# Fake News Detection using Machine Learning

**Mr. Harshwardhansinh K. Chauhan[1] and Dr. Sheshang Degadwala[2]**

Student, Department of Computer Engineering[1]

Associate Professor & Head of Department, Department of Computer Engineering[2]

Sigma Institute of Engineering, Gujrat Technological University, Gujrat, India[1]

Sigma University, Vadodara, Gujarat, India[2]

hwkc392@gmail.com[1] and sheshang13@gmail.com[2]

**Abstract:** *The approach of the Internet and the quick reception of online entertainment stages (like Facebook and Twitter) prepared for data scattering that has never been seen in mankind's set of experiences previously. With the ongoing use of online entertainment stages, shoppers are making and sharing more data than at any other time in recent memory, some of which are deluding with no importance to the real world. Mechanized characterization of a text article as deception or disinformation is difficult. Indeed, even a specialist in a specific space needs to investigate various viewpoints before deciding on the honesty of an article. In this work, we propose to utilize an AI-gathering approach for robotized order of news stories. Our review investigates different printed properties that can be utilized to separate phony items from genuine ones. By utilizing those properties, we train a mix of various AI calculations utilizing different troupe strategies and assess their exhibition on 4 genuine world datasets. Exploratory assessment affirms the unrivaled presentation of our proposed outfit student approach in contrast with individual students. Counterfeit news discovery is bit by bit happening to principal significance to our society to keep away from the alleged reality dizziness and safeguard specifically the less instructed people. Different AI procedures have been proposed to resolve this issue. This article presents a thorough exhibition assessment of eight AI calculations for counterfeit news identification/characterization*

**Keywords:** Counterfeit news discovery is bit by bit happening to principal significance to our society to keep away from the alleged reality dizziness and safeguard specifically the less instructed people. Different AI procedures have been proposed to resolve this issue. This article presents a thorough exhibition assessment of eight AI calculations for counterfeit news identification/characterization.

## I. INTRODUCTION

**What is Fake Information**

Counterfeit news alludes to data content that is bogus, deluding or whose source can't be checked. This content might be created to purposefully harm notorieties, misdirect, or acquire consideration

**What is Phony Information?**

A kind of sensationalist reporting, counterfeit news embodies bits of information that might be deceptions and is for the most part spread through virtual entertainment and other web-based media. This is frequently finished to further or force specific thoughts and is frequently accomplished with political plans. Such news might contain misleading and misrepresented guarantees and may turn out to be virtualized by calculations, and clients might wind up in a channel bubble.

**How to Make a Fake News Identification Framework?**

To make a Fake news identification framework and to make the framework practical, python gives a lot of libraries. To comprehend how to make a framework utilizing Python and make it practical for the Fake News recognition framework

### About Identifying Counterfeit News with Python

This best-in-class Python venture of distinguishing counterfeit news manages phony and genuine news. Utilizing ski-kit learn, we fabricate a TfidfVectorizer on our dataset. Then, we introduce a PassiveAggressive Classifier and fit the model. Eventually, the precision score and the disarray grid let us know how well our model tolls.

Fake News identifier is a characteristic language handling task that includes recognizing and grouping news stories or different kinds of text as genuine or counterfeit. The objective of phony news recognition is to foster calculations that can consequently recognize and signal phony news stories, which can be utilized to battle deception and advance the scattering of precise data.

### Distinguishing Counterfeit News with Python

To construct a model to characterize a piece of information as Genuine or Counterfeit precisely.

### The Fake News Dataset

The dataset we'll use for this Python project-we'll call it news.csv. This dataset has a state of 7796×4. The principal section distinguishes the news, the second and third are the title and text, and the fourth segment has marks indicating whether the news is Genuine or Counterfeit.

The central meaning of fake news is data that leads individuals wrong. These days, counterfeit news gets out far and wide, and individuals share it without affirming it. This is often finished to progress or authorize explicit convictions and is habitually achieved through political plans.

