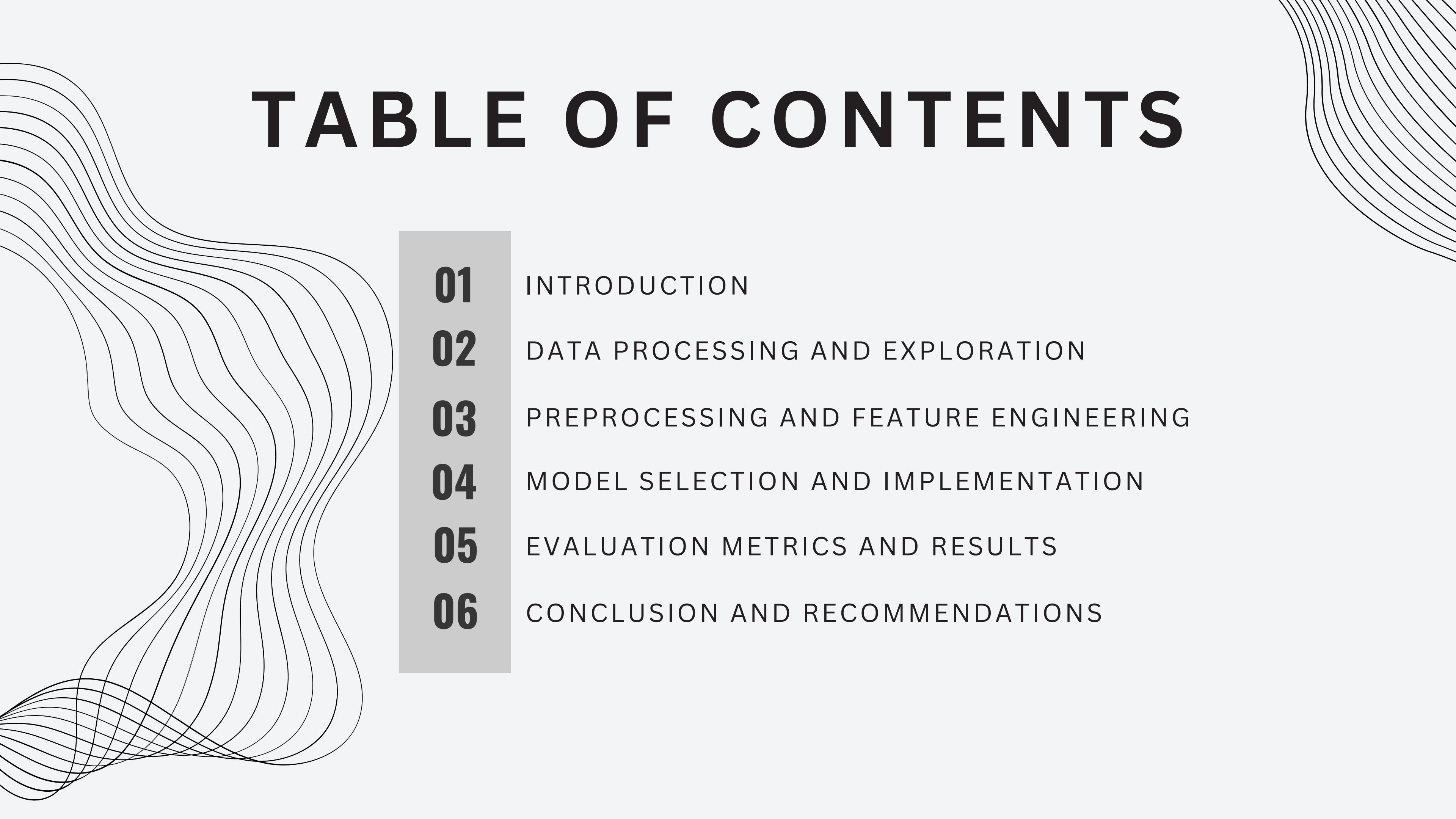


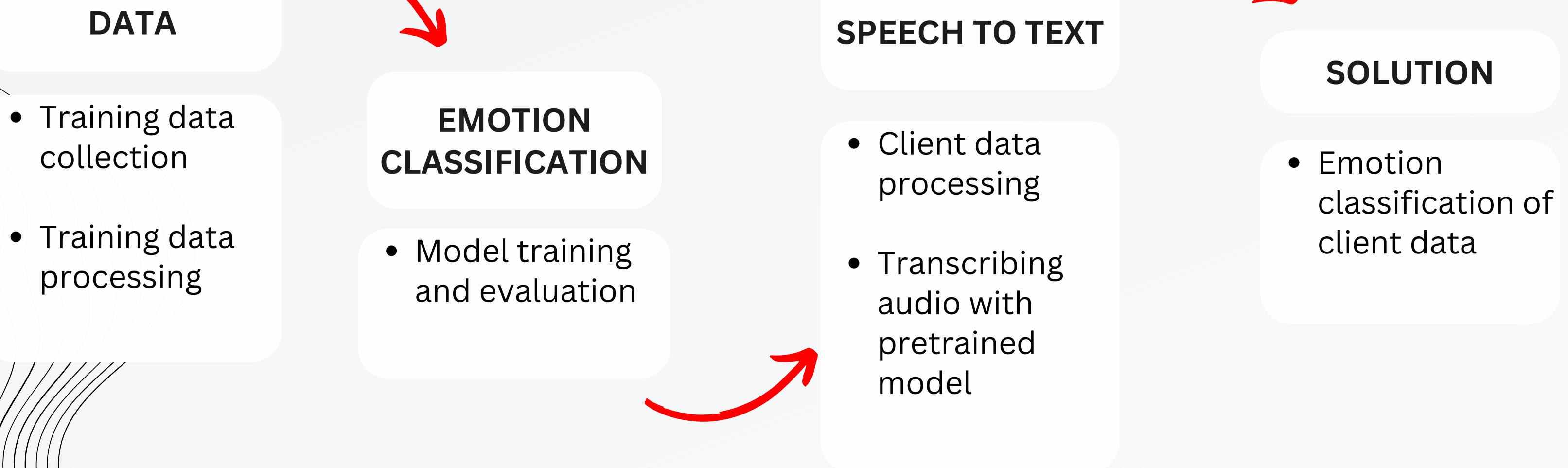
# **NATURAL LANGUAGE PROCESSING**

KORNELIA FLIZIK  
FEDOR CHURSIN  
MATEY NEDYALKOV  
PANNA BLANKA PFANDLER

# TABLE OF CONTENTS

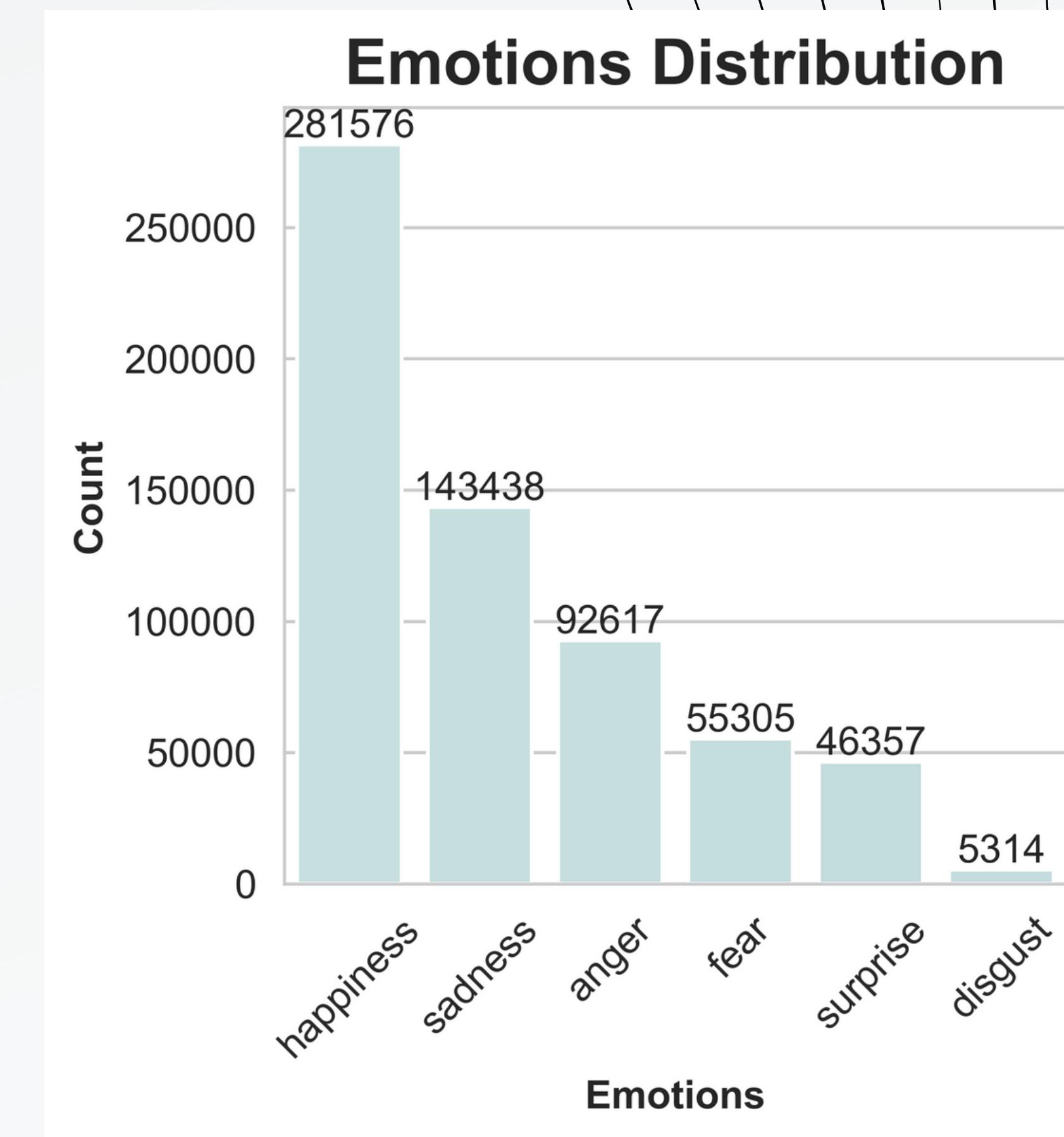
- 
- 01** INTRODUCTION
  - 02** DATA PROCESSING AND EXPLORATION
  - 03** PREPROCESSING AND FEATURE ENGINEERING
  - 04** MODEL SELECTION AND IMPLEMENTATION
  - 05** EVALUATION METRICS AND RESULTS
  - 06** CONCLUSION AND RECOMMENDATIONS

# PROJECT PIPELINE



# Data collection

Dataset name	Source
GoEmotions	Reddit comments
Smile Twitter Emotion	Twitter mentions
Friends Emotion-Labeled Dialogues	"Friends" TV show utterances
MELD dataset	"Friends" TV show utterances from dialogues
Carer dataset	English tweets



# EDA

- Annotation
- Lemmatization
- Merging all the data
- Tokenization
- Label Encoding

