



# CORPORATE STOCK REVIEW !



Introduction to Artificial Intelligence - IT3160E



Nguyen Chi Long - Ngo Xuan Bach  
Le Xuan Hieu - Dinh Viet Quang

## PROBLEM STATEMENT

Stock analysis, which has always been a crucial aspect of any organization, is a method for investors and traders to make buying and selling decisions.

By studying and evaluating past and current data, investors and traders attempt to gain an edge in the markets by making informed decisions

## OUR TARGETS

Evaluate how the sentiment impacts stock performance through article text.

Assign a large set of corporate stock reviews to the several levels of polarity of opinion as accurately as possible.

[Back to Agenda Page](#)



# APPROACH

**Sentiment Analysis Techniques**

1	Explore the our collected data
2	Choose the metric to score our models: accuracy and confusion maxtrix
3	Compare the results of each model and find the best model
4	Load the models for demonstration and deploiment

# Data Gathering Phases

Scrapy



1

Data crawling

Doccano



2

Data annotation

3

Data processing and  
Model training



Python

Demo and deployment

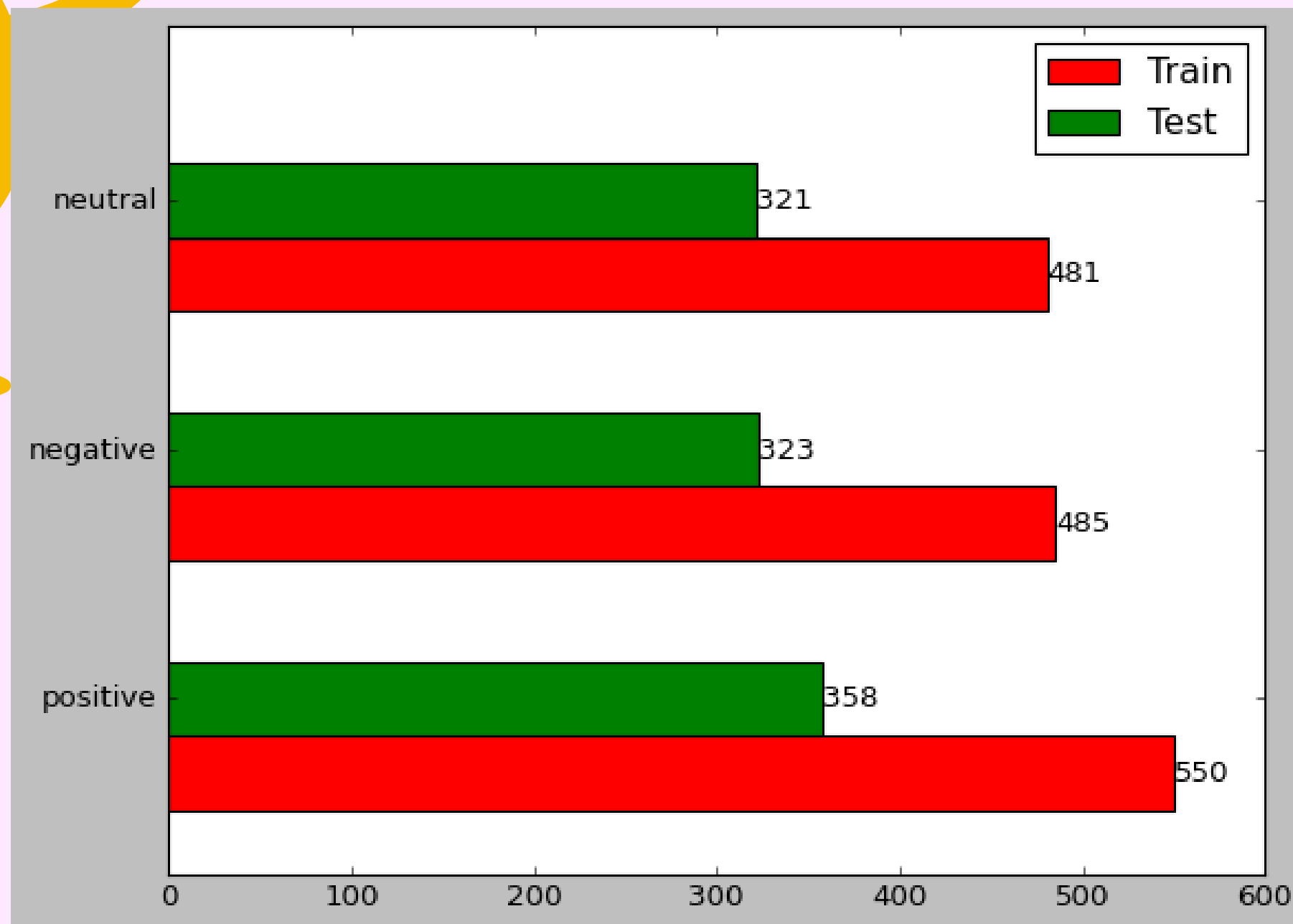
jupyter



Flask

4

[Back to Agenda Page](#)



# OUR DATASET

Collect corporate stock reviews from various sources (CafeF, VnExpress etc.)

