



# *Attention-Based Deep Learning Models for Detecting Long-Term Effects Misinformation of COVID-19*

**Advisor : Che-Lun Hung**

**Chun-Ying Wu**

**Student : Jian-An Chen**

# Outline

- Introduction
- Methods
  - Data collection
  - Data preprocessing
  - Data analysis
  - Models
- Experiments
- Discussions & Conclusion

# Background

## COVID-19 Pandemic (2019-2022)

1. The novel coronavirus (SARS-CoV-2) was identified in December 2019.
2. Governments and health organizations worldwide implemented various measures to control the pandemic.
3. During the crisis, people sought reliable information on symptoms, prevention, and treatment.

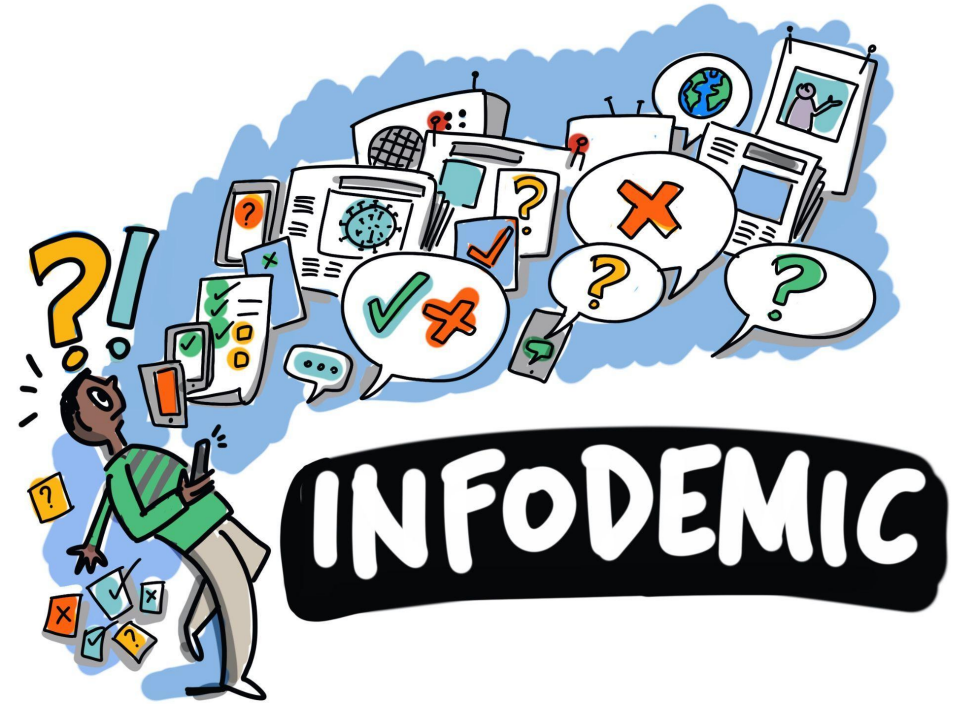


COVID-19 Dashboard, Johns Hopkins University

# Background

## Challenges of Infodemic

1. An infodemic<sup>1</sup> refers to an oversupply of information, including false or misleading content, during a disease outbreak.
2. It can cause confusion, risk-taking behaviors, and mistrust in health authorities.
3. Reliable health information detection systems play a crucial role.



1. [Infodemic \(who.int\)](https://www.who.int/infodemic)

# Challenges beyond pandemic control

A decorative graphic in the top right corner consisting of a network of interconnected nodes and lines, resembling a molecular or biological structure.

## Long-COVID

While the distribution of COVID-19 vaccines contributed to the gradual control of the pandemic, the virus persisted, giving rise to post-infection symptoms known as long COVID, confirmed in at least 10% of those who contracted the virus<sup>2</sup>.

## Reinfection

Instances of reinfection after initial recovery were observed, with research from the U.S. Department of Veterans Affairs indicating increased risks of mortality, hospitalization, and post-symptomatic conditions for reinfected patients<sup>3</sup>.

2. [Long COVID: major findings, mechanisms and recommendations.](#)

3. [Acute and postacute sequelae associated with SARS-CoV-2 reinfection | Nature Medicine](#)

# Related work

Title	Published date	Methods / Techniques
Fighting an Infodemic: COVID-19 Fake News Dataset <sup>4</sup>	Nov, 2021 CCIS,volume 1402	TF-IDF + ML Algorithms
A Heuristic-driven Ensemble Framework for COVID-19 Fake News Detection <sup>5</sup>	Jan, 2021 CONSTRAINT Workshop, AAI'21,	Pre-trained language models + Soft Voting Method
Cross-SEAN: A cross-stitch semi-supervised neural attention model for COVID-19 fake news detection <sup>6</sup>	August, 2021 Applied Soft Computing Volume 107	Multiple features extracted from tweets for detection + RNN

4. [Fighting an Infodemic: COVID-19 Fake News Dataset](#)
5. [A Heuristic-driven Ensemble Framework for COVID-19 Fake News Detection](#)
6. [Cross-SEAN: A cross-stitch semi-supervised neural attention model for COVID-19 fake news detection](#)

# Objectives

01

**Detect long COVID misinformation with deep learning approaches.**

02

**Enhance the performance by using different ensemble methods.**

03

**Compare performance with SOTA LLM-based embedding models.**

# Methods

A decorative network diagram in the top right corner, consisting of several circular nodes connected by thin lines, forming a web-like structure.

- Introduction
- **Methods**
  - **Data collection**
  - **Data preprocessing**
  - **Data analysis**
  - **Models**
- Experiments
- Discussions & Conclusion



# Data collection

2 fact-check websites



2 government bodies



5 open-source datasets

PolitiFact

WHO

CTF

Snopes

CDC

Fighting an Infodemic

CoAID

FibVID

FaCOV



All data from various sources was filtered using keywords related to the long COVID and reinfection, ensuring relevance to the topic.

# Fact-check websites

## Database

	label	tweet
0	1	"COVID-19 is targeted to attack Caucasians and...
1	0	COVID-19 wasn't targeted to spare Jewish and C...
2	1	"A European study has found COVID vaccines cou...
3	0	Study on possible COVID-19 brain effects looke...
4	1	There have been no new COVID-19 variants since...
...	...	...
82	0	Masks for COVID-19 are effective, as a six-par...
83	1	"New autopsy reports suggest Jeffrey Epstein m...
84	0	There's no new autopsy report linking Jeffrey ...
85	1	"People Of Color May Be Immune To The Coronavi...
86	0	Melanin doesn't protect against coronavirus

## Article on PolitiFact

### Claim:

To receive assisted suicide in Germany, you must first be fully vaccinated against Covid-19.