Emotion	Label
happiness	0
anger	1
disgust	2
sadness	3
surprise	4
fear	5

# MODEL SELECTION

MACHINE LEARNING

NAIVE BAYES

LOGISTIC REGRESSION

SUPPORT VECTOR MACHINE

NEUTRAL NETWORKS

MULTI-LAYER PERCEPTRON

LSTM

TRANSFORMERS

BERT

ALBERT

ROBERTA

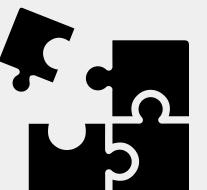
# ROBERTA MODEL



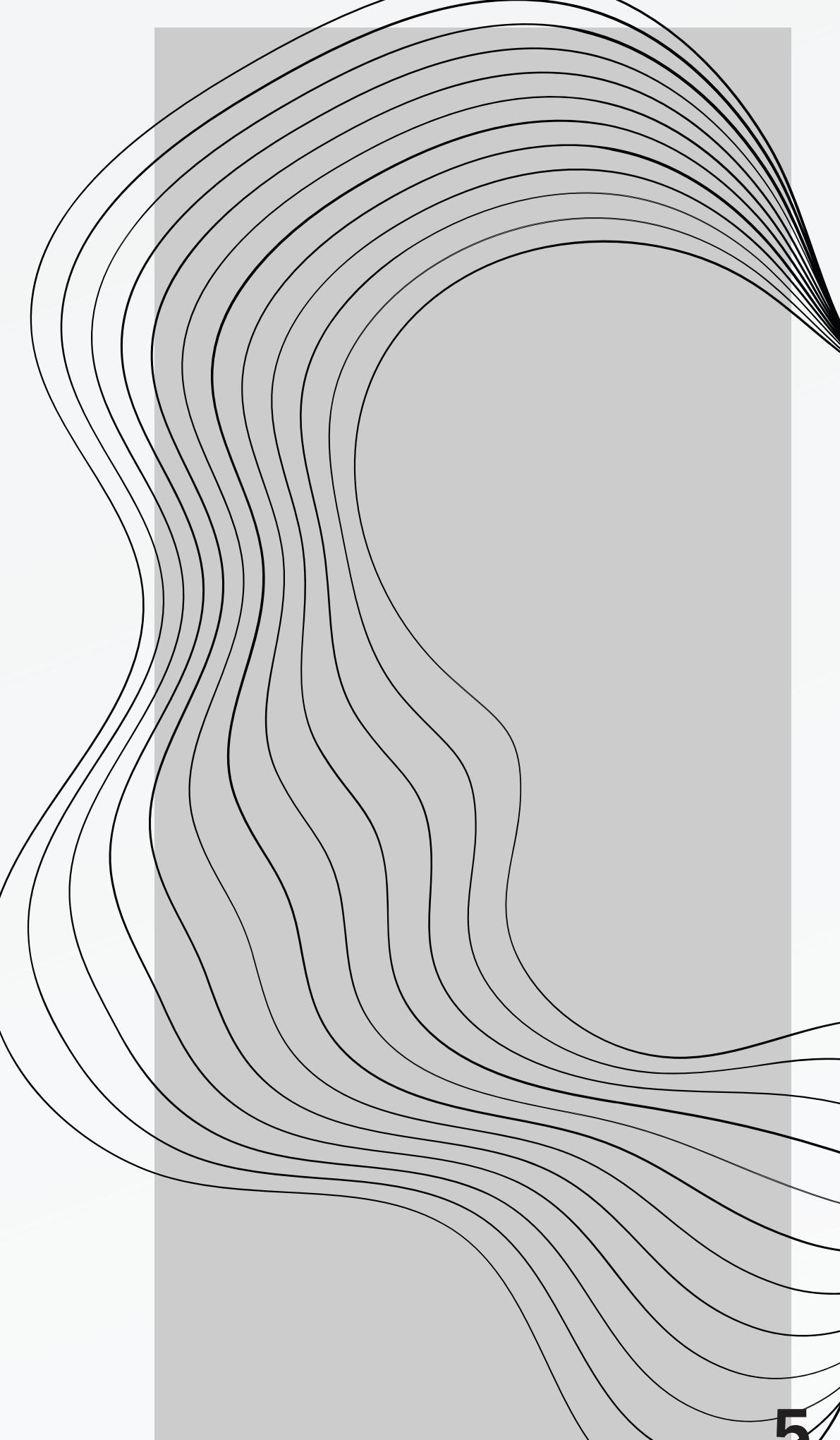
Transformer model with 24 layers



Pretained on a combination of five massive datasets resulting in a total of 160 GB of text data.



Employs a dynamic masking strategy during training for more reliable and adaptable word representations.



# EMOTION CLASSIFICATION MODEL

## Input

X: tokens using  
Roberta Tokenizer

Y: one of 6 classes

## Compiling

Sparse categorical  
cross-entropy

Adam with learning  
rate: 5e-5

## Training

Training time: 4 epochs

Batch size of 128

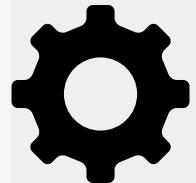
Early stopping  
monitoring validation  
loss



# WHISPER MODEL



Created by OpenAI and first released as open-source software in September 2022.



Transformer model that processes audio chunks into log-Mel spectrograms.