The dataset contains 2500 additional unlabeled data. The number of positive, negative and neutral sentiment is more 800

Data split: 60% Train - 40% Test

[Back to Agenda Page](#)

# EXAMPLE

“**Positive**” review: “Nhiều mã ngành này lọt vào nhóm góp mức tăng lớn nhất cho toàn thị trường, thị giá kịch trần ngay từ đầu buổi sáng”

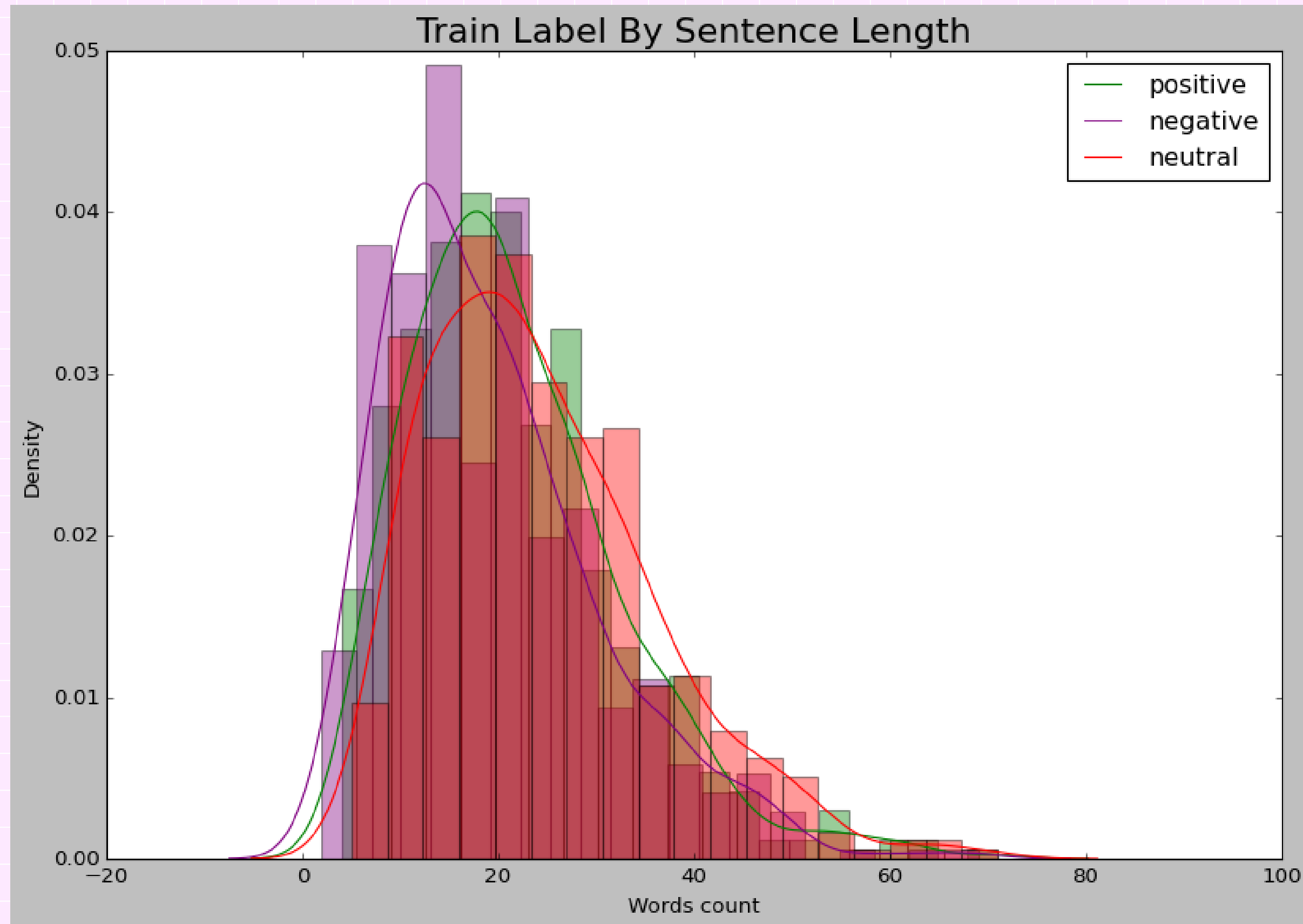
“**Negative**” review: “Trong thông báo hôm nay, lãnh đạo FLC đã trình bày một loạt khó khăn khiến doanh nghiệp không thể công bố các thông tin theo đúng quy định sau khi cựu Chủ tịch HĐQT Trịnh Văn Quyết bị bắt cuối tháng 3/2022”

“**Neutral**” review: “Đóng cửa phiên giao dịch ngày 20/4, cổ phiếu NT2 có giá 29.400 đồng/cổ phiếu”

[Back to Agenda Page](#)