### Rating:



**False**

[About this rating](#)

### Context

Germany's government has never mandated such a regulation for assisted suicide, but a nongovernmental assisted suicide organization in Germany — Verein Sterbehilfe — did mandate, in November 2021, that access to their facility be limited to those who received a COVID-19 vaccination or had recently recovered from the disease.

[Clarify statement](#)

# Dataset

Source	periods	samples	fake	genuine
CTF	up to 2021	1292	1130	162
Fighting an Infodemic	up to 2021	218	62	156
CoAID	up to 2020	70	0	70
FibVID	up to 2020	615	318	297
FaCOV	up to 2021	811	811	0
PolitiFact	up to 2023	87	42	45
Snopes	up to 2023	15	9	6
CDC+WHO	up to 2023	58	0	58
Total		3166	2372	794

# Preprocessing

## Steps

**1.Remove duplicate entries**

**2.Remove emojis**

**3.Remove URLs**

## Original

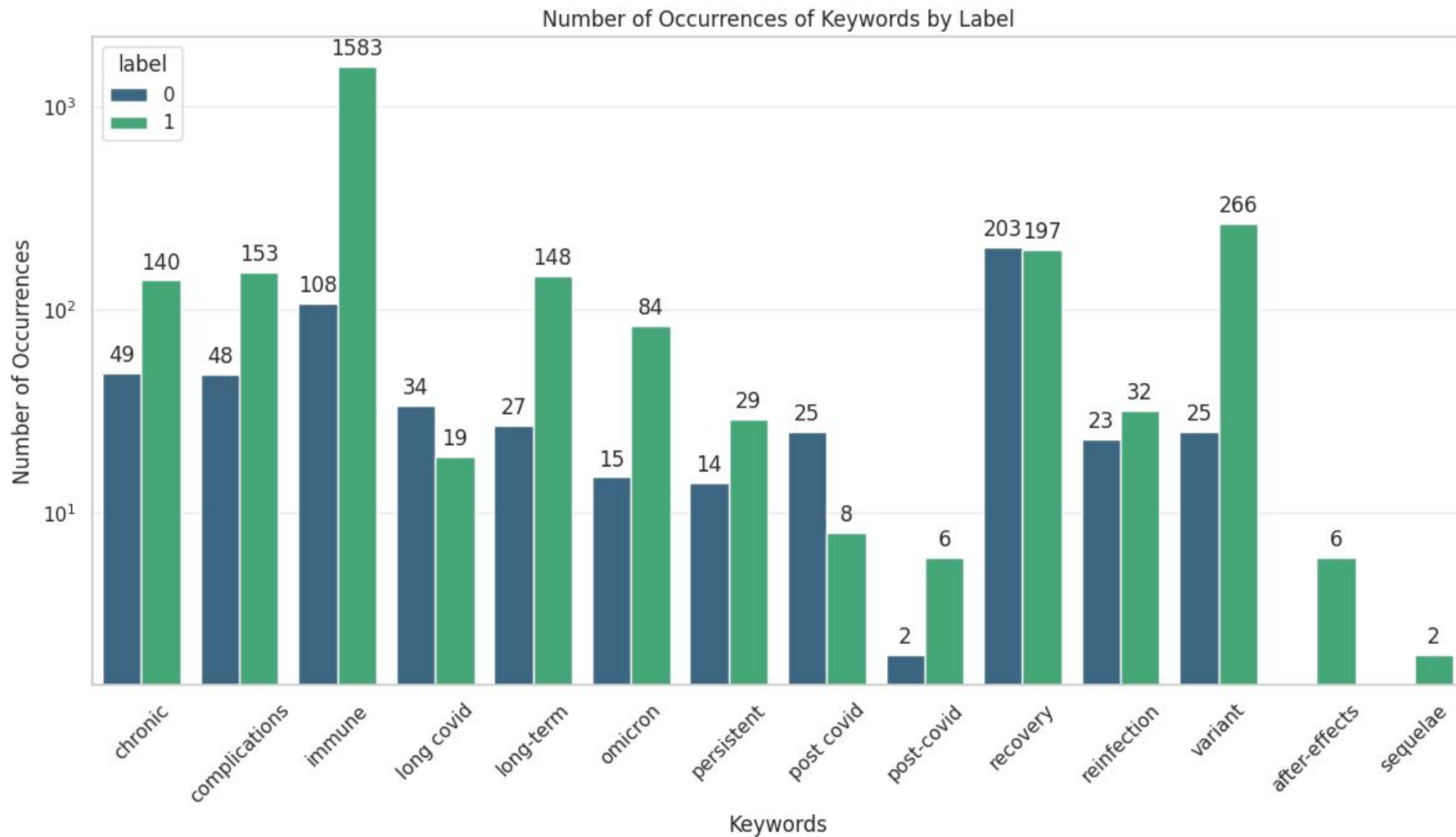
"📢#CoronaVirusUpdates:  
✅India's #COVID19  
recovery rate improves to  
76.98% as on September  
02 2020 📍 Steady  
improvement in India's  
COVID-19 recovery rate  
since #lockdown initiation  
on March 25 2020  
#IndiaFightsCorona  
@ICMRDELHI Via  
@MoHFW\_INDIA  
<https://t.co/I7aQjSrudh>"

preprocess

## Refined

"#CoronaVirusUpdates:  
India's #COVID19 recovery  
rate improves to 76.98% as  
on September 02 2020  
Steady improvement in  
India's COVID-19 recovery  
rate since #lockdown  
initiation on March 25 2020  
#IndiaFightsCorona  
@ICMRDELHI Via  
@MoHFW\_INDIA"