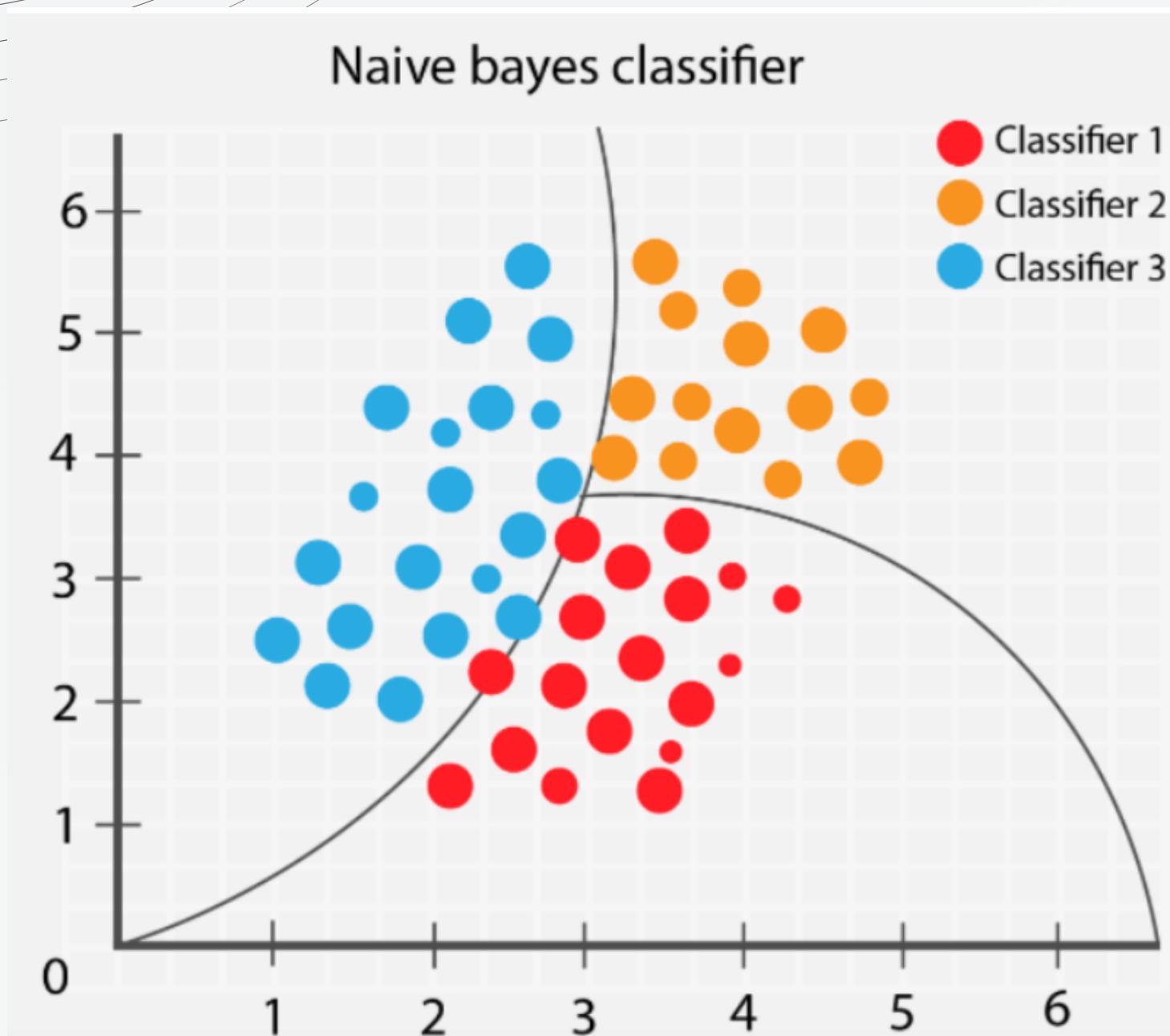
It is capable of transcribing speech in several languages, and translating speech into English.

# TRAINING RESULTS

Model Type	Loss	Accuracy	F1 Score
Naive Bayes	0.62	80%	0.565
Logistic Regression	0.8	65%	0.521
Multilayer Perceptron	1.09	76%	0.459
LSTM	1.07	63%	0.573
roBERTa	0.34	87%	0.733

# TRAINING RESULTS - NAIVE BAYES

9



Model Type	Loss	Accuracy	F1 Score
Naive Bayes	0.62	80%	0.565

Strengths:

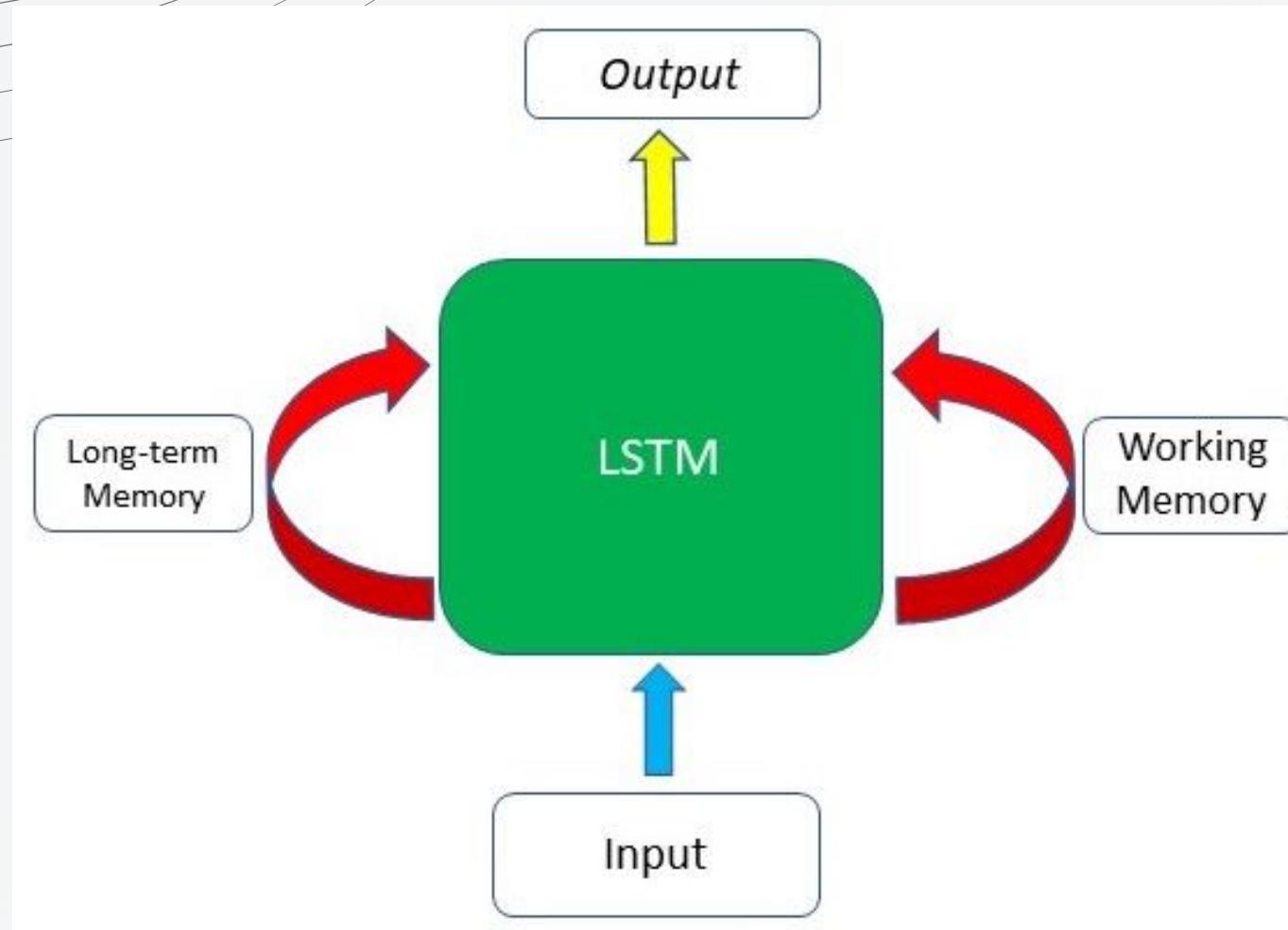
- Simplicity & Speed
- Baseline performance