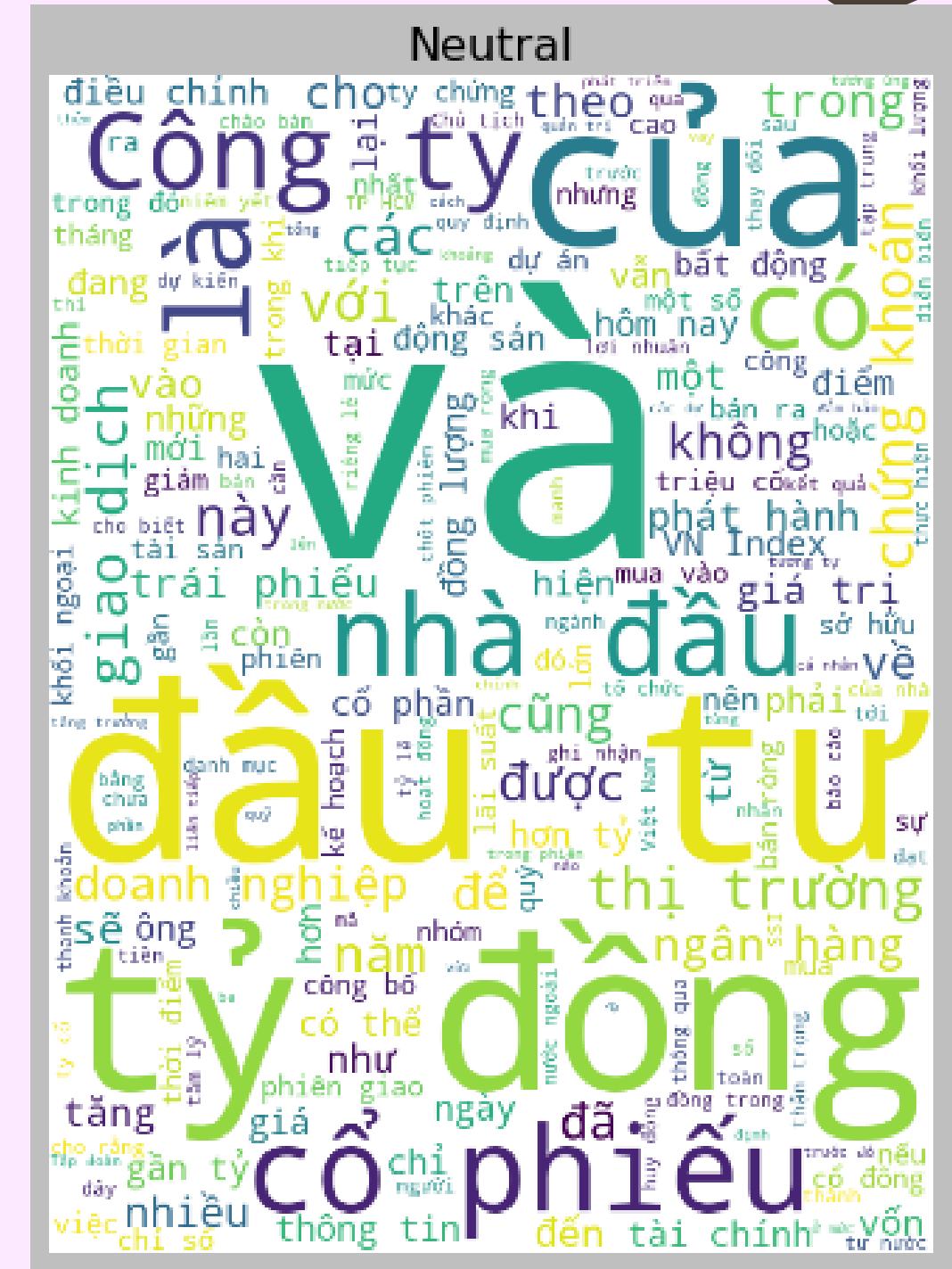
# SENTENCE LENGTH

Good  
Job!!





# WORD FREQUENCY



[Back to Agenda Page](#)





# DATA PREPROCESSING

## TF - IDF



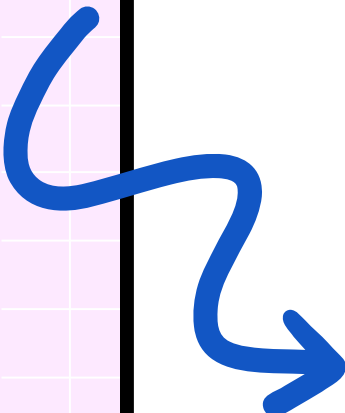
$$\text{TFIDF}(t, d, N) = \text{tf}(t, d) \times \log(N/\text{df}(d, t)+1)$$



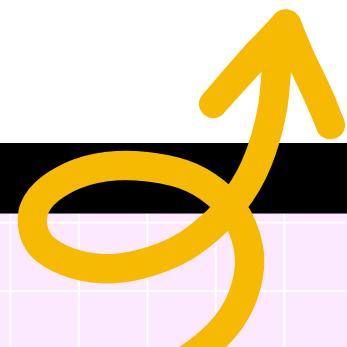
TF - IDF

TERM  $t$  WITHIN DOCUMENT  $d$

$\text{tf}(t, d)$  = frequency of  $t$  in  $d$   
 $\text{df}(d, t)$  = number of documents  $d$  containing  $t$ .  
 $N$  = total number of documents.



Term frequency - Inverse document frequency: is a numerical statistic that is intended to reflect how important a word is to a document in a collection or corpus.



