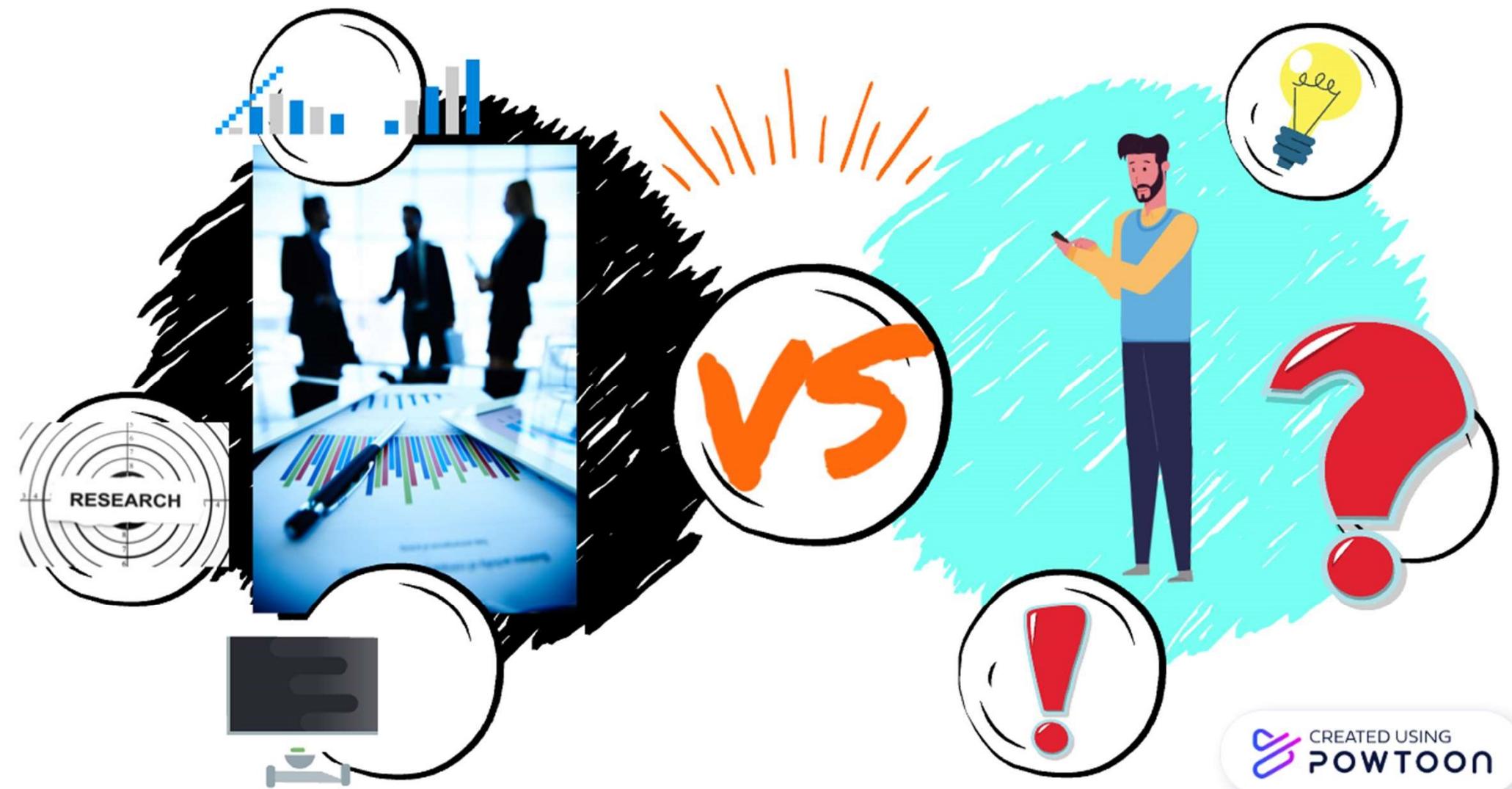
*The Stock Traders of
Reddit and TikTok*
-- The New York
Times

*"Robinhood's
Addictive App
Made
Trading a Pandemic
Pastime"*

-- Bloomberg

*The
Wall
Street
Bets*
-- Financial Times







Market Color
Investor Sentiment

Bullish

Panic

SELL

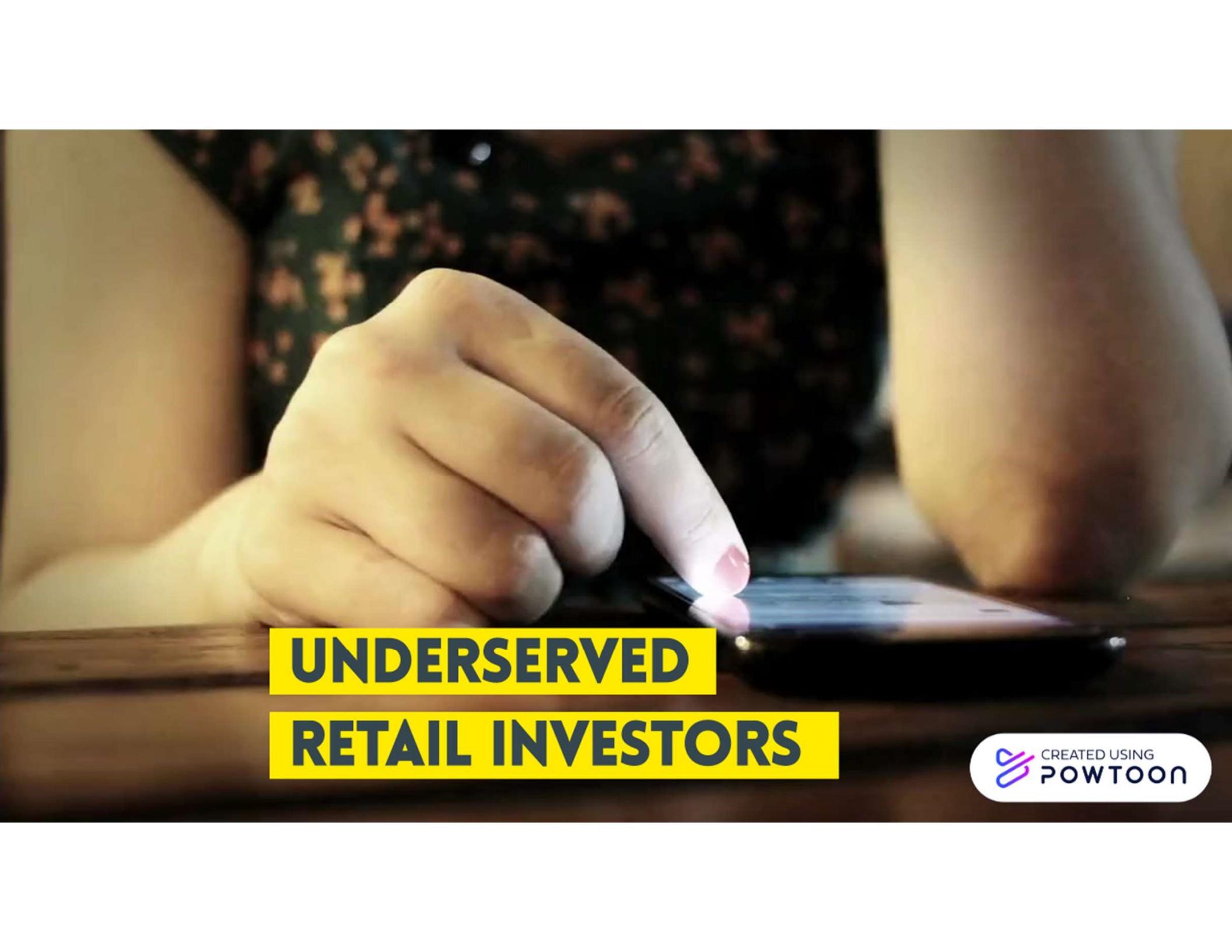
Strong Buy

Hold

Two ~~Gaps~~

Opportunities

Identified...



UNDERSERVED RETAIL INVESTORS



Institutional Investors

Social Media Sentiment



Now let's formally define the problem...

Problem Formulation

"What is the current sentiment on social media towards stock X?"