Limitations:

- Feature independence
- Unseen features

# TRAINING RESULTS - LSTM

10



Model Type	Loss	Accuracy	F1 Score
LSTM	1.07	63%	0.573

Strengths:

- Temporal dynamics

- Memory

Limitations:

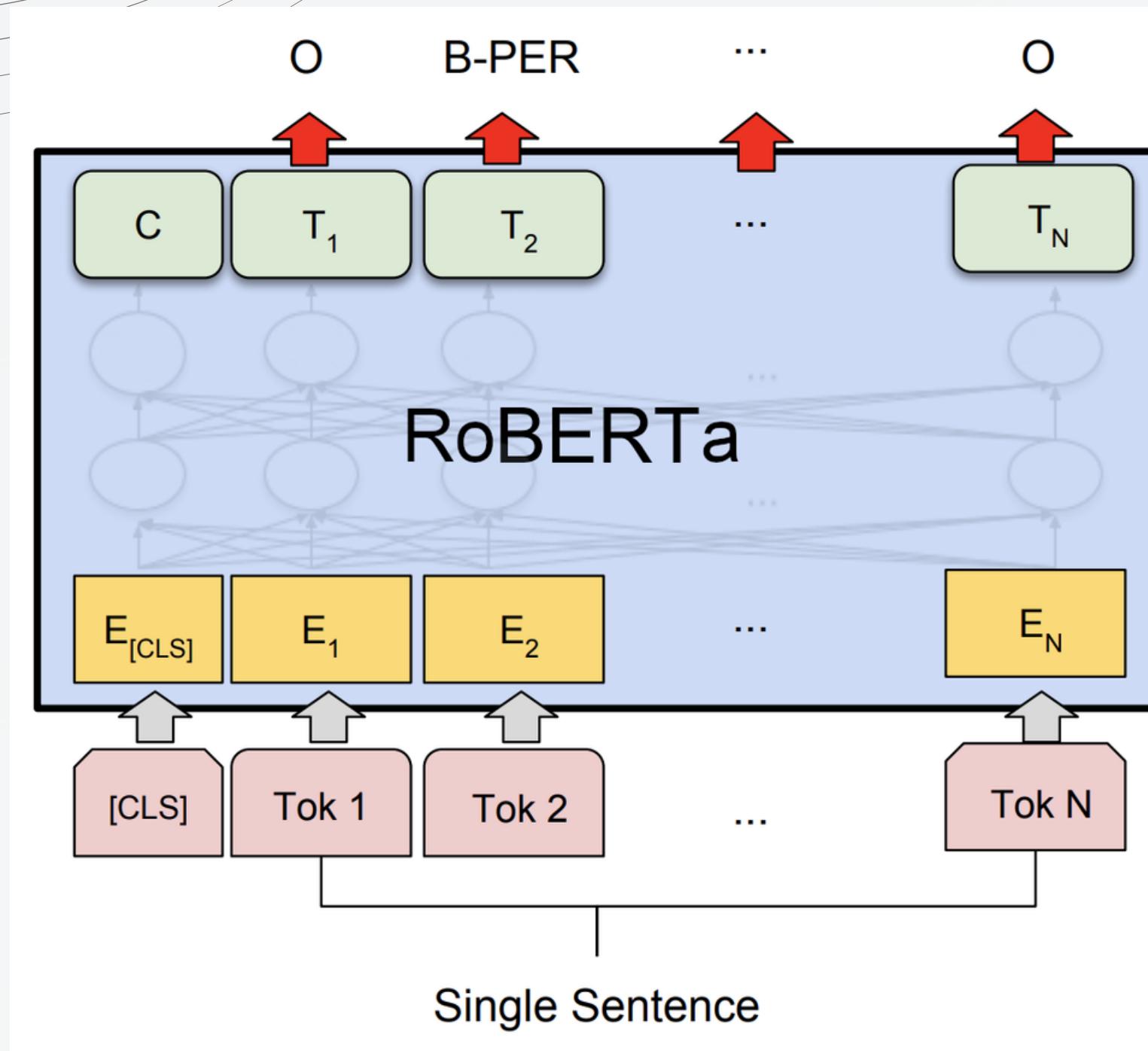
- Computational cost

- Tuning complexity

# TRAINING RESULTS -

## ROBERTA

11



Model Type	Loss	Accuracy	F1 Score
roBERTa	0.34	83%	0.733

Strengths:

- State-of-the-art performance
- Pre-trained knowledge

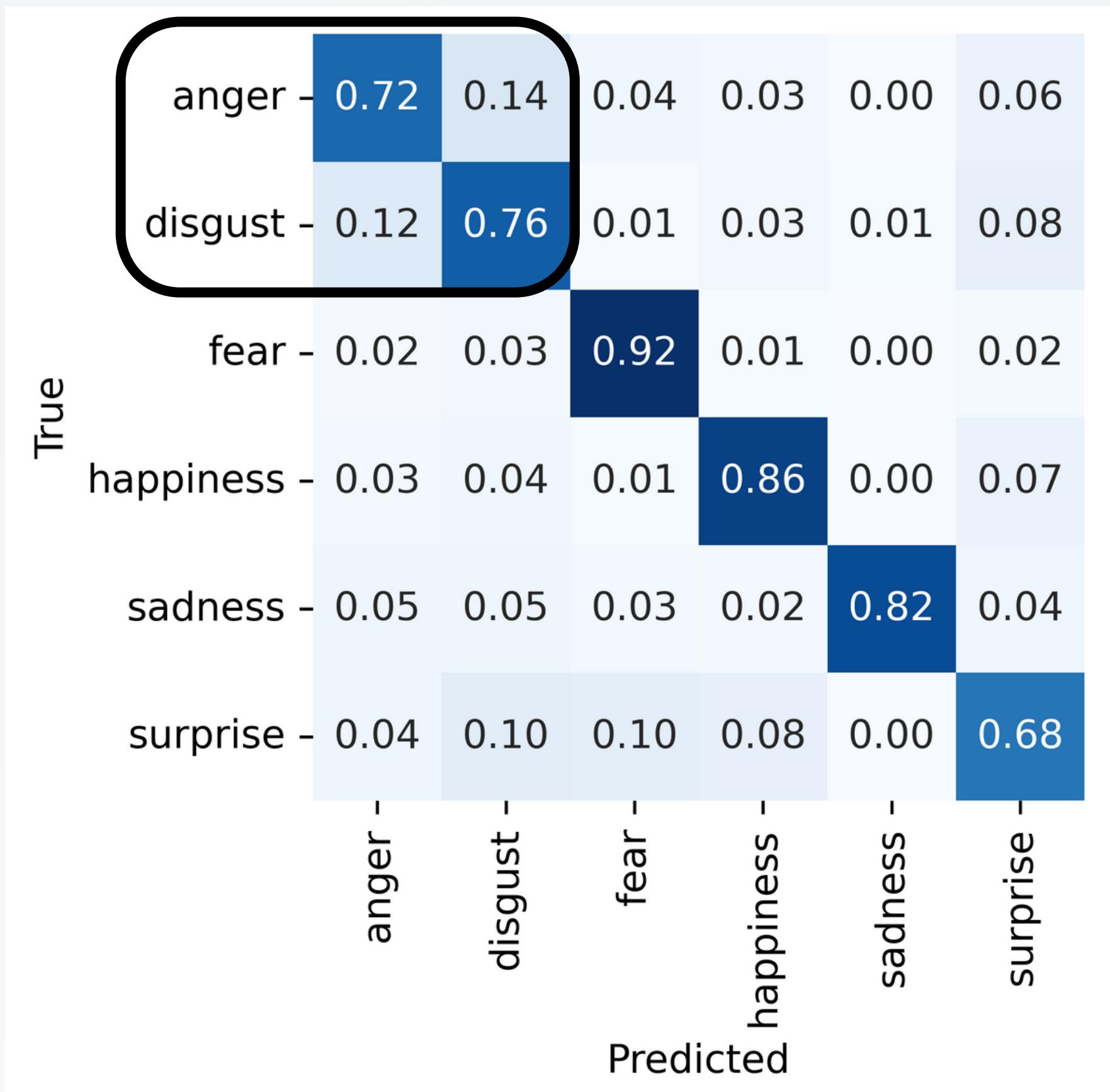
Limitations:

- Computational resources
- Black box nature

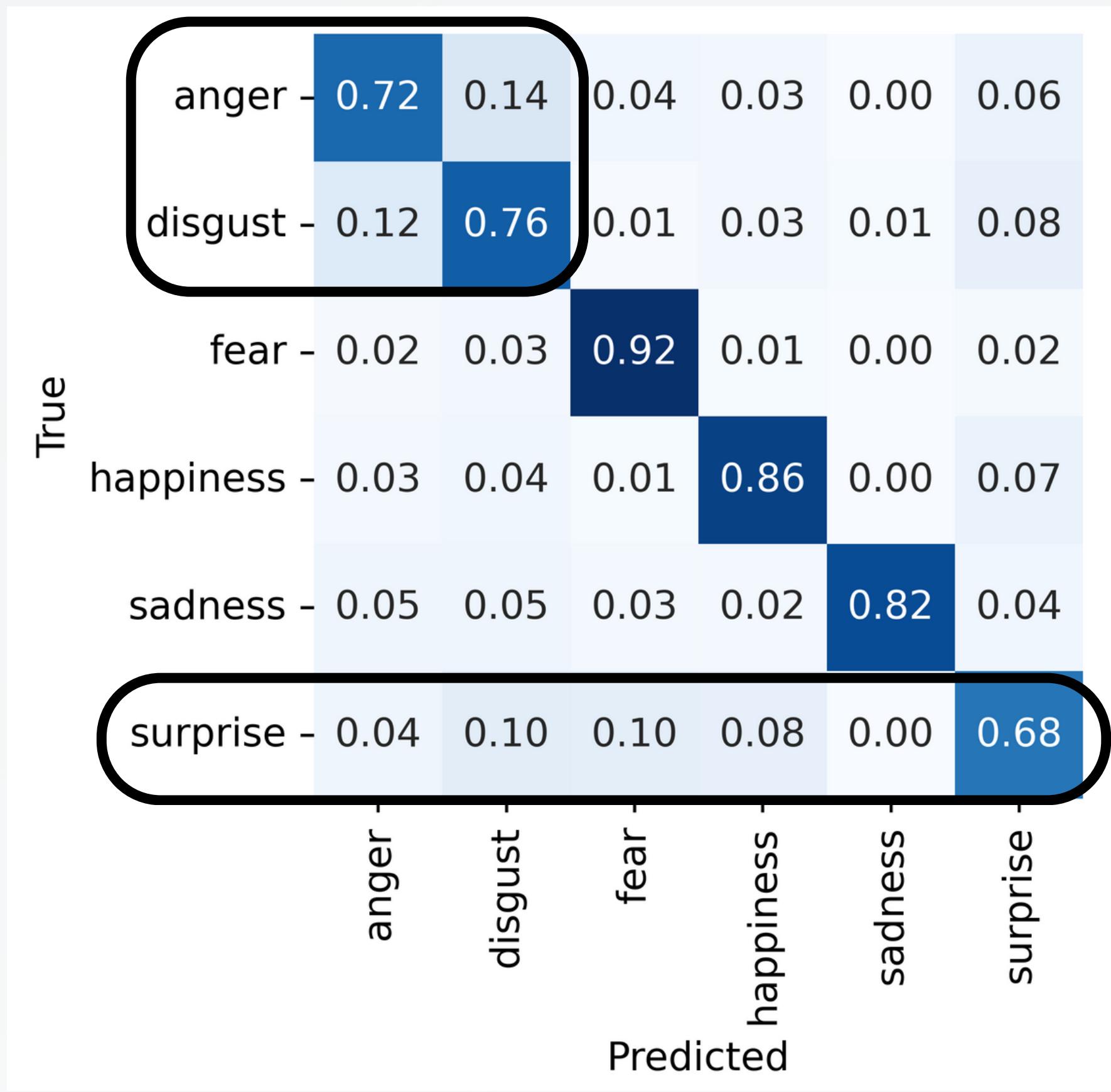
# Error Analysis

True	anger -	disgust -	fear -	happiness -	sadness -	surprise -
Predicted	anger -	disgust -	fear -	happiness -	sadness -	surprise -
anger -	0.72	0.14	0.04	0.03	0.00	0.06
disgust -	0.12	0.76	0.01	0.03	0.01	0.08
fear -	0.02	0.03	0.92	0.01	0.00	0.02
happiness -	0.03	0.04	0.01	0.86	0.00	0.07
sadness -	0.05	0.05	0.03	0.02	0.82	0.04
surprise -	0.04	0.10	0.10	0.08	0.00	0.68