# MODELING

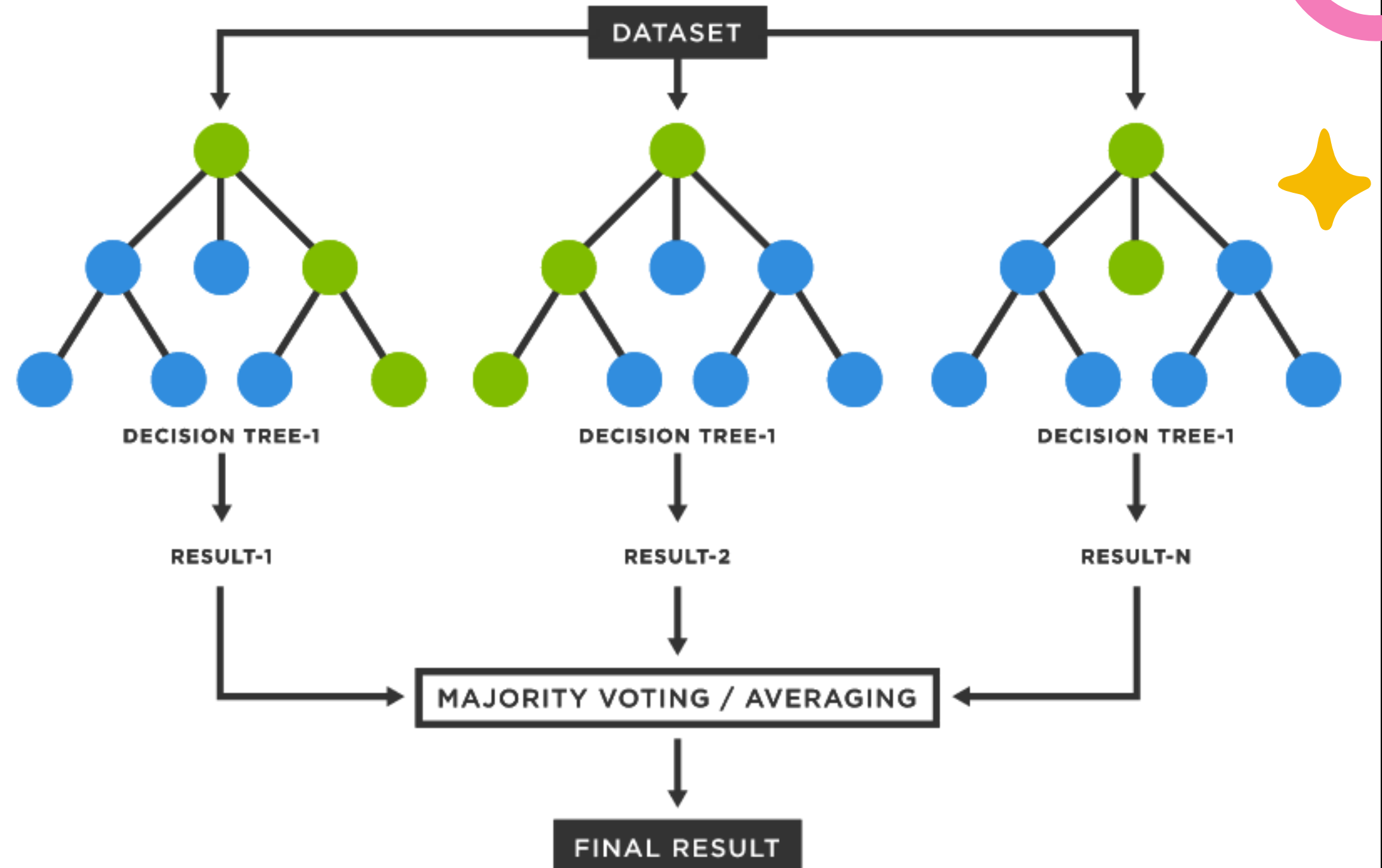
## Random Forest

Grows multiple classifier trees with bootstrap sampling with replacement to train each tree with a different part of the dataset.

Three main hyperparameters:  
+ Depth of the trees: 50.  
+ Max features: "sqrt".  
+ Random state

Criterion = 'Gini'

[Back to Agenda Page](#)