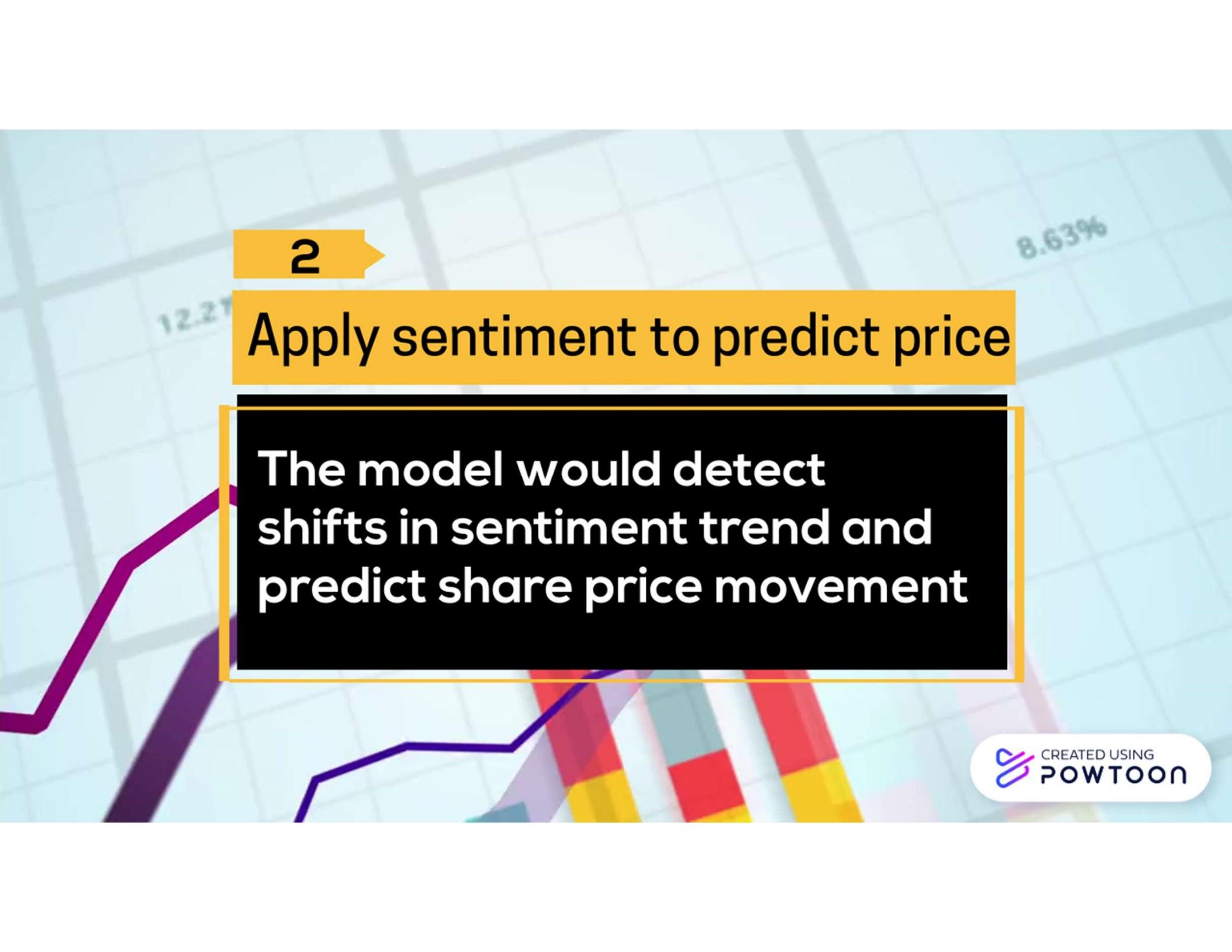
"How would th

Project Objectives

1

Create a prototype pipeline

... that programmatically collects Twitter feed data on specific companies and automatically extracts sentiments from them.



2

Apply sentiment to predict price

The model would detect shifts in sentiment trend and predict share price movement

V e], U (" y { } m S 9' b1 - ^ 0 MC B2 s " > 9 6 w % H
9 2 b1 S A , a q] Z - f q | 6 B #) i @ 5 g M - 3 I g
j { * ; j A (, K + o } : l S ! ? u F Q Q O P W _ P z u
0 . j { ' - W M > L k W Z a l r l 3 2 J p m I o) 0 { ^ W
0 m + l d @ 3) c b) p n e l , m y * [0 I X 9 Y p 1
?) w c * > ^ 1 I I , i | F o F , N ? P } : Z , - j I T R
> q } Z n X < i > + ? m 4 " l i e K R = A D T } 6 _ k 6 J
u # b S - [T w & A o o l = f g E S { ^ > L ; 0 I 6
R ! b , C " \ k ' t ' z 0 z 1 J w q /
" U (T M N N # G F R u 6 7 3 \ ? \$ 6
0 ; I o p . f 7 g } / % r - % @ 5 K h U P
[I H I S P & S ! X z A m - 6 W : > d @ ^ r c 6 y d : T
i ; m N J . 2 J h ' [

```
|64 bytes from 173.194.115.2: icmp_seq=220 ttl=57 t  
|ime=11.971 ms  
|64 bytes from 173.194.115.2: icmp_seq=221 ttl=57 t  
|ime=9.943 ms  
|64 bytes from 173.194.115.2: icmp_seq=222 ttl=57 t  
|ime=9.904 ms  
|64 bytes from 173.194.115.2: icmp_seq=223 ttl=57 t  
|ime=11.735 ms  
|64 bytes from 173.194.115.2: icmp_seq=224 ttl=57 t  
|ime=9.866 ms  
|64 bytes from 173.194.115.2: icmp_seq=225 ttl=57 t  
|ime=11.284 ms
```

Explore NLP and ML tools

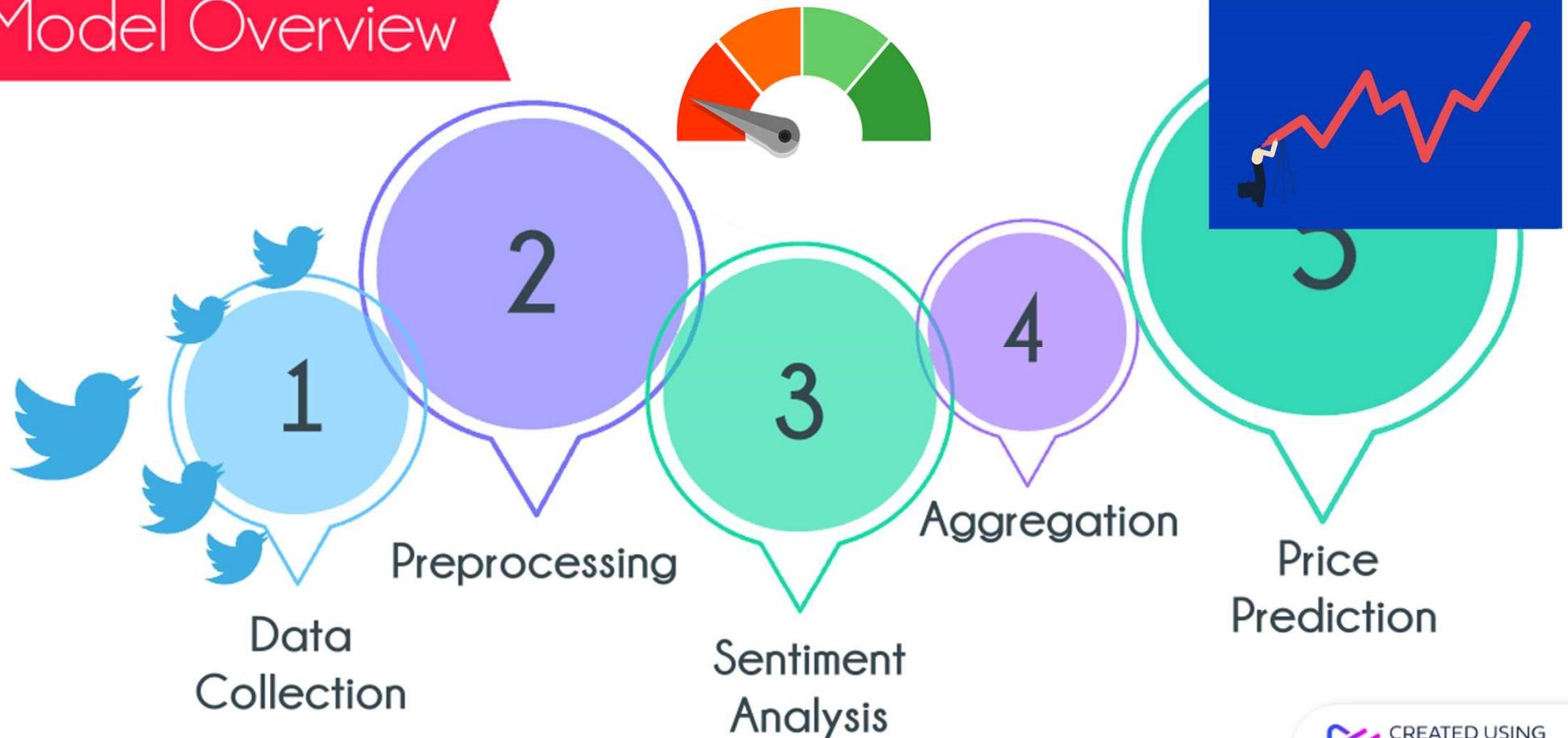
... and assess their applicability as a feasible, practical solution to a real world problem.

```
139          n W 9 s u   p s u 0 R g    ?    M J CPU uso2e: 55.35% user5, 2.38% sys, 18.90% idle !  
140          U 9 8 -   h k   >   K @ C , d CPU usage: 55.55% user, 25.46% sys, 18.98% idle !  
141          [ < n Ce = 3 r o ^ ! S R U a " ) T 609 28 916 1 85  
142          [ 7 G 4 & 9 g { 9 c l ] $ e z ) X 0832 5 35 2304 107M unuse6  
143          < w s Y C & x E L Z t ) g : h = K 1 5 594 6  
<44          y 6 Y D & x E g R 3 1 H o 4 3 38  
145          r C < i > Q W b i L q z , G ` U v 54 56 17 44 *  
146          Z l y g [ > Q 2 j ] w T N 1 60 32  
147          6 7 p s i F q f z j . R } ) 1  
148          D p s v o D . t = o d @ , c ) g ^  
#          Q x = D 3 % t z 3 S 4 J  
Y          $ K J a 3 S  
          z b -6
```