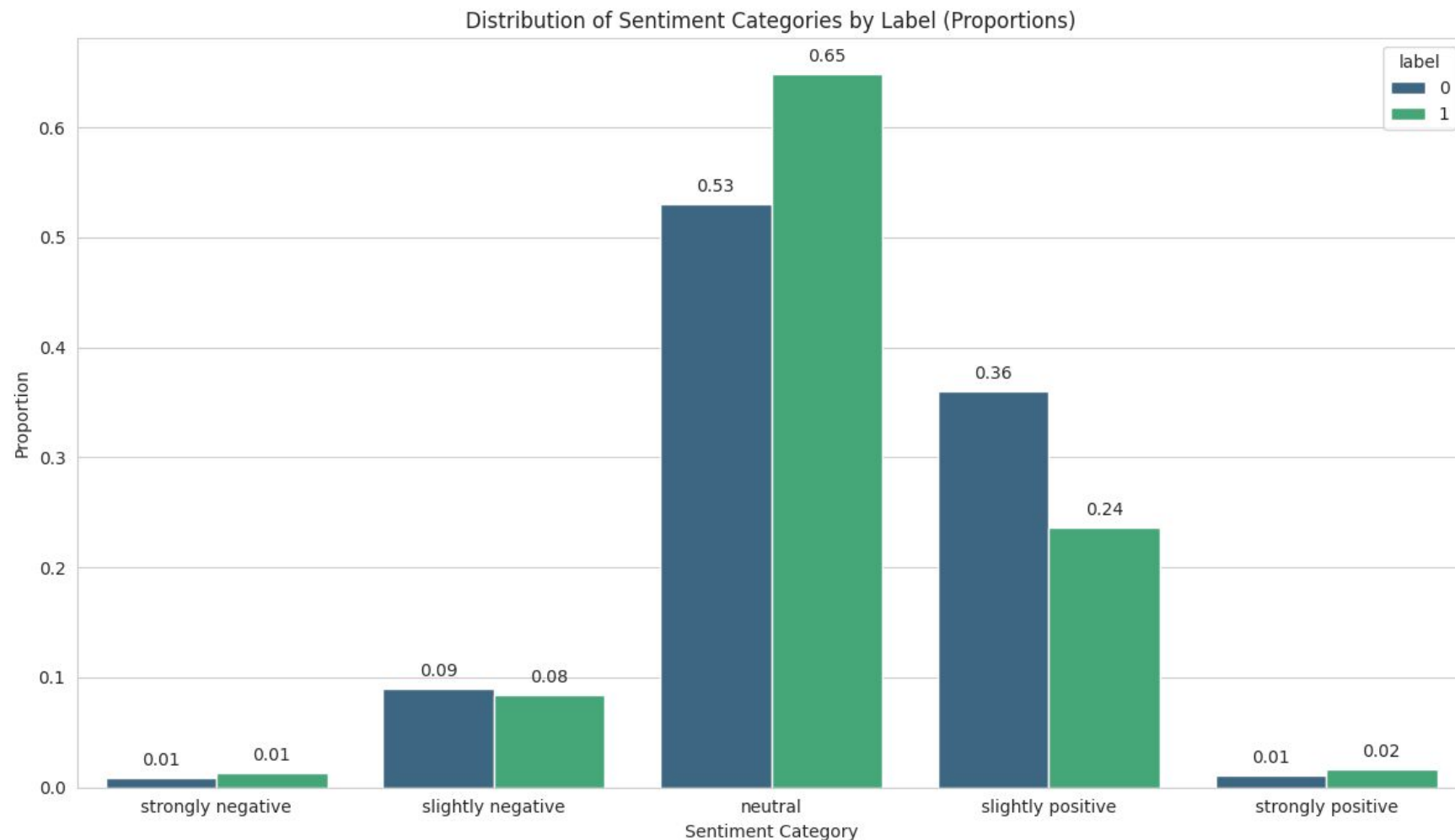
# Keywords occurrences



0 = genuine

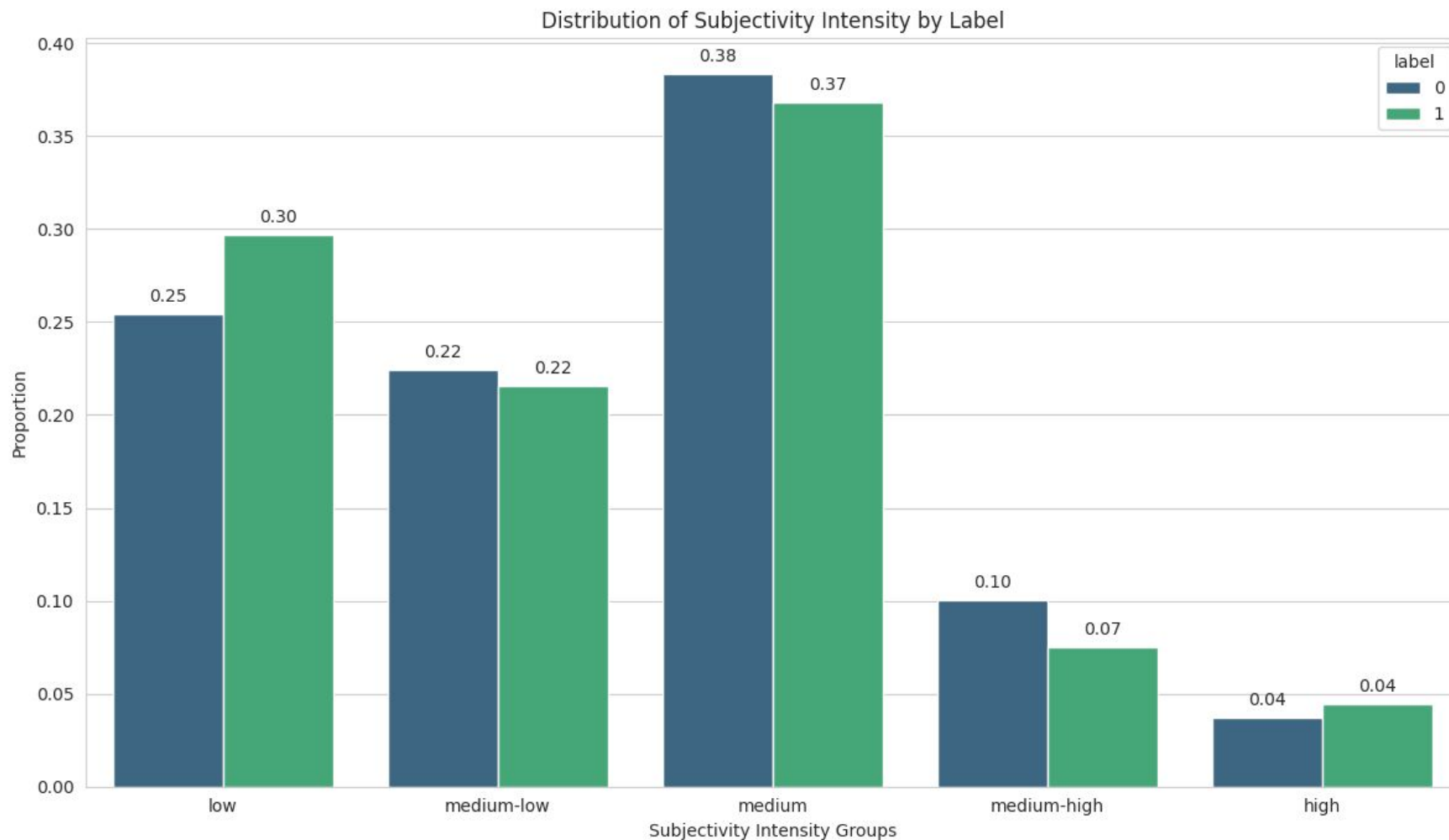
1 = fake

# Sentiment analysis (polarity)



	Polarity
Strongly negative	-1 to -0.5
Slightly negative	-0.5 to -0.1
Neutral	-0.1 to 0.1
Slightly Positive	0.1 to 0.5
Strongly Positive	0.5 to 1

# Sentiment analysis (subjectivity)



	Subjectivity
Low	0 to 0.2
Medium-Low	0.2 to 0.4
Medium	0.4 to 0.6
Medium-High	0.6 to 0.8
High	0.8 to 1