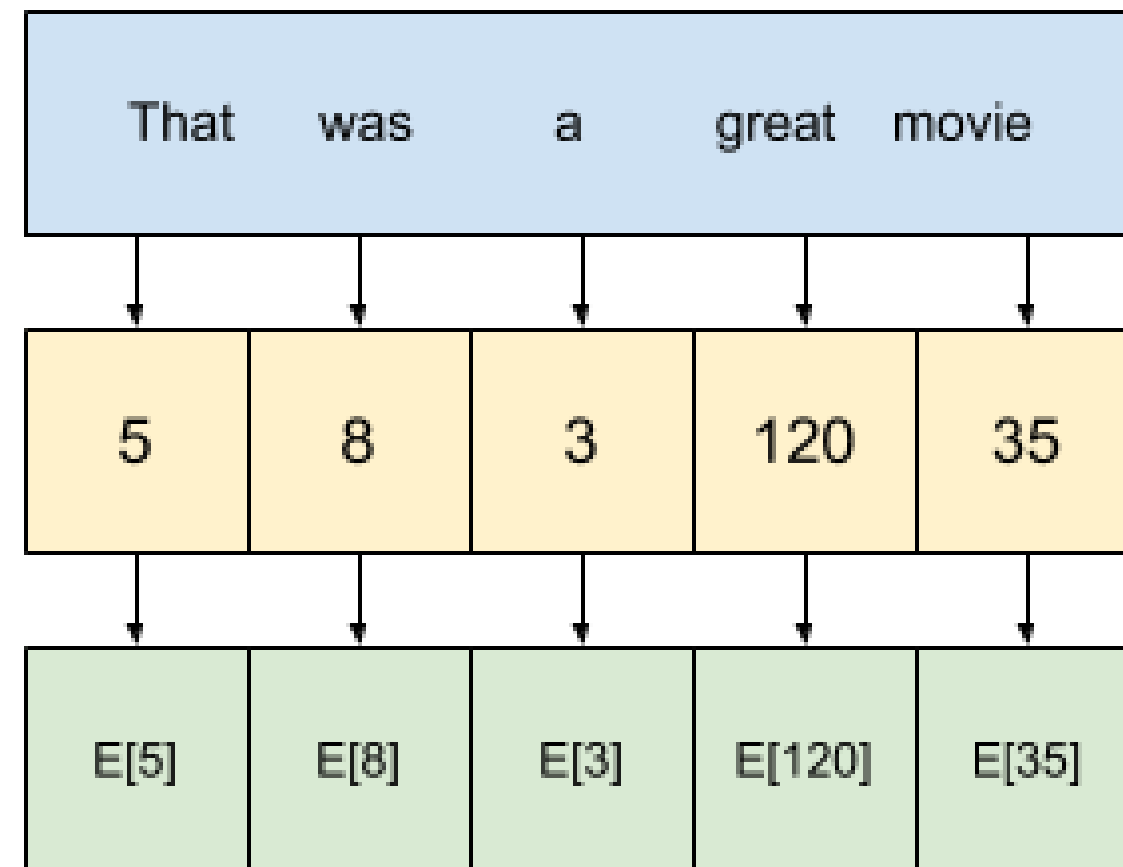
# MODELING

## Neural Network

### PREPROCESSING: WORD EMBEDDING

Lower and strip punctuation,  
turn word into index

Turn positive integers (indexes)  
into dense vectors of fixed size



Input

TextVectorization

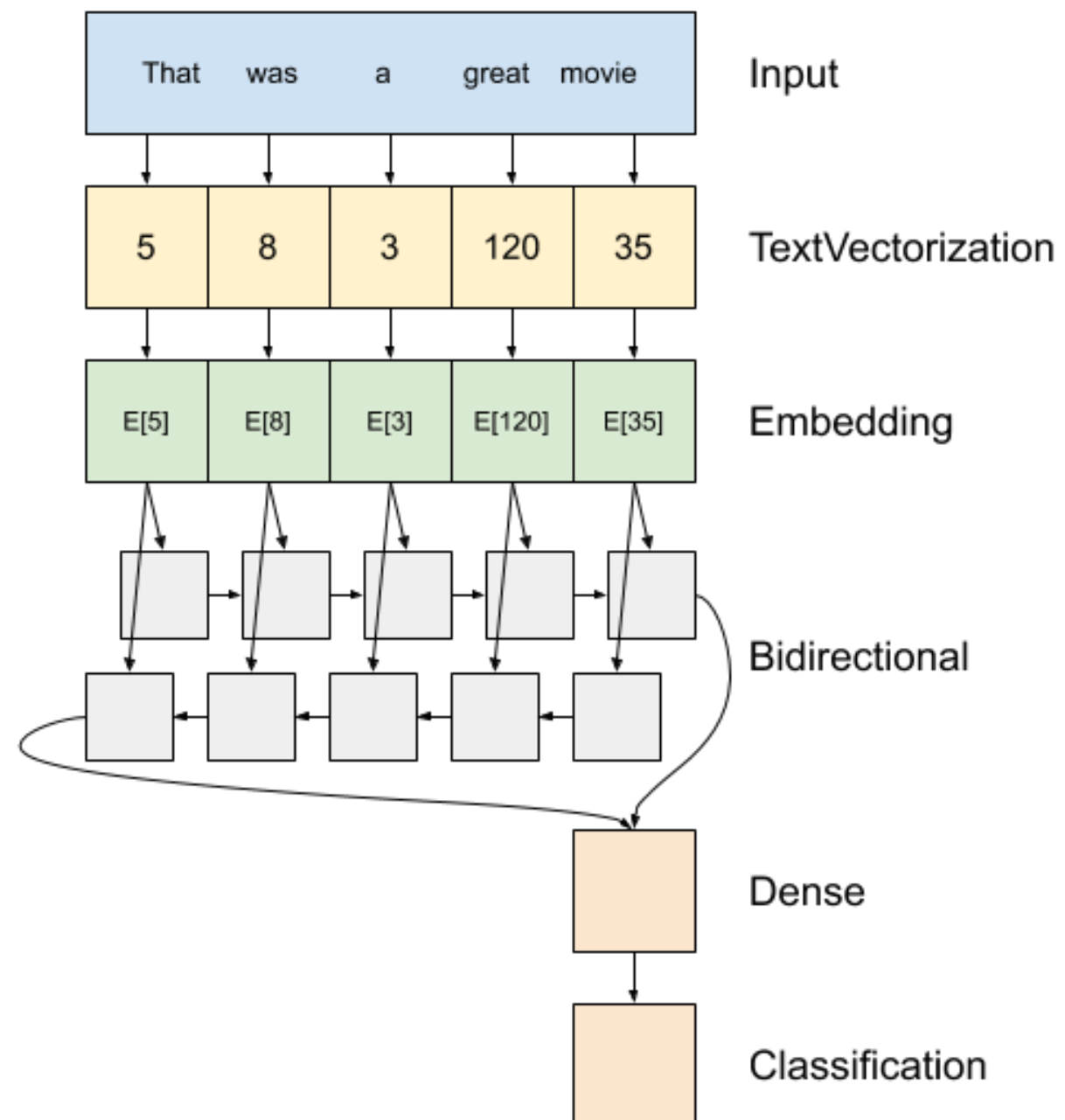
Embedding

Word embeddings are a form of word representation designed to capture semantic relationships, wherein words that are closer together in the vector space are generally assumed to have similar meanings