Model Overview

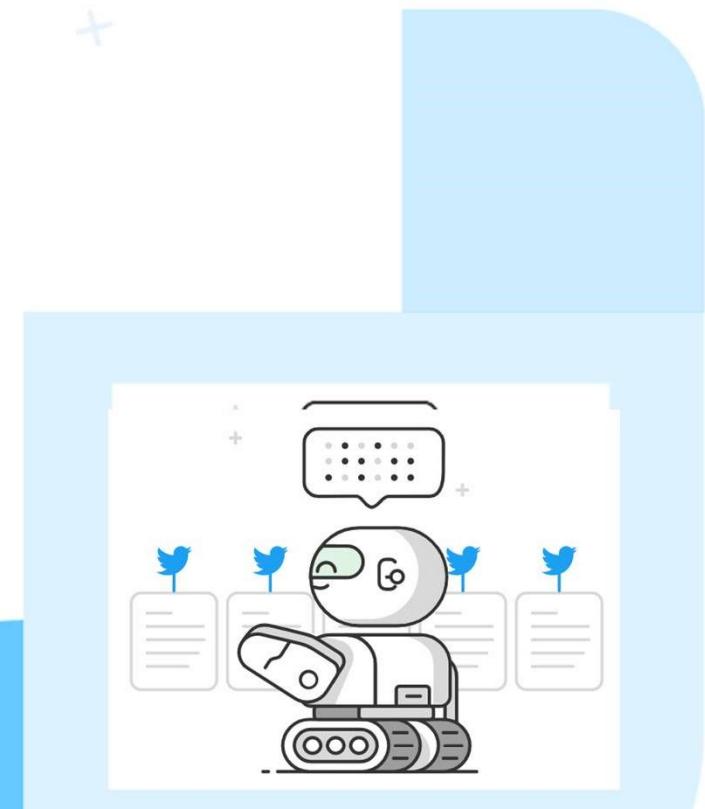
Model Overview



01

Data Collection

- Raw Tweets: Twitter API
- Annotated Dataset: SemEval 2016
- Share Price: Yahoo!Finance



02

Preprocessing

- url, @username, #
- negation
- elongated words
- emojis and emoticons
- acronyms and slangs
- special characters (non-punctuation)
-
- spelling, stopwords, stemming



03

Sentiment Analysis

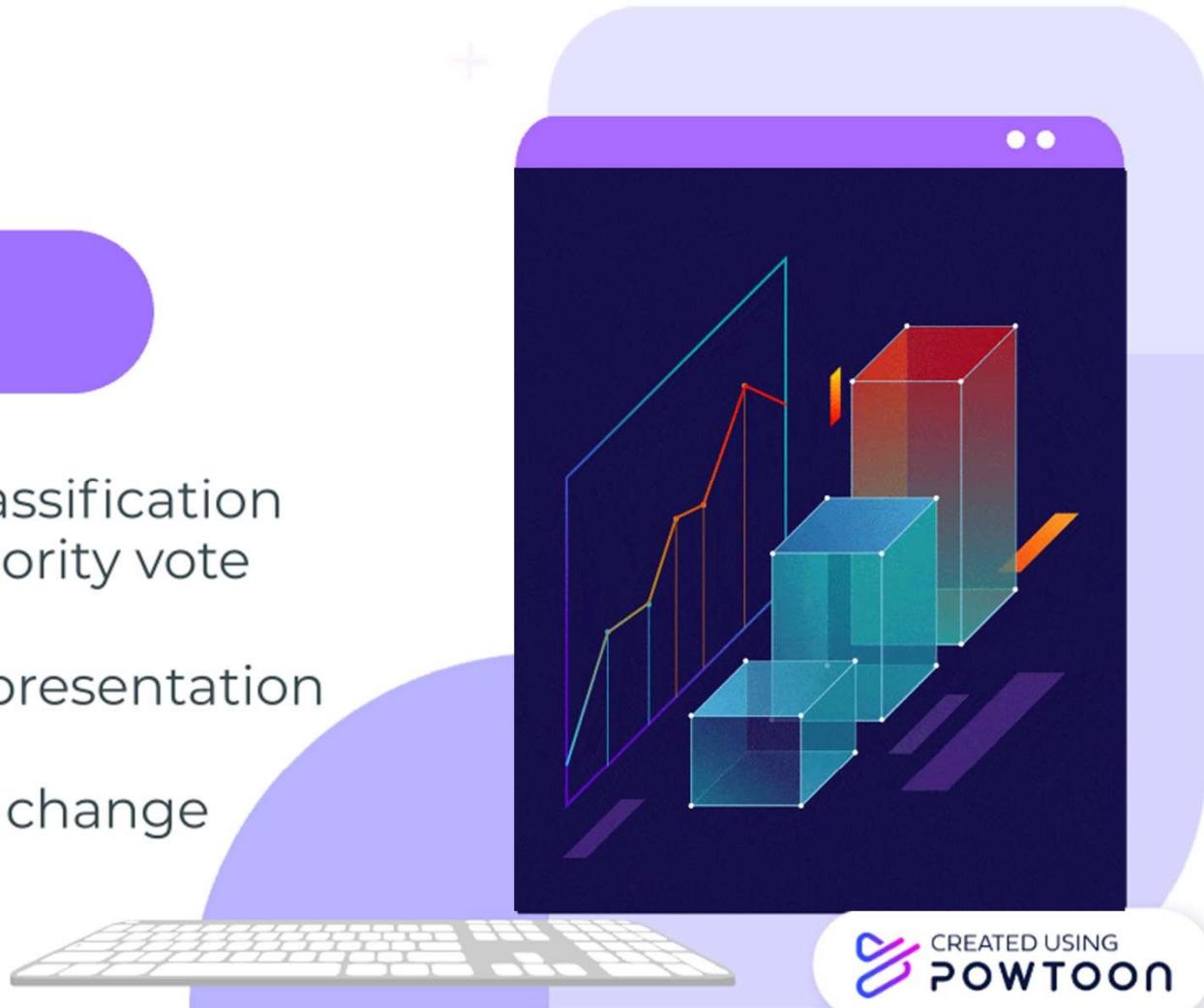
- Lexicon / Rule-based
- Classical ML-based:
 - * Naive Bayes
 - * Logistic Regression
 - * Support Vector Machine
- BERT with fine-tuning