# Models

## Baseline

TF-IDF + SVM

Engineered features +  
XGBoost

## Deep models

HAN

BERT

RoBERTa

DeBERTa

XLNet

## LLM embedding

OpenAI's text-ada-002

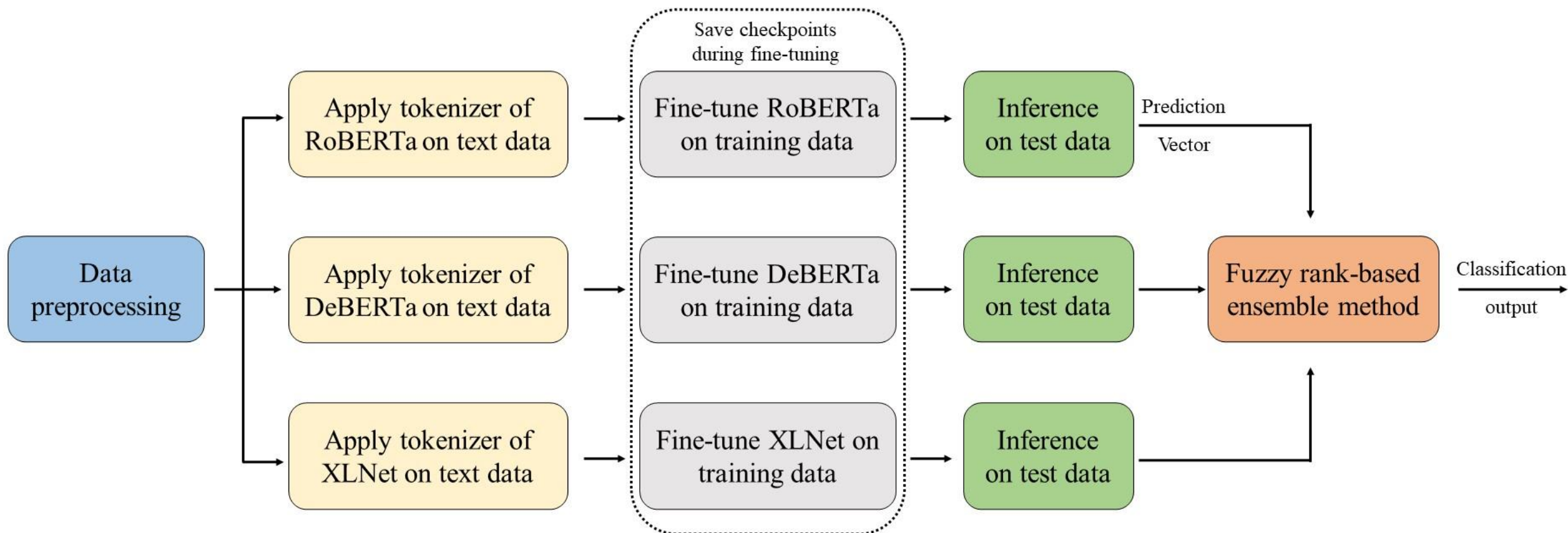
Google's Gemini

## LLM

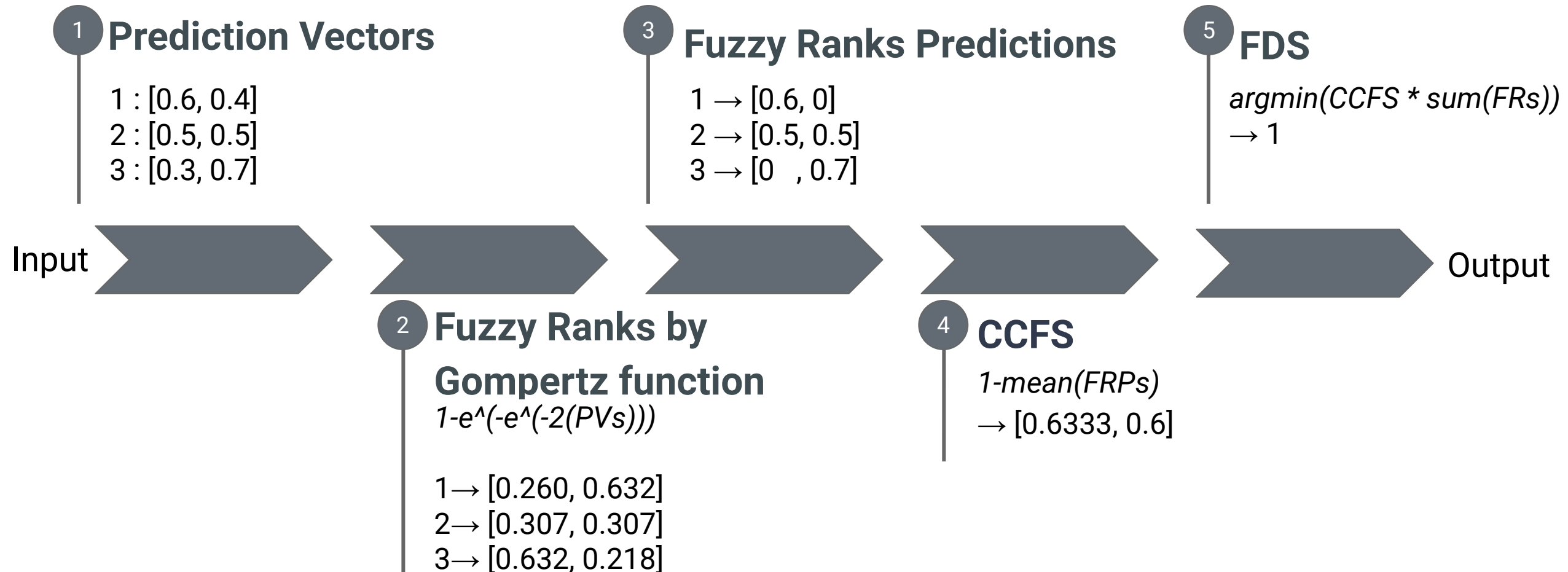
GPT - 4



# Ensembling



# Fuzzy rank method with Gompertz function



# Experiments

- Introduction
- Methods
  - Data collection
  - Data preprocessing
  - Data analysis
  - Models
- **Experiments**