# MODELING

## Neural Network

### Deep Bidirectional Architecture



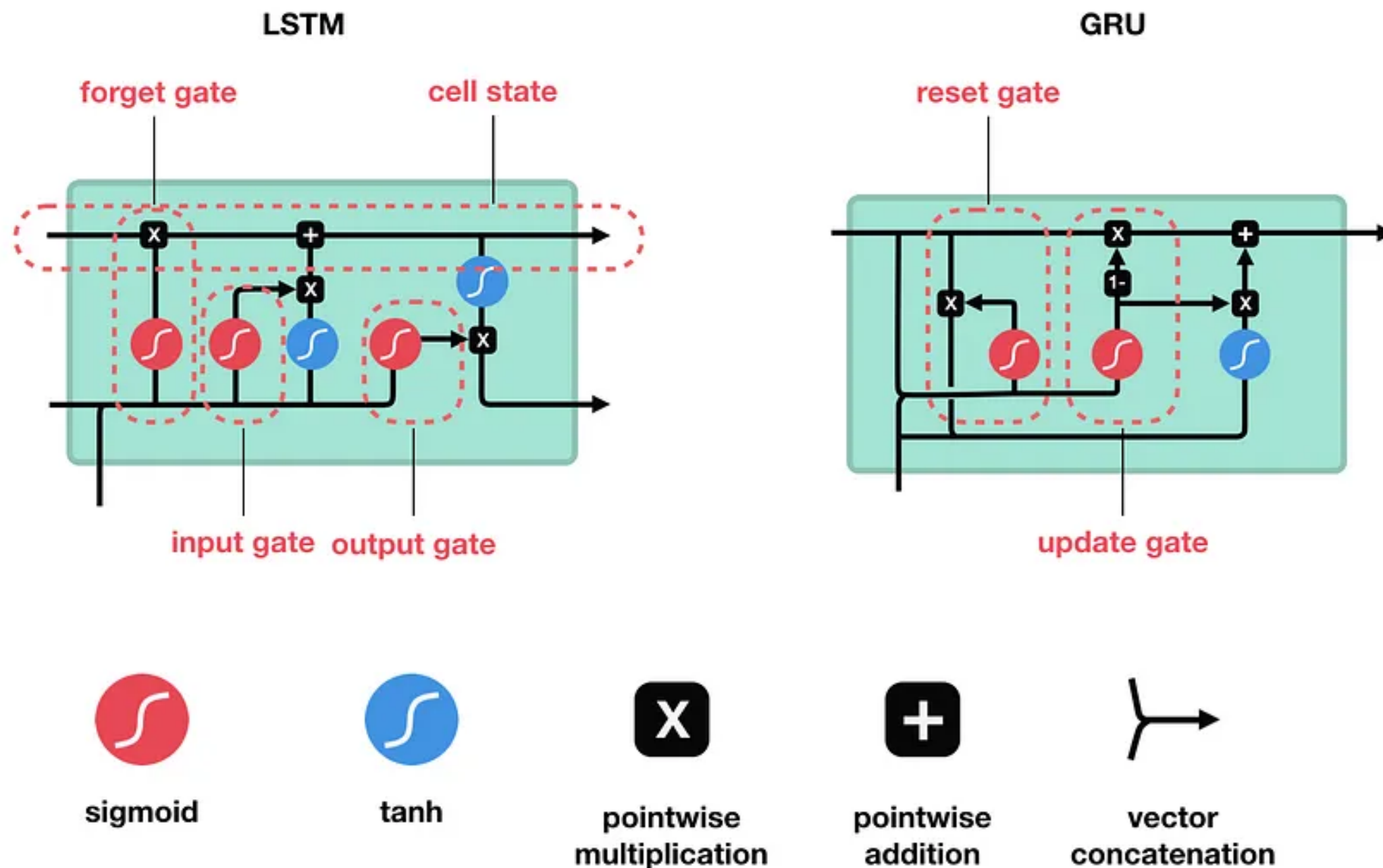
Word embedding

Feature Extraction

# MODELING

## Neural Network

### Long Short-term Memory (LSTM) and Gated Recurrent Unit (GRU)



Newer generation of Recurrent Neural networks.

The core concept of LSTM are the cell state, and its various gates.

LSTM cell replacement: GRU (Similar to an LSTM cell).

GRU's got rid of the cell state and used the hidden state to transfer information.