# Error Analysis



# Error Analysis



# Error Analysis

Slang-heavy dataset



Hallucinations

Model Type	Loss	Accuracy	F1 Score
roBERTa	0.34	87%	0.733

# Error Analysis

Formal language dataset



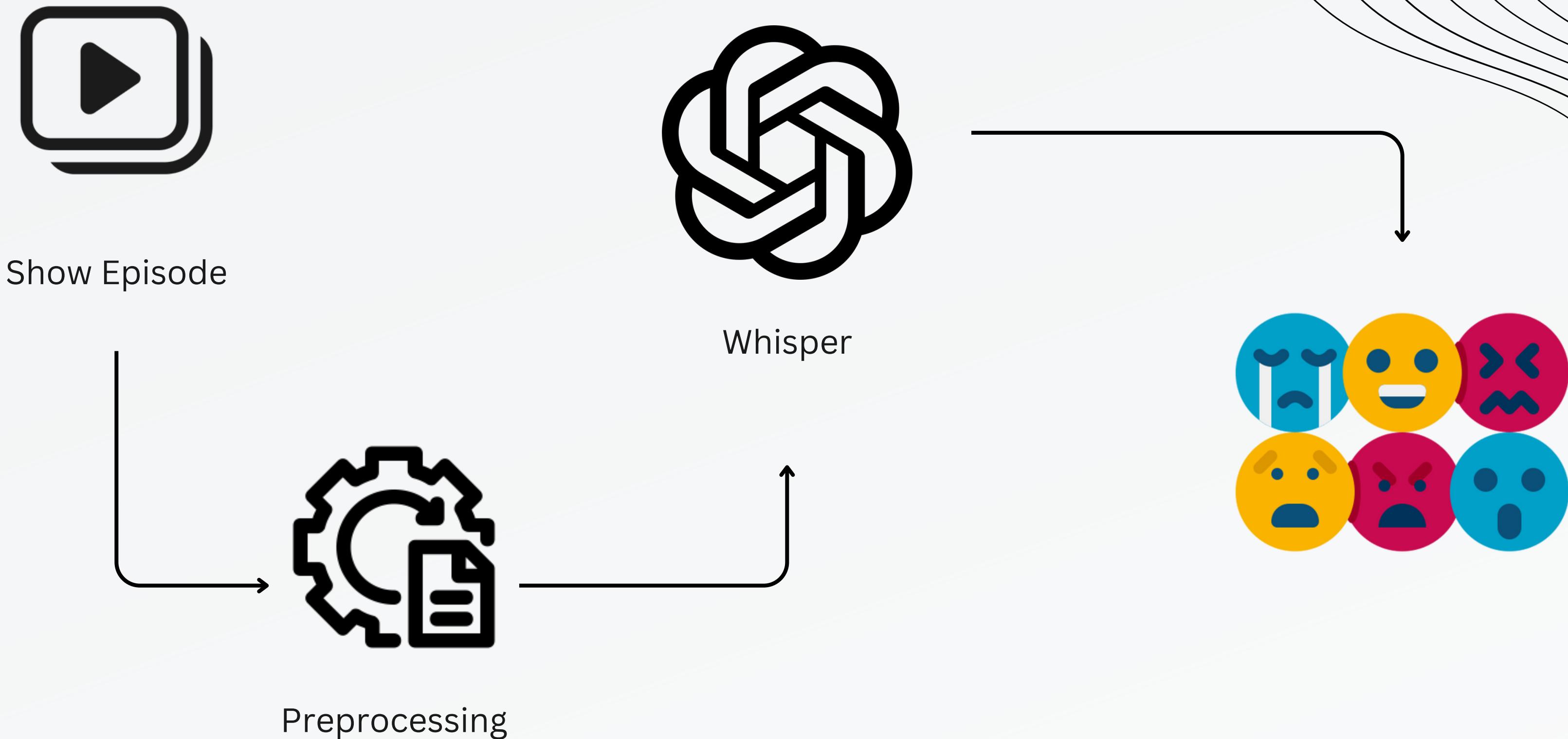
Efficient learning process

Model Type	Loss	Accuracy	F1 Score
roBERTa	0.07	96%	0.580

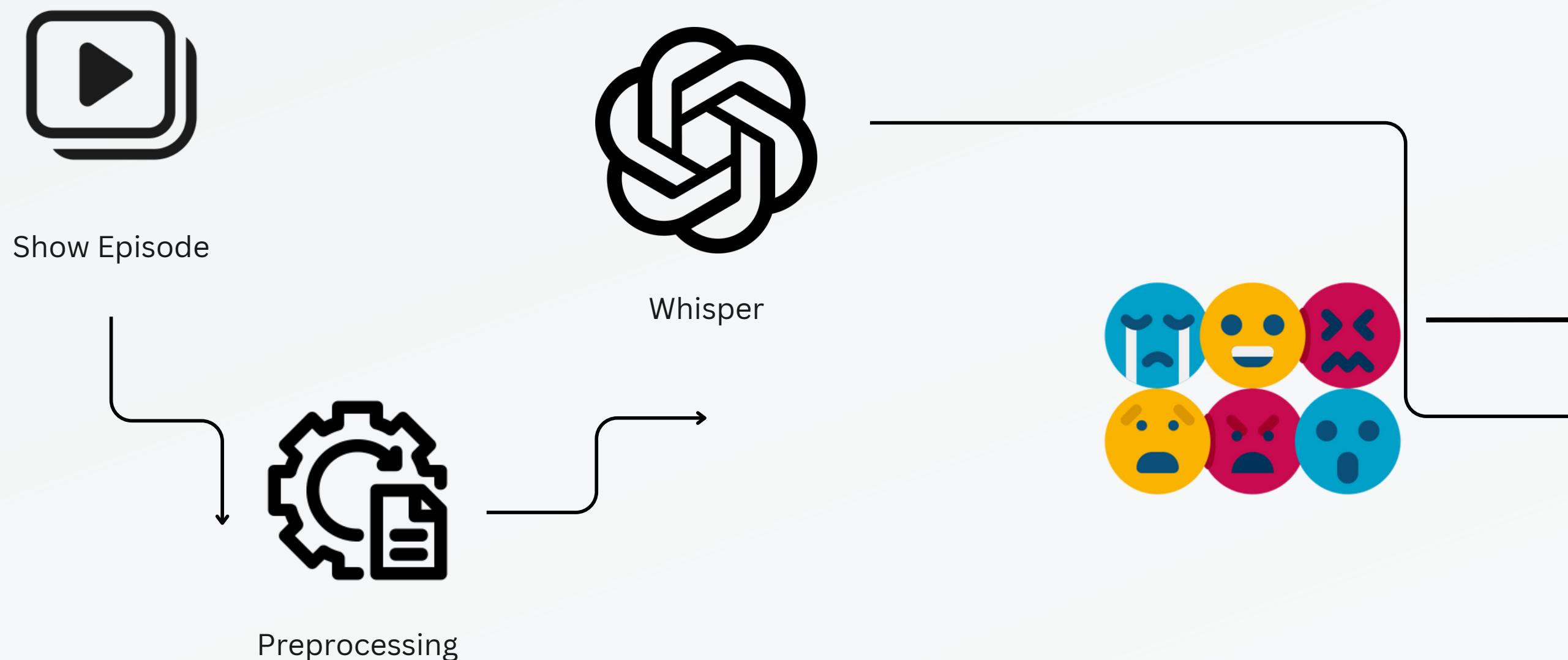
# Error Analysis

True	anger -	disgust -	fear -	happiness -	sadness -	surprise -
Predicted	anger -	disgust -	fear -	happiness -	sadness -	surprise -
anger -	0.91	0.00	0.04	0.01	0.04	0.00
disgust -	0.36	0.21	0.07	0.07	0.21	0.07
fear -	0.00	0.00	0.87	0.00	0.05	0.08
happiness -	0.00	0.00	0.00	0.99	0.00	0.01
sadness -	0.01	0.00	0.00	0.00	0.99	0.00
surprise -	0.02	0.00	0.00	0.01	0.01	0.96

# Application



# Application Limitations



# RECOMMENDATIONS



IMPROVED DATA  
LABELING



MODEL TUNING



DOMAIN-SPECIFIC  
ADAPTATION

# **THANK YOU FOR YOUR ATTENTION**

*Any questions?*