04

Aggregation

- Select best performing classification models; ensemble by majority vote
- Sentiment vector, shift representation
- Relative percentage price change



05

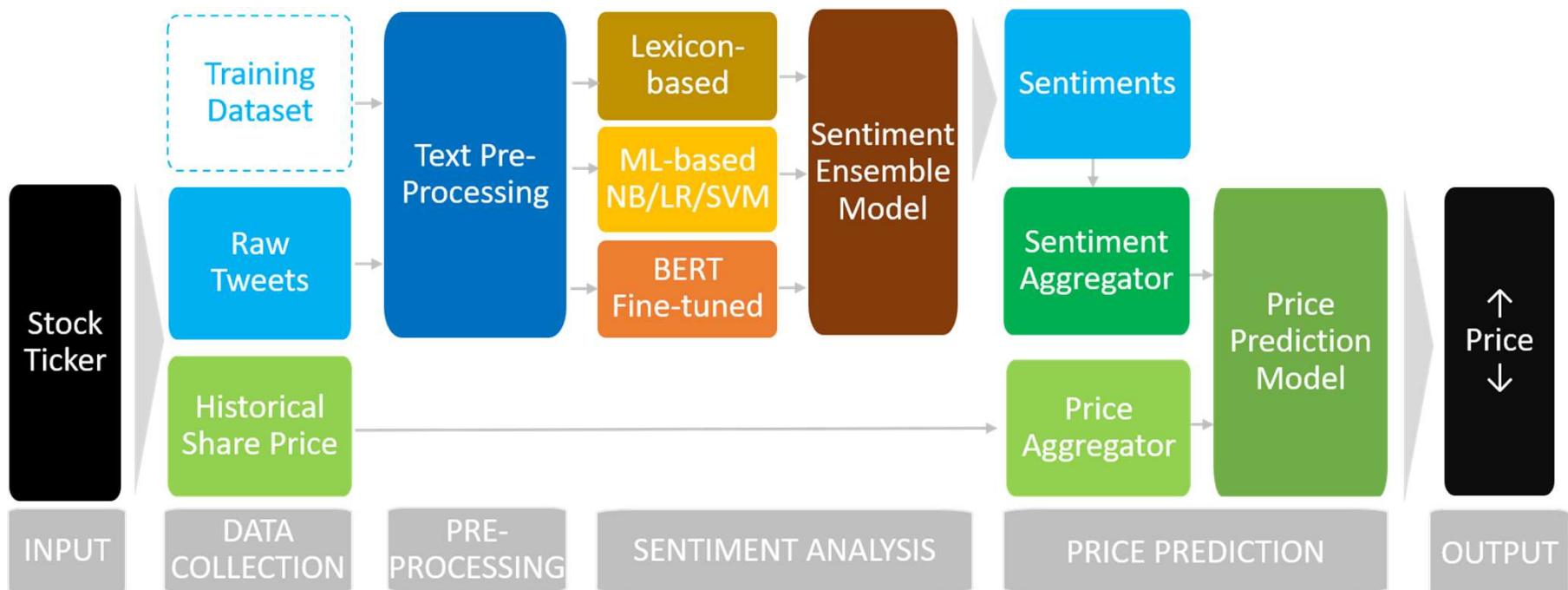
Price Prediction

- RNN / LSTM / CNN
- Work-in-Progress
- Only 39 trading days

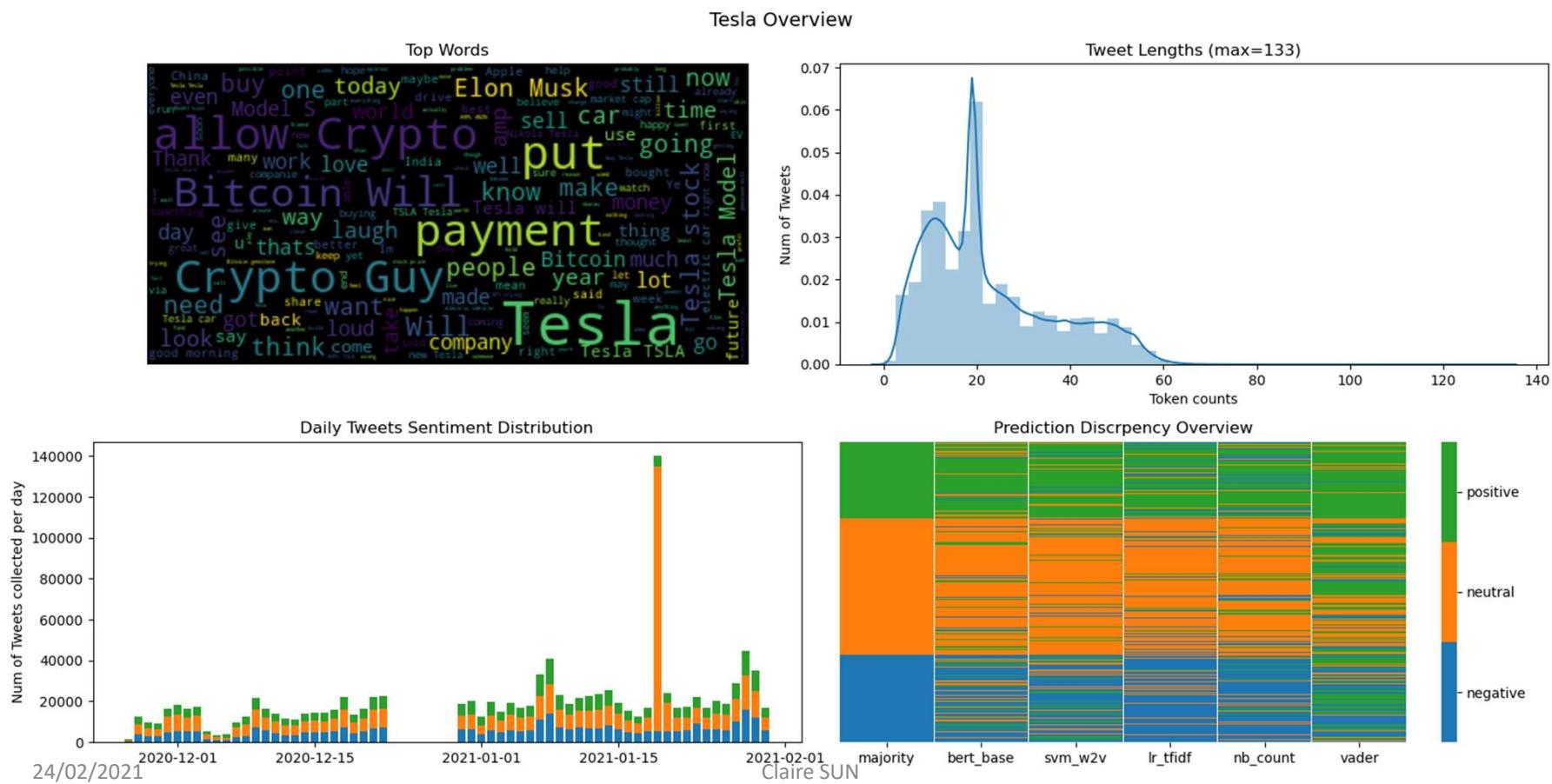


Experiment Setup

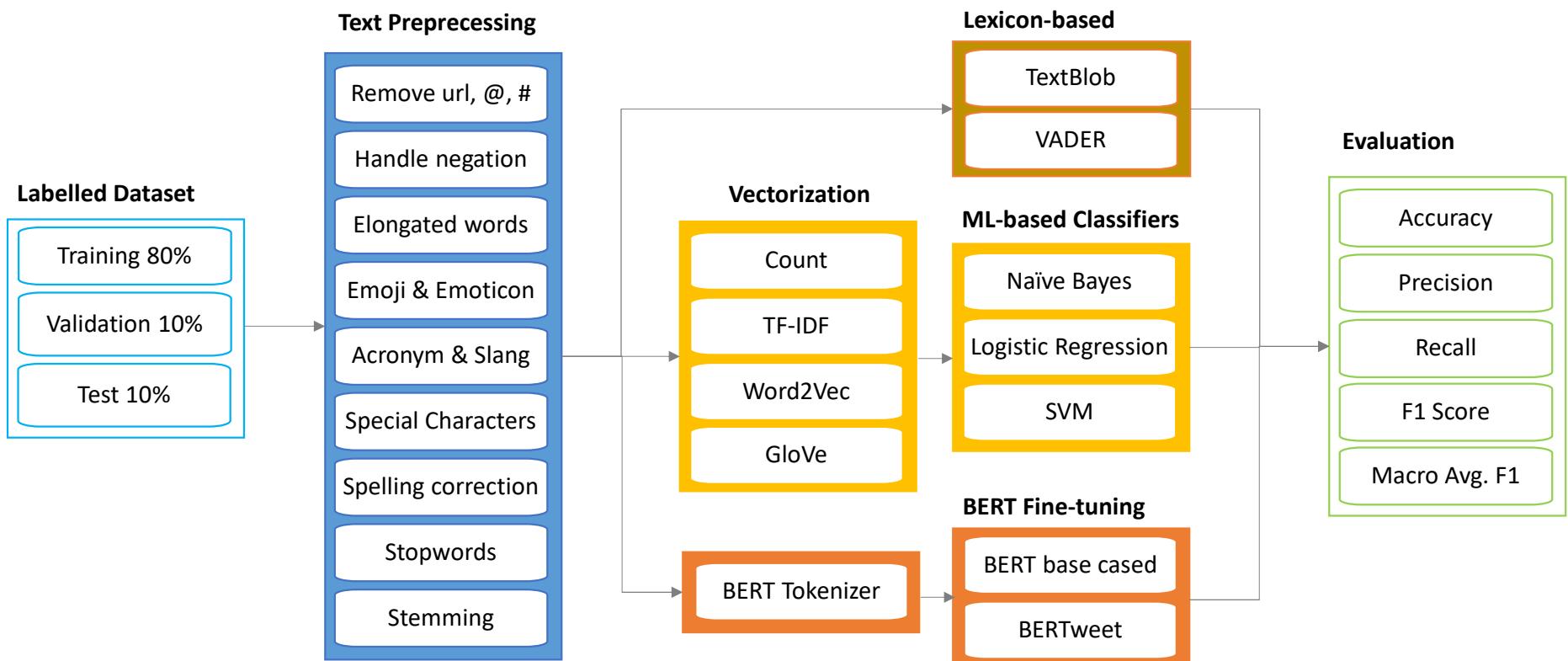
Project Planning



Twitter Data Visualization

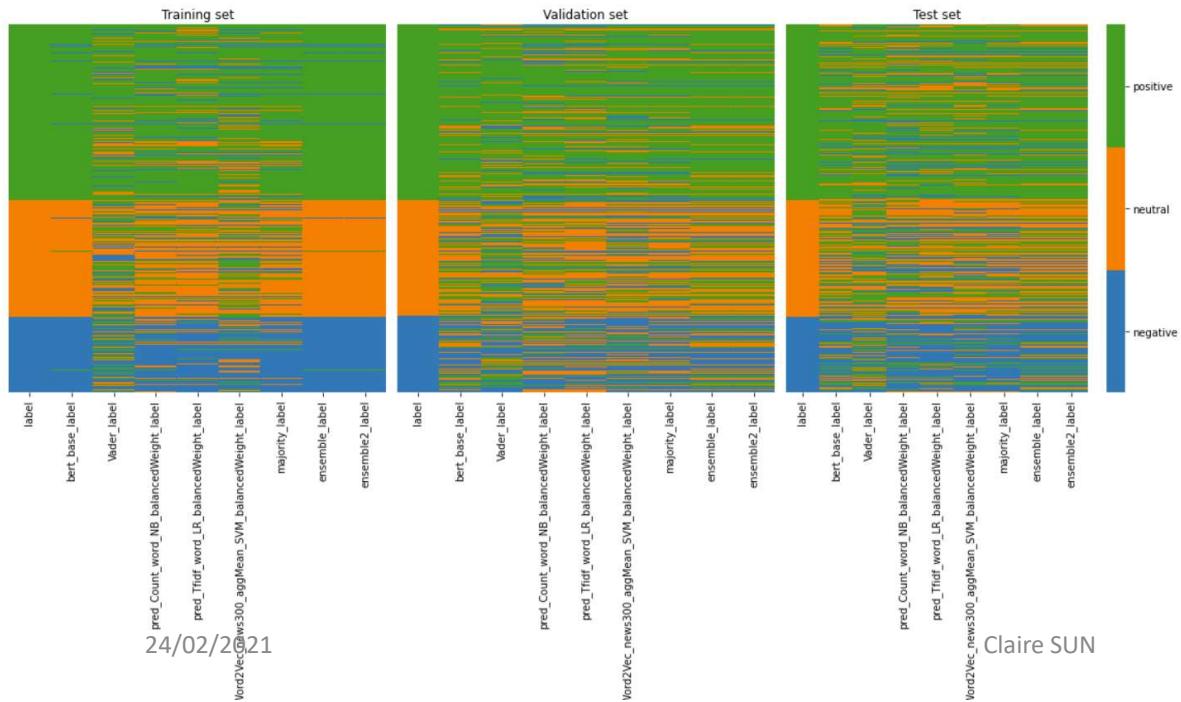


Experiment Setup



Model Training & Selection