- Discussions & Conclusion

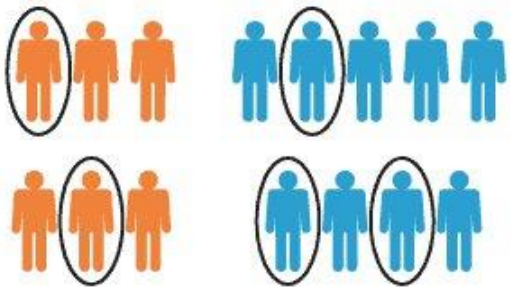
# Experiments

## Split data

Training set : 90%

Test set : 10%

## Stratified strategy



## Evaluation metrics

Accuracy

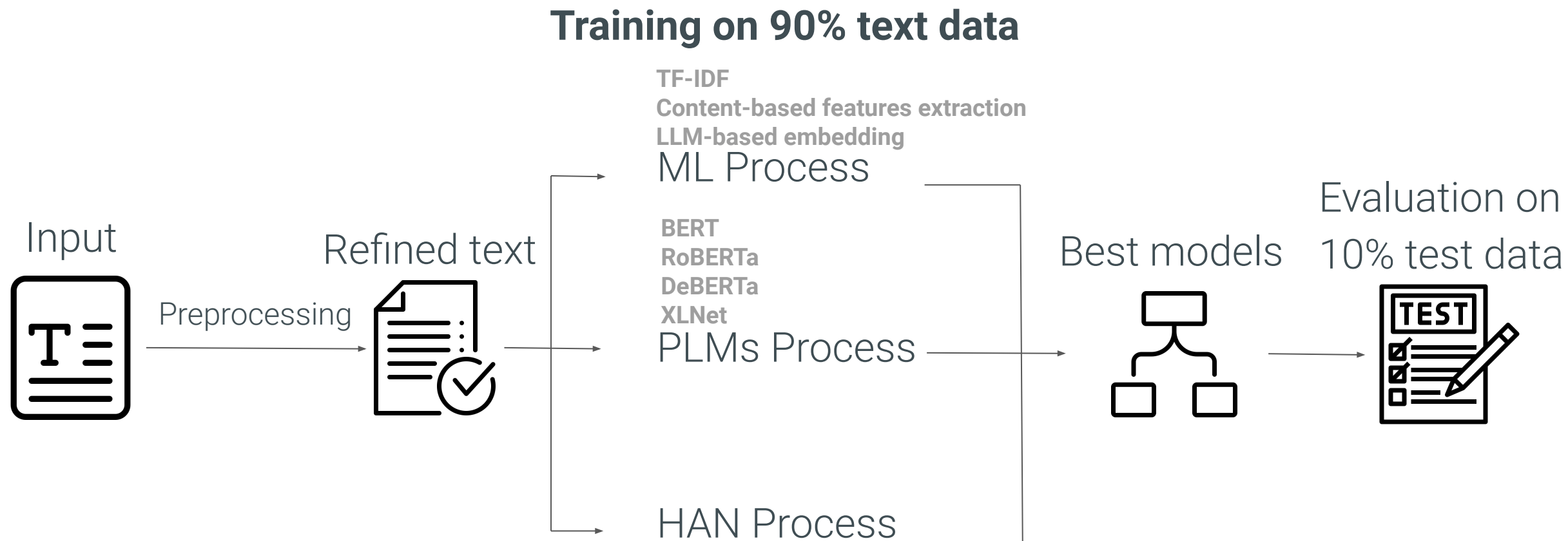
Precision

Recall

F1-score

AUC

# Workflow



# Implementation

Model	Configuration	Hardware
ML	5-fold Cross-Validation	Not specified
PLMs	AdamW optimizer, 20 epochs , lr=2e-5, batch_size = 64, max_length = 256	RTX A5000
HAN	Original settings, 20 epochs	Tesla T4