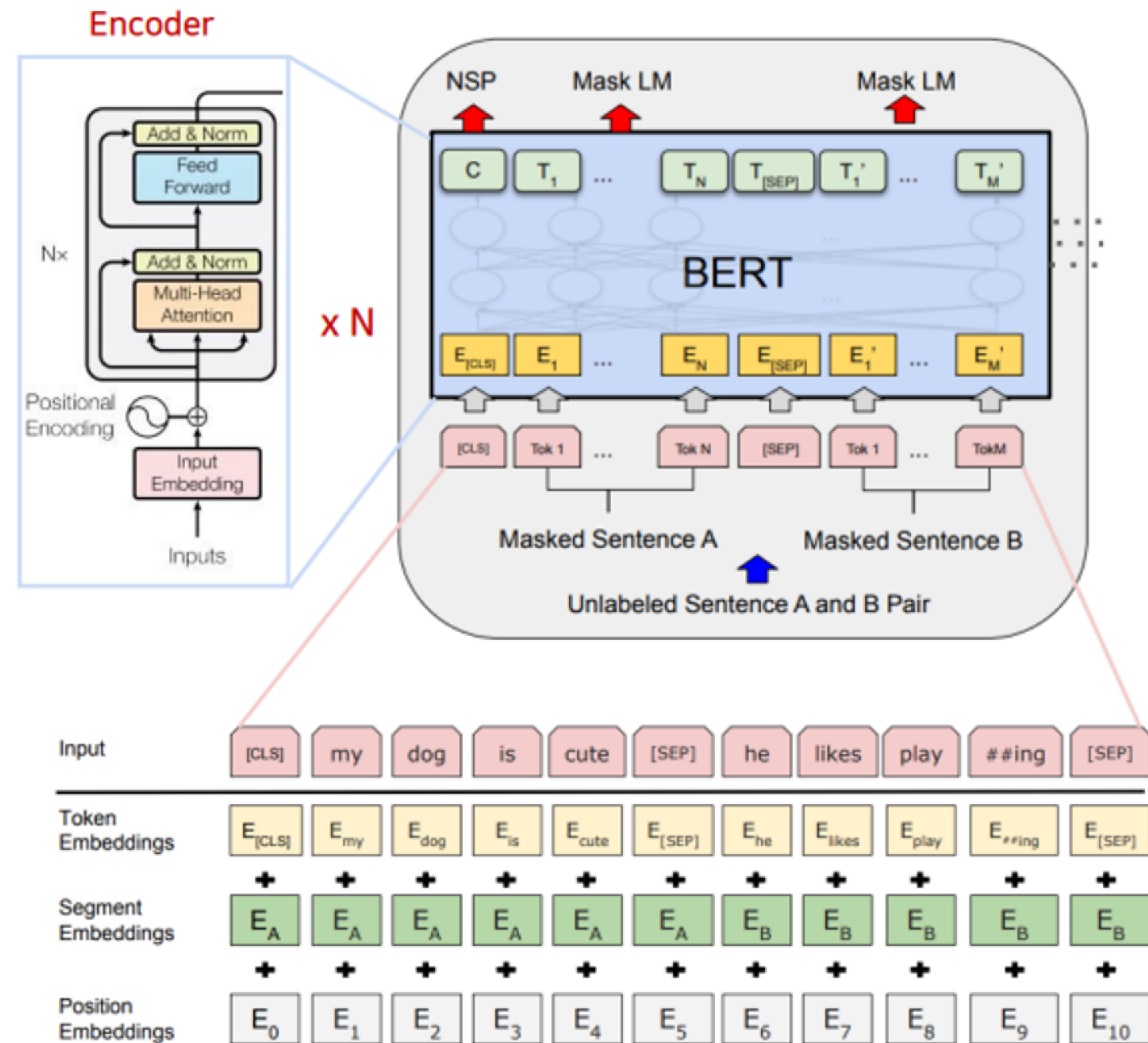
# MODELING

## phoBERT Archiecture

- 1 Build the dataset containing **3 kinds of sentiment**.
- 2 Fine-tuning PhoBERT with the dataset above.
- 3 Use the model for evaluation.



Back to Agenda Page





# MODELING

## XLM - Roberta

1

3 pre-trained language modelling objectives are used: CLM, MLM, TLM.

2

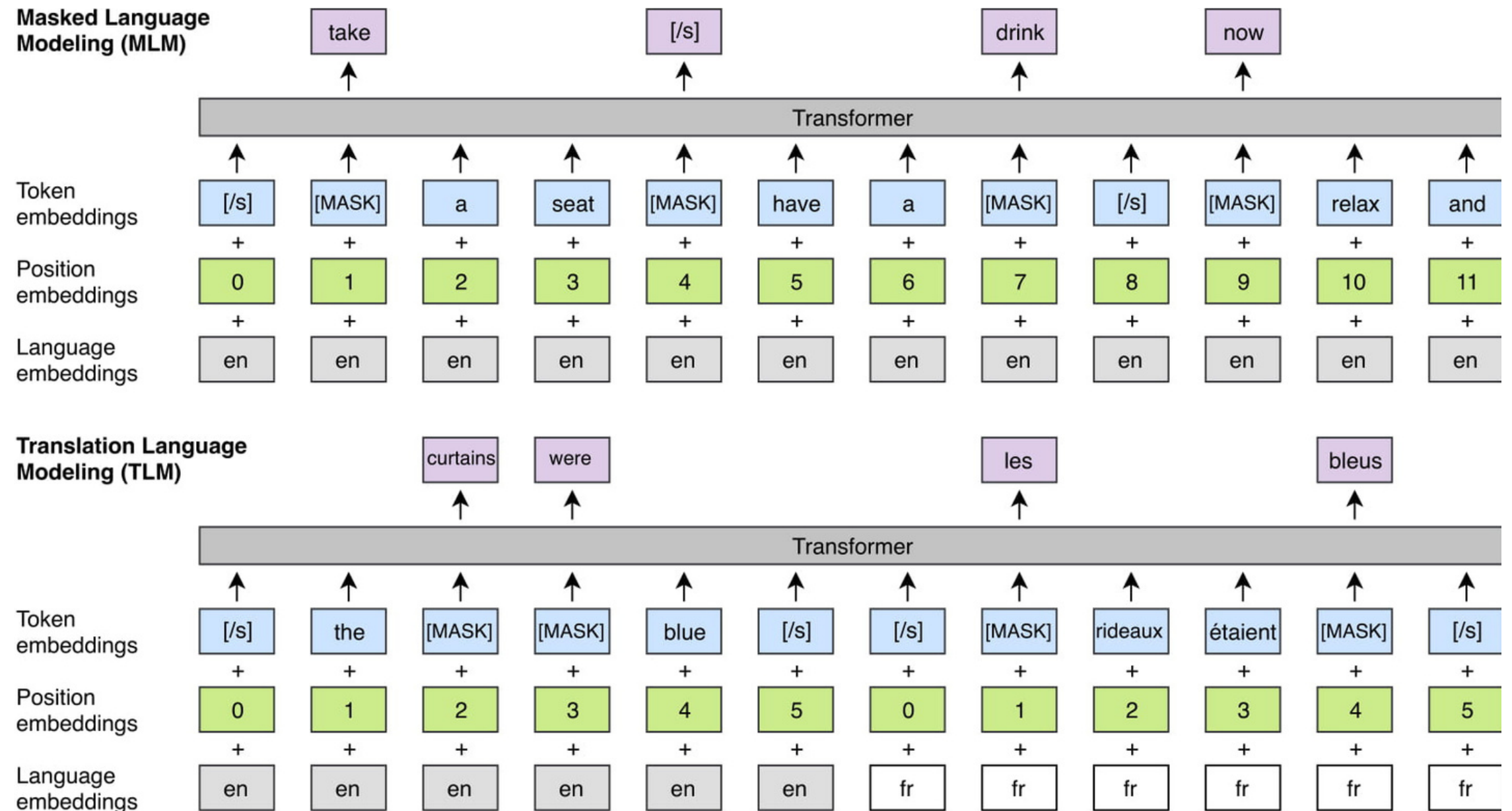
The model is pre-trained over RoBERTa architecture.