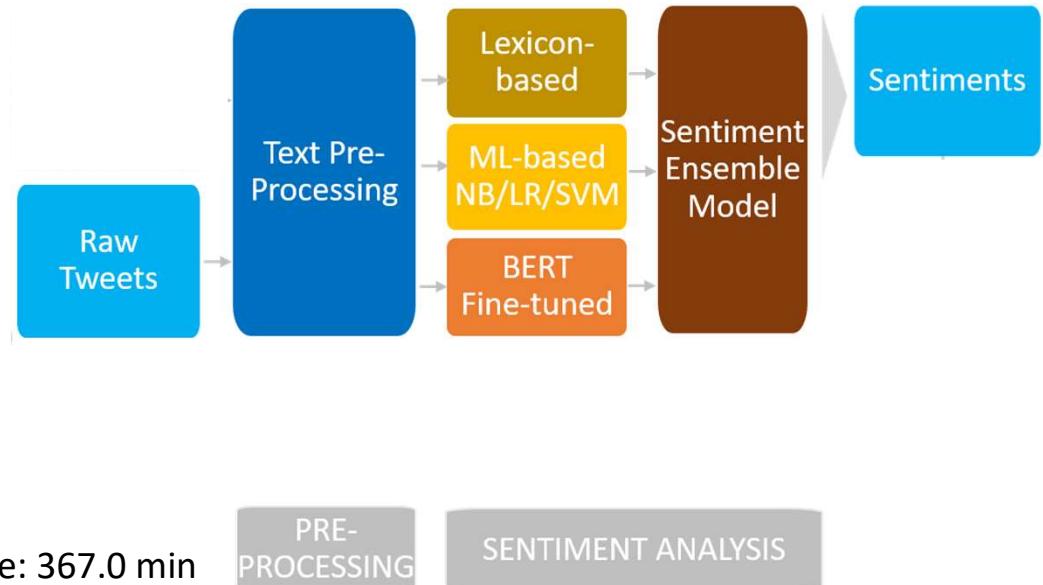
Sentiment	Count	%Total	id	text	sentiment
POSITIVE	9,388	47.8%	634633252020555781	"b""HERE today, gone tomorrow- but still here! A short note on Nokia's patent deals with @Microsoft and @Alcatel_Lucent http://t.co/y5wFgUygFD """	neutral
NEUTRAL	6,232	31.7%			
NEGATIVE	4,033	20.5%	637308387353432064	b@joebelfiore @GabeAul @Microsoft @Lumia @satyanadella plz add L1520 in the 1st wave of windows10 phones release.plz dont hurt ur diehardfans'	neutral
TOTAL	19,653	100.0%			