The capacity to attract clients to media associations' sites is expected to make web-based promoting income. Thus, perceiving fake news is indispensable.

### Different Sorts of Fake news include

Misleading content. Frequently attractive substance to catch perusers to the detriment of being genuine.

Parody/spoof. This sort of satisfied is viewed as tomfoolery and funny accordingly viewed as engaging, yet a few perusers might decipher the substance as reality.

Misleading publicity. This is content intended to misdirect and impact the peruser.

One-sided/sectarian/hyper-hardliner. Frequently this is one-sided political substance professing to be fair-minded.

Untrustworthy news. Columnists might distribute news whose sources are unconfirmed, or without completing any type of reality actually taking a look at themselves.

### Fake News Functions

Online entertainment stages are unbelievably persuasive. As per web live details, the assessed everyday number of tweets is around 500 million. These stages are omnipresent. They are the go-to climate to share considerations, sentiments, assessments, and aims. This furnishes ideal circumstances to convey news with insignificant rules and limitations.

In this day and age, it is typical to get news from online sources like web-based entertainment. News is frequently emotional to perusers. We frequently decide to ingest content that requests the various feelings we have. Thus, taking into account this, the data that gets the most reach may not be genuine or exact information. Moreover, genuine news might be wound in transmission. A peruser may wind up with various variants of a similar news. This might prompt data over-burden.

### Why you Ought to Mind

At the point when the globe is characterized by a pandemic, general well-being relies upon dependable data. However, we gaze intently at the barrel of an infodemic. An infodemic is the blend of the word data and scourge. It is an unreasonable measure of data about an issue that makes the arrangement more troublesome. It likewise characterizes a wide and quick spread of falsehood.

This implies that our singular well-being is an aggregate liability. It is attached to the way of behaving of others since news impacts the way of behaving of the crowd. The World Wellbeing Association has featured the risks of a

Coronavirus was driven infodemic. It presents as much risk as the actual infection. As per WHO, counterfeit news gets out quicker and more effectively than the infection**.**

**Robotized Counterfeit News Identification**

Mechanized location frameworks offer some benefit regarding mechanization and versatility. There are different strategies and approaches executed in counterfeit news discovery research. What's more, it is actually important that these methodologies frequently cross-over contingent upon point of view. According to my own viewpoint, I decide to talk about just two methodologies.

These two methodologies center around the strategies utilized, instead of the substance being investigated. They may likewise both include Regular Language Handling (NLP) in their approach.

Counterfeit news on various stages is spreading generally and involves serious worry, as it causes social conflicts and super durable breakage of the bonds laid out among individuals. A ton of exploration is as of now going on zeroed in on the order of phony news.Here we will attempt to settle this issue with the assistance of AI in Python.Prior to beginning the code, download the dataset by tapping the connection.

Counterfeit news identification/grouping is slowly happening to vital significance to out society to keep away from the purported reality dizziness, and safeguard specifically the less taught people. Different AI methods have been proposed to resolve this issue. This article presents an exhaustive presentation assessment of eight AI calculations for counterfeit news recognition/grouping.

Virtual entertainment for news utilization is a situation with two sides. From one viewpoint, its minimal expense, simpleaccess, and fast scattering of data lead individuals to search out and consume news fromsocial mediaOn the other hand, it empowers the boundless of \fake news", i.e., bad quality newswith deliberately misleading data.

The broad spread of phony news has the potential forincredibly adverse consequences on people and society. Thusly, counterfeit news identification on friendlymedia has as of late turned into arising research that is drawing in colossal consideration. Counterfeitnews identification via virtual entertainment presents one-of-a-kind qualities and difficulties that make existingrecognition calculations from customary news media ineffective or not material. In the first place, counterfeit news ispurposefully written to delude perusers to accept bogus data, which makes it troublesome andnontrivial to distinguish in view of information content; hence, we really want to incorporate helper data,like client social commitment via online entertainment, to assist with making an assurance.

Second,taking advantage