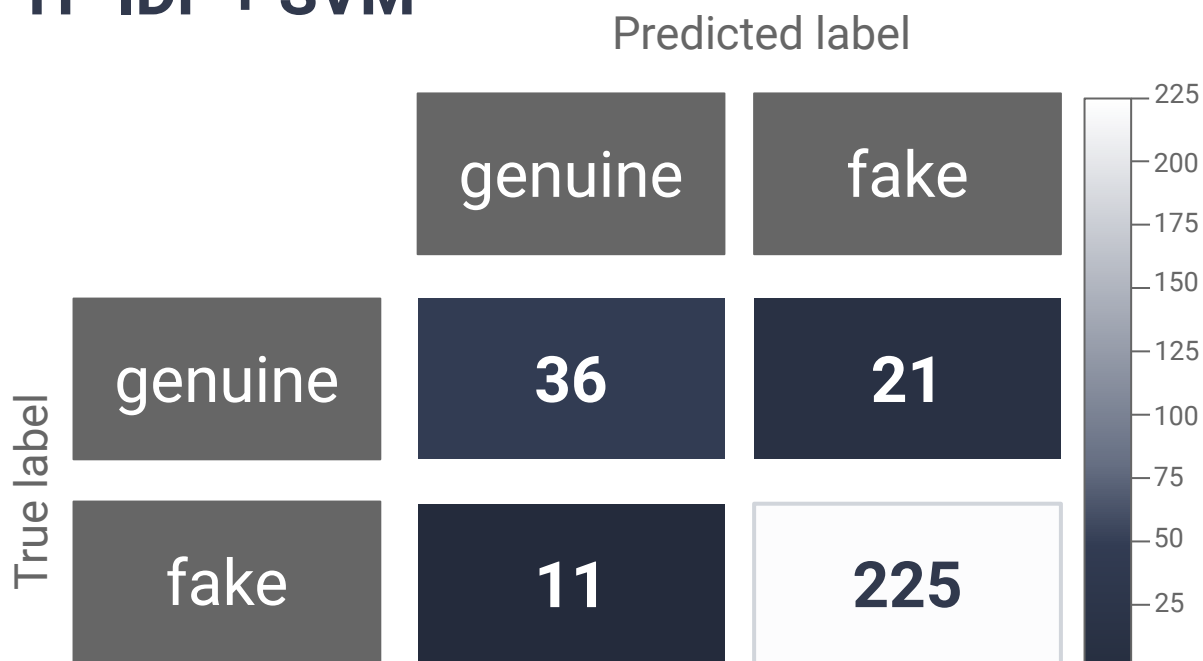
# Results

Model	Accuracy	Precision	Recall	F1-score	AUC
TF-IDF + SVM	89.08%	91.46%	95.34%	93.36%	92.02%
Engineered features + XGBoost	82.59%	83.88%	97.03%	89.98%	78.95%
HAN	90.78%	94.09%	94.49%	94.29%	93.09%
BERT	91.13%	93.75%	95.34%	94.54%	96.31%
RoBERTa	91.81%	94.54%	95.34%	94.94%	96.56%
DeBERTa	91.81%	94.17%	95.76%	94.96%	95.75%
XLNet	92.83%	94.24%	97.03%	95.62%	96.12%
GPT-4	82.25%	91.82%	85.59%	88.60%	N/A
Text-ada-002 + SVM	92.83%	93.52%	<b>97.88%</b>	95.65%	94.88%
Gemini + SVM	91.47%	92.71%	97.03%	94.82%	93.27%
Soft voting	93.17%	94.63%	97.03%	95.82%	97.14%
Soft voting with features	92.49%	93.15%	<b>97.88%</b>	95.45%	96.52%
Fuzzy ranks	<b>93.52%</b>	<b>94.65%</b>	97.46%	<b>96.03%</b>	<b>97.15%</b>

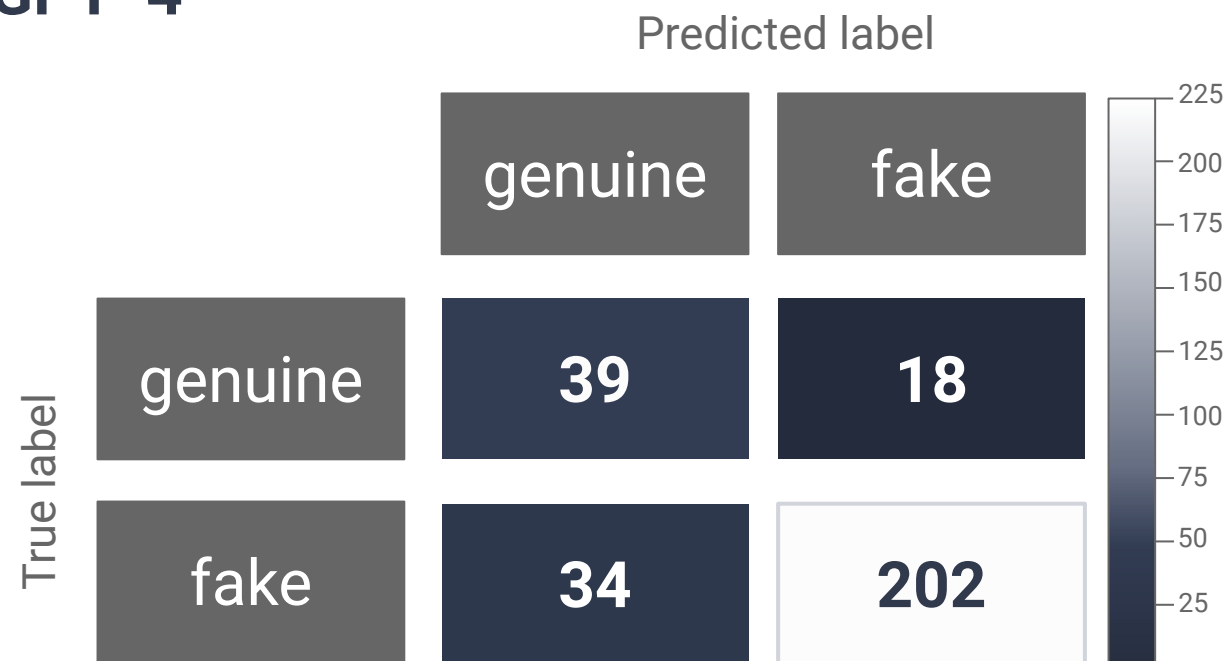
# Confusion Matrix (TF-IDF + SVM v.s. GPT-4)



**TF-IDF + SVM**



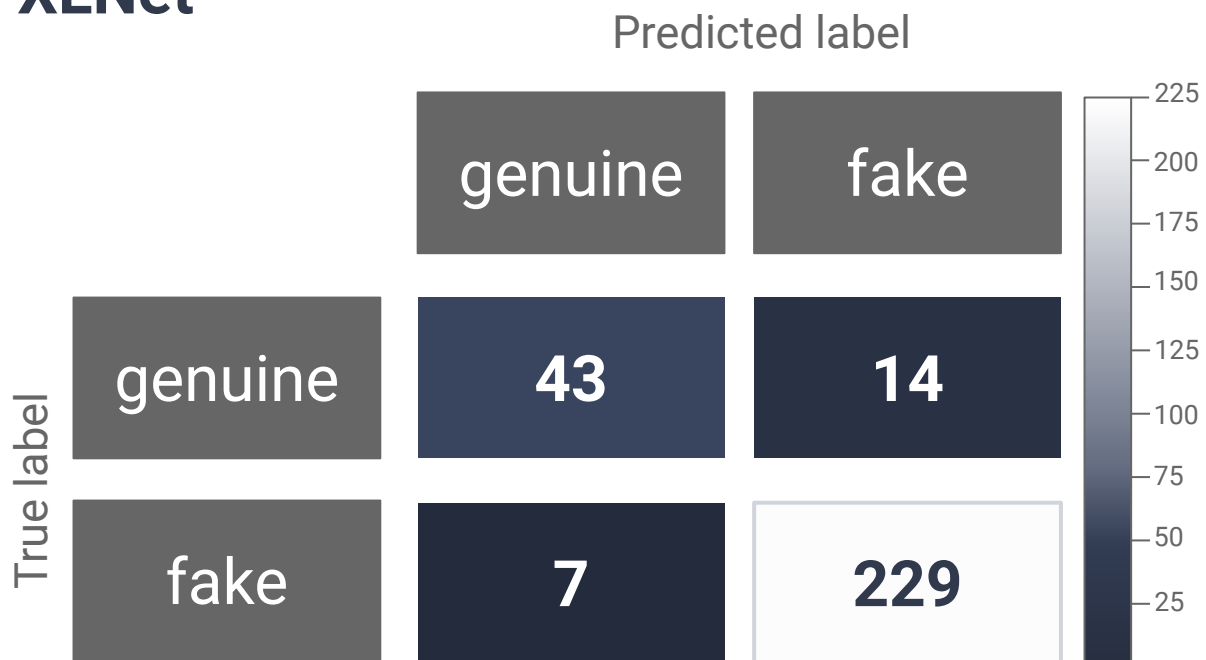
**GPT- 4**



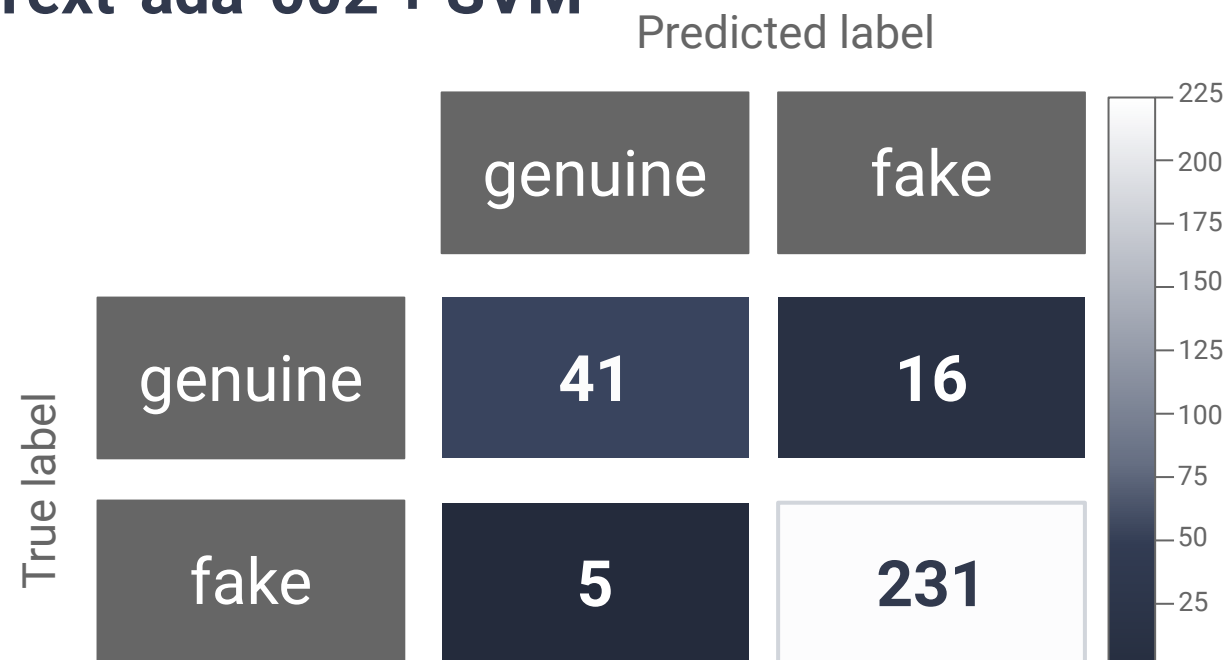


# Confusion Matrix (XLNet v.s. Text-ada-002 + SVM)

**XLNet**



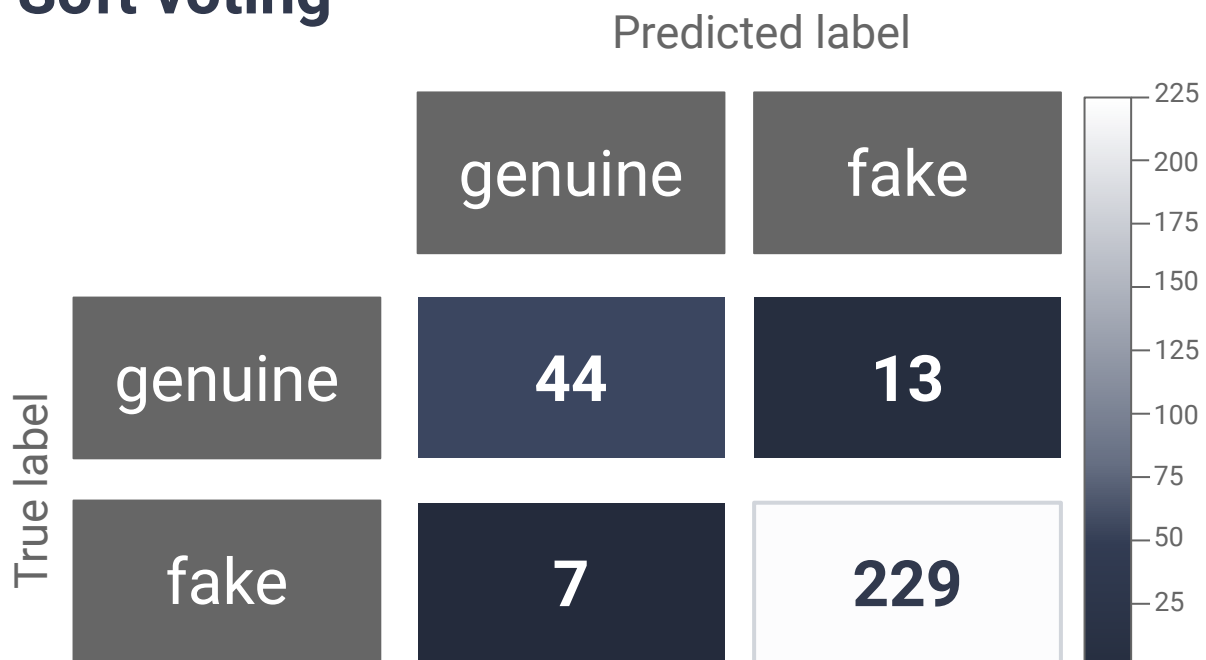
**Text-ada-002 + SVM**



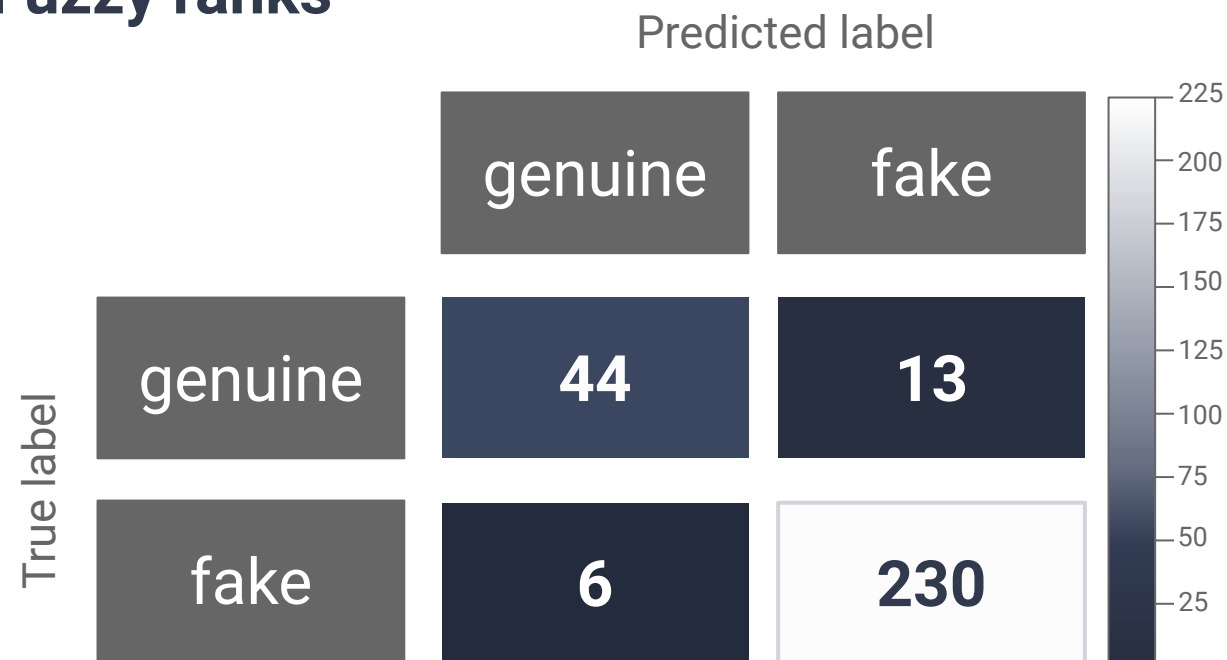
# Confusion Matrix (Soft voting v.s. Fuzzy ranks)



## Soft voting



## Fuzzy ranks



# Real case inference using fuzzy method

Content / Mainly claim	ori_length	tokenize_length	prediction	ground_truth
A German study has revealed long COVID is linked to the vaccine.	12	15	fake	fake
Covid vaccination before infection strongly linked to reduced risk of developing long covid	13	15	genuine	genuine
Long COVID's causes and risk factors remain a subject of ongoing research, with potential factors including reactivation of SARS-CoV-2 particles, overactive immune responses, and the development of autoantibodies attacking organs. Certain groups, such as those with severe COVID-19 history, underlying health conditions, or lacking vaccination, are at higher risk, alongside other factors like sex, age, initial immune response, and viral variants. Health inequities may also contribute, especially affecting racial or ethnic minority groups and individuals with disabilities.	269	367	genuine	genuine
While Omicron's subvariants find new ways to evade vaccines and destabilize immune systems, another pandemic has overwhelmed officials who are supposed to be in charge of public health. In any case, COVID, a novel virus that can wreak havoc with vital organs in the body, continues to evolve at a furious pace. In response officials have largely abandoned any coherent response, including masking, testing, tracing and even basic data collection. Yes, the people have been abandoned. So don't expect "normal" to return to your hospital, your airport, your nation, your community or your life anytime soon.	343	469	fake	fake

# Discussion

## Principal Findings

01

Attention-based models outperformed traditional baseline model.

02

Model architecture and optimization techniques are essential.

Model	Parameters
HAN	2,343,202
BERT	109,483,778
RoBERTa	124,647,170
DeBERTa	139,193,858
XLNet	117,310,466
GPT4, Gemini, text-ada-002	unknown

# Discussion

A decorative network diagram in the top right corner, consisting of several circular nodes connected by thin lines, forming a web-like structure.

## Principal Findings

03

The similar sentiment distribution across genuine and fake texts may contribute to the limitations observed in the content-based feature approach.

04

The ensemble methods demonstrated the effectiveness in the experiments.

# Discussion



## Regarding embedding models

05

Text-ada-002 performs slightly better than Gemini.

06

GPT-4 falls short compared to training SVM on vector-transformed training data using embedding models.

07

LLM-based embedding models performed well, but accessing them demand charges.

# Discussion

## Overall

07

Ensemble methods can combine open-source PLMs to achieve even better results.

08

Fuzzy fusion-based technique allows for determining ensemble model weights for each test case, resulting in superior performance.

# Limitation

A network diagram consisting of several circular nodes connected by lines, forming a web-like structure, located in the top right corner of the slide.

## **Data imbalance**

The imbalance in the data suggests a potential bias towards the prevalence of fake information on the internet.

## **Single label**

There may need to be more than a single label for classification to accurately distinguish between truthful and false segments within an article.



# Conclusion

A decorative network diagram in the top right corner, consisting of several circular nodes connected by thin lines, forming a web-like structure.

01

The study has shown the strength of attention-based models compared to the baseline model and state-of-the-art LLM.

02

The fuzzy rank-based ensemble technique with PLMs presented an approach, offering potential improvements.

03

Experimental results indicated that training solely on textual content can achieve high performance.



*Thank you*

SMART COMPUTING