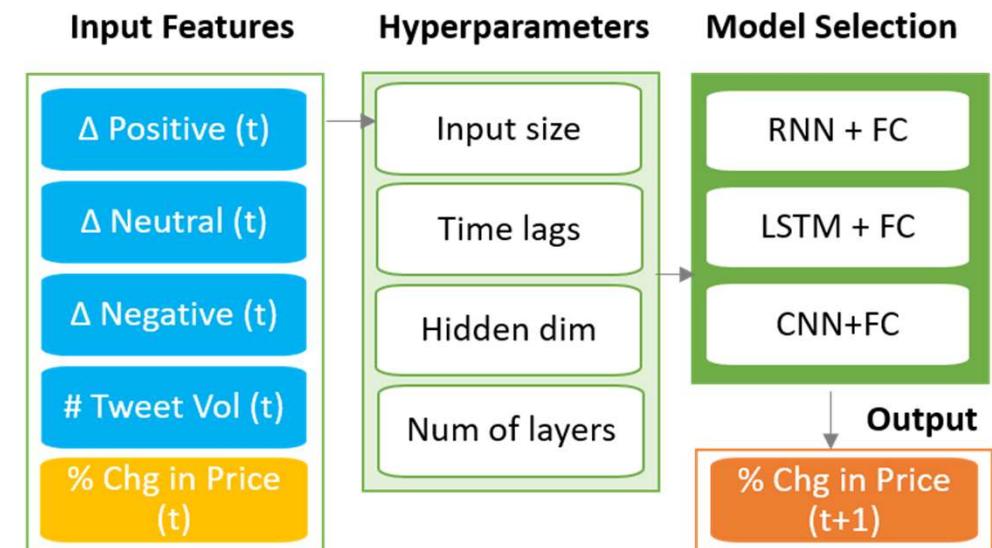
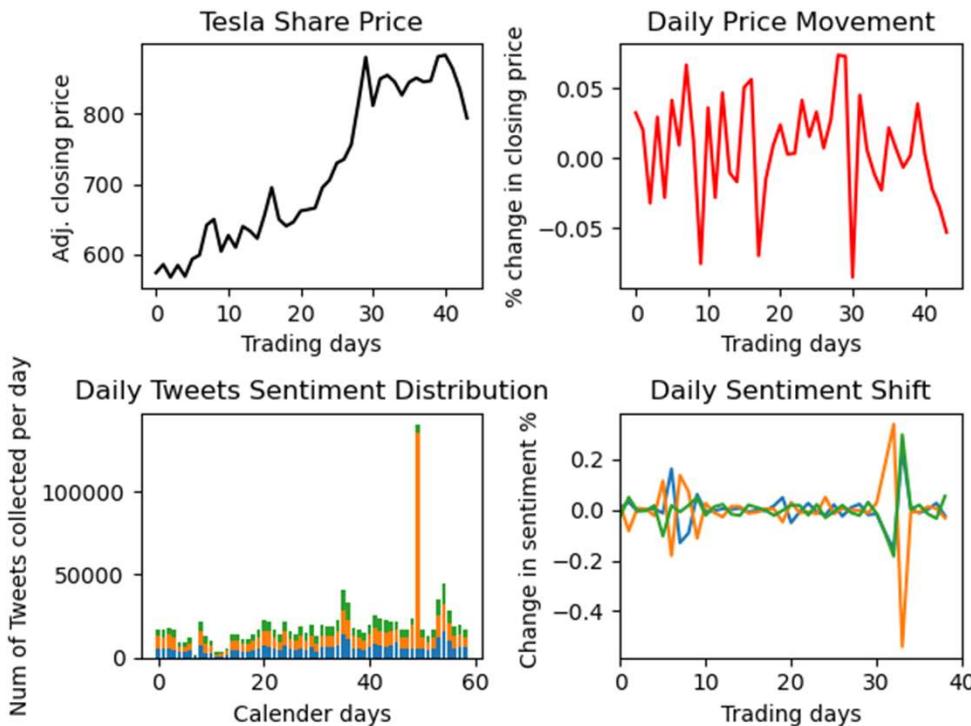
Classification Method	Macro Average		
	Precision	Recall	F1-score
Vader	54%	54%	53%
<i>TextBlob</i>	48%	48%	48%
<i>Count_NB (SMOTE)</i>	58%	60%	59%
<i>Count_LR (SMOTE)</i>	55%	56%	55%
<i>Count_SVM (SMOTE)</i>	53%	53%	53%
Count_NB (balanced_weight)	59%	60%	59%
<i>Count_LR (balanced_weight)</i>	57%	58%	57%
<i>Count_SVM (balanced_weight)</i>	52%	53%	53%
<i>Count_NB_unigram,bigram,trigram</i>	57%	58%	57%
<i>Count_LR_unigram,bigram,trigram</i>	57%	57%	57%
<i>Count_SVM_unigram,bigram,trigram</i>	53%	52%	52%
TFIDF_NB (balanced_weight)	59%	58%	58%
TFIDF_LR (balanced_weight)	60%	61%	60%
TFIDF_SVM (balanced_weight)	57%	57%	57%
<i>TFIDF_NB_bigram</i>	52%	53%	53%
<i>TFIDF_LR_bigram</i>	52%	53%	52%
<i>TFIDF_SVM_bigram</i>	50%	50%	50%
<i>GloVe-Twitter-25_NB</i>	48%	49%	44%
<i>GloVe-Twitter-25_LR</i>	51%	53%	51%
<i>GloVe-Twitter-25_SVM</i>	51%	52%	51%
<i>GloVe-Twitter-200_NB</i>	51%	51%	48%
<i>GloVe-Twitter-200_LR</i>	55%	57%	55%
<i>GloVe-Twitter-200_SVM</i>	55%	56%	56%
<i>Word2Vec-News-300_NB</i>	51%	52%	48%
<i>Word2Vec-News-300_LR</i>	60%	62%	60%
Word2Vec-News-300_SVM	61%	62%	61%
BERTweet	29%	30%	20%
BERT Base Cased	64%	64%	64%
<i>Ensemble_FC_Softmax</i>	64%	64%	64%
<i>Ensemble_MajorityVote</i>	64%	66%	28
Ensemble_MajorityVote	64%	66%	64%

Processing Time

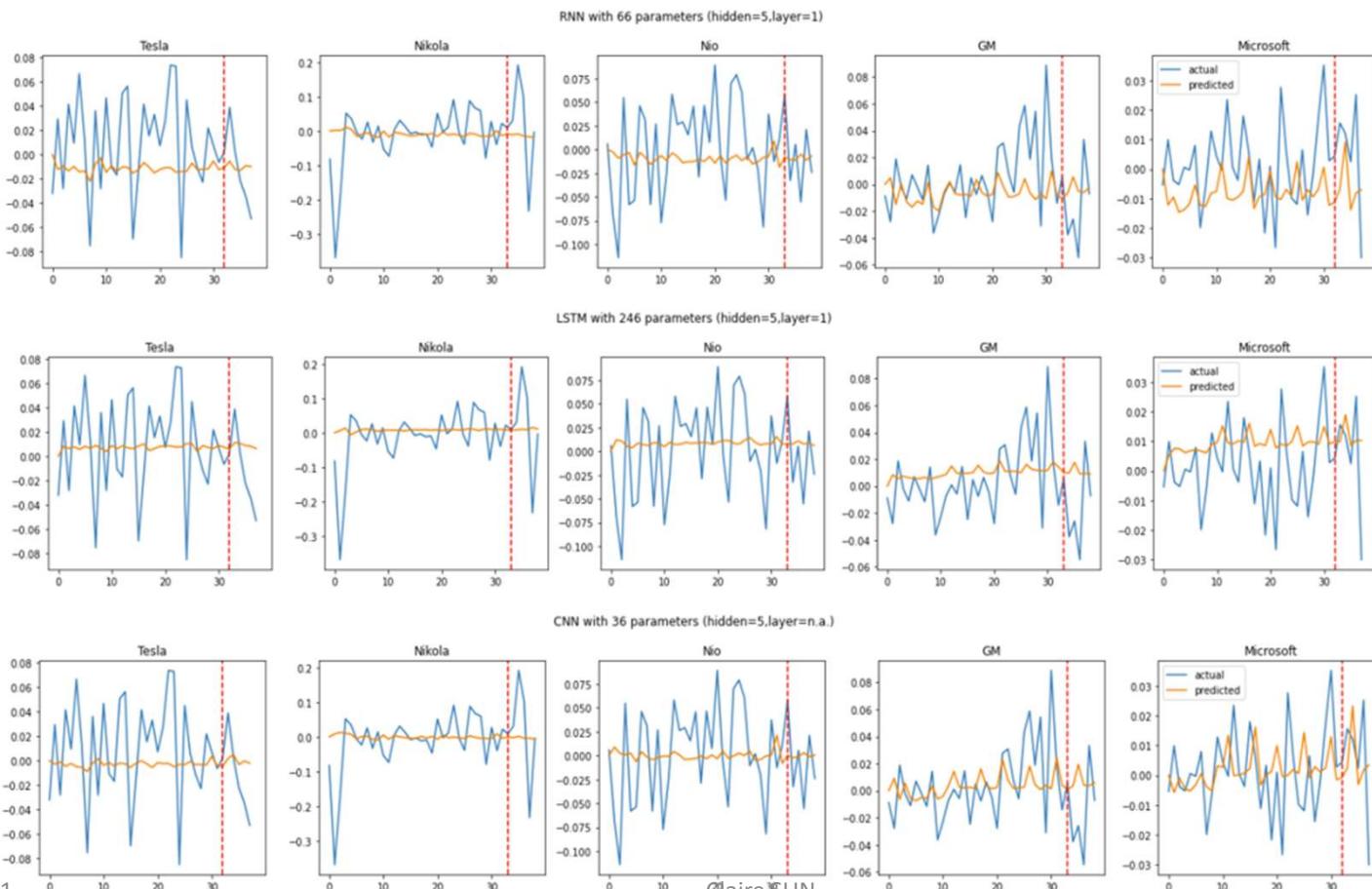
- ### Tesla ###
- Loading raw tweets from 9 files:
 >> loading: 4.896 s
- Start processing 1181205 raw tweets:
 >> text preprocessing: 1414.389 s
- Start sentiment analysis on 1181205 tweets:
 >> vader: 125.264 s
 >> nb_count: 18.192 s
 >> lr_tfidf: 18.299 s
 >> svm_w2v: 114.864 s
 >> bert: 19473.192 s (cuda)
 >> majority: 346.132 s
- Sentiment analysis completed! Total processing time: 367.0 min