3

Use the model for evaluation.



[Back to Agenda Page](#)



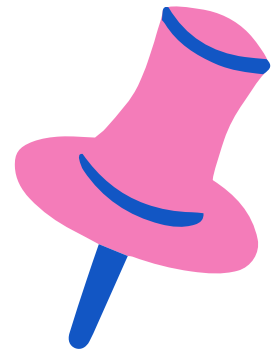
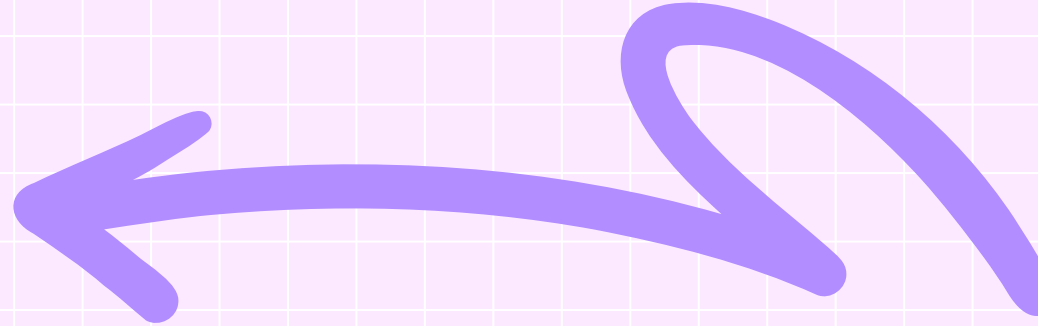
# EVALUATIONS

Model	Train Accuracy (%)	Test Accuracy (%)	Inference Time	Input Fields
PhoBERT	99.86	81.24	7.43h	Text Fields Only
RF + TF - IDF	99.61	68.76	3.2s	Text Fields Only
DBLSTM	99.60	62.28	300.1s	Text Fields Only
DBGRU	99.80	63.67	412.4s	Text Fields Only
XLM - RoBERTa	<b>99.92</b>	<b>82.04</b>	3.30h	Text Fields Only

The classification results for each model

[Back to Agenda Page](#)


# DIFFICULTIES AND LIMITATIONS



## Difficulties

Do not have common knowledge about the problem.

Do not find the effective way to preprocess the data.



## Limitations

Limited hardware resources.

Limited knowledge about the theoretical properties of the models



# CONCLUSIONS

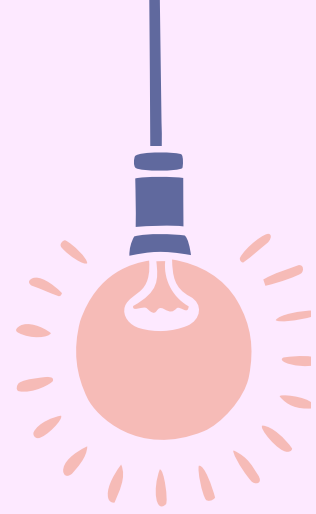
## SUMMARY

- To gain success in solving a problem, having great insight about the problem is crucial.
- By using various methods and tools, we have trained five different models on our dataset.
- **The best model among all is XLM - RoBERTa.**

## POSSIBLE EXTENSIONS

- Understanding the dataset and collect more data for training accuracy, finding better training strategy for each models.
- Research and implement new models like ViDeBERTa.
- Try out the application of self-supervised and semi-supervised learning methods.

[Back to Agenda Page](#)



THANKS FOR  
WATCHING