Price Prediction



Price Prediction



Demo

Input

Stock ticker: **TSLA** Prediction date: **2021-02-19**

Data Collection

- ✓ Share price downloading completed in **1.8s**
- ✓ Tweets downloading completed in **26.8 min**

Preprocessing

- ✓ Start processing **50,159** raw tweets:
- ✓ Text preprocessing completed in **63.2s**

Sentiment Analysis

- ✓ VADER Analyzer: **5.5s**
- ✓ Naïve Bayes_CountVectorizer: **0.8s**
- ✓ Logistic Regression_TFIDF: **0.8s**
- ✓ SVM_Word2Vect: **17.4s**
- ✓ BERT base cased: **915.7s**
- ✓ Majority Vote: **17.1s**

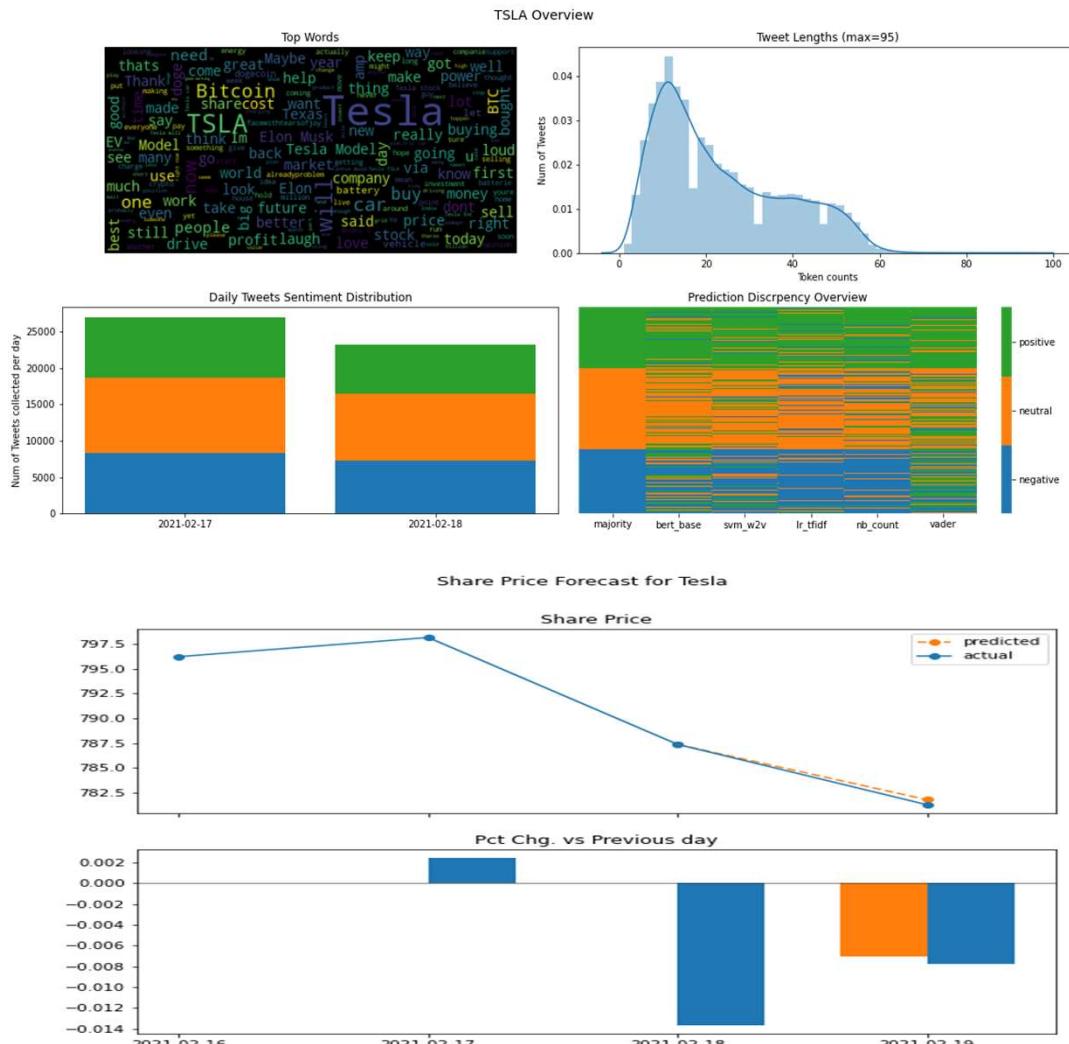
- Sentiment analysis completed! Total processing time: **17.0 min**

Price Prediction

- ✓ Share price prediction completed in **0.3s**

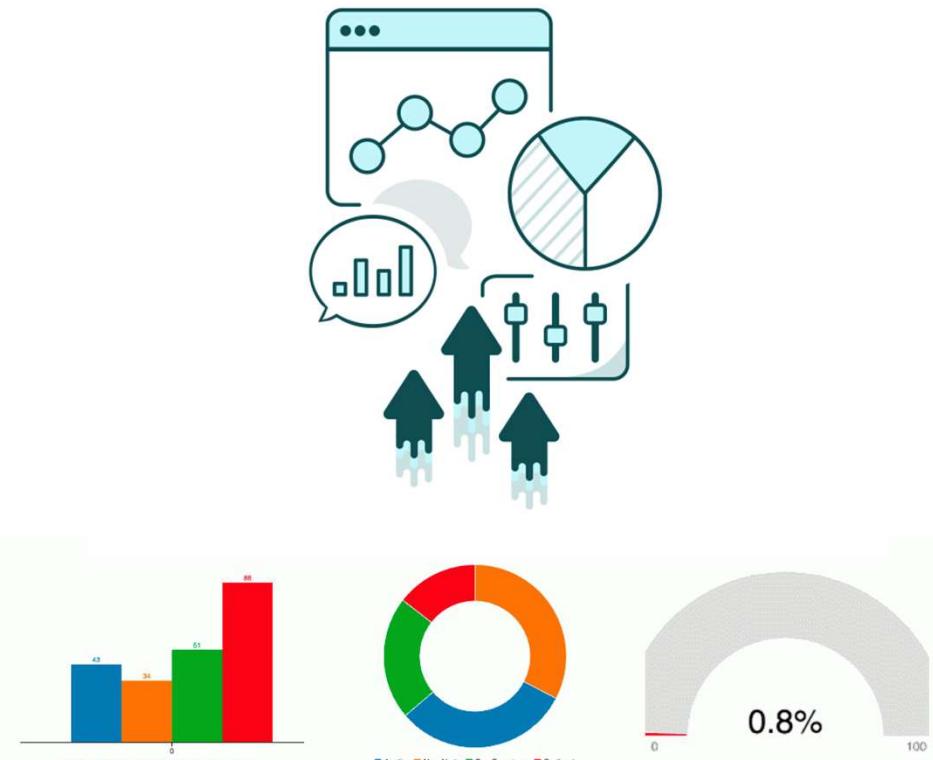
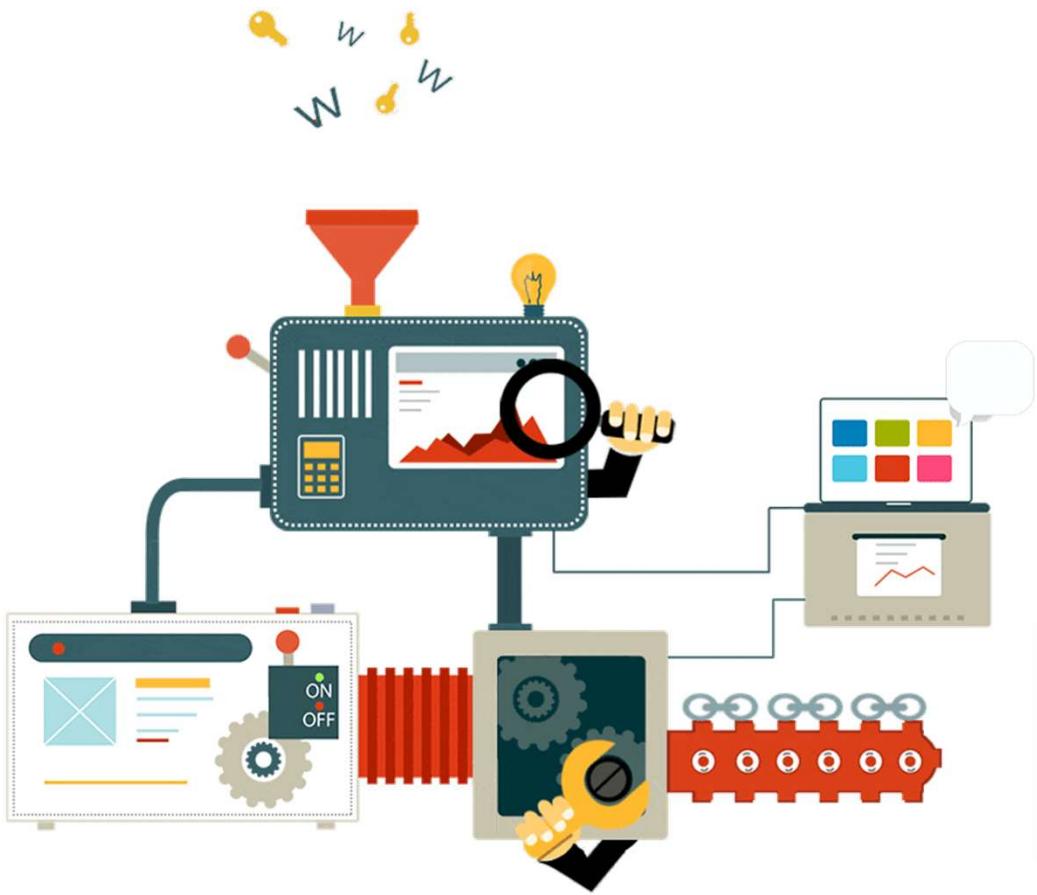
Output

- Predicted next day price movement for Tesla: **-0.71%**
- Prediction: **\$781.80/share vs actual closing \$781.30/share**



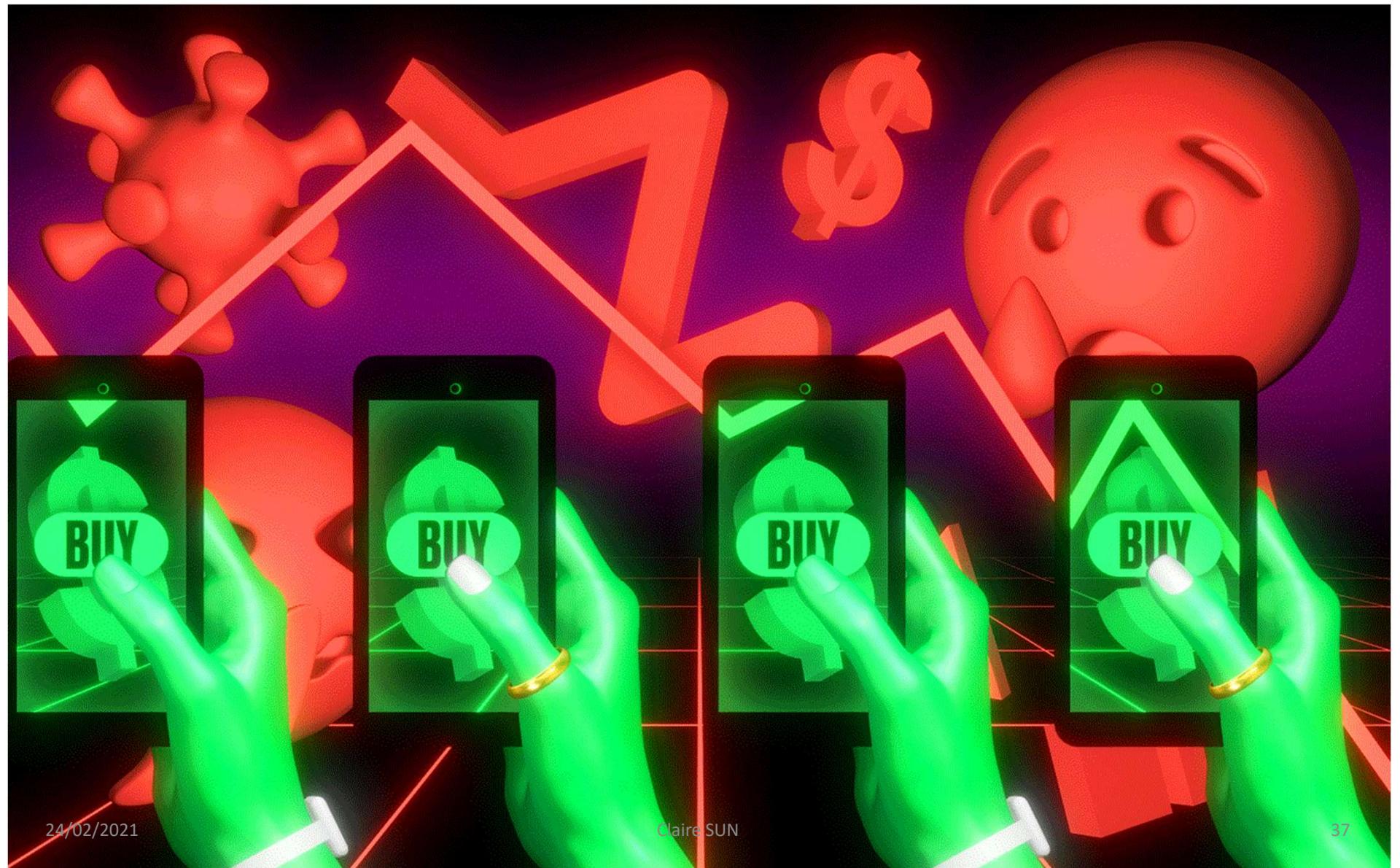
Challenges & Key Learnings

- **Pilot project to evaluate feasibility**
 - Demo is a toy example; price prediction model is still WIP
- **Domain-specific Training Dataset:**
 - Potential benefit from developing proprietary, annotated Twitter dataset vs. time/costs
- **Data collection over a longer period of time:**
 - Or subscription to Twitter's Enterprise API for some proper back-testing
- **Model speed vs accuracy:**
 - BERT model is c. 4-5% more accurate at sentiment prediction during training than NB/LR but c. 1000x slower at inferencing
 - Could potentially try other lighter-version of BERT such as DistilBERT
- **Potential improvement areas in the next iteration:**
 - Incorporating extra features: Not all collected data is consumed by the current model. For example, we have trading volume data as well as Twitter metadata that we could apply to potentially enhance the model architecture
 - Iteratively refining and optimizing the pipeline in order to accelerate processing speed and increase data-handling capacity
 - Improve design for user-interface



Conclusion

- From research findings into a real world problem solution
- End-to-end sentiment analysis and price prediction pipeline
- Key concepts and methodologies in NLP and ML
- Challenges of applying research techniques to a practical problem
- Improvement roadmap for the next iteration



24/02/2021

Claire SUN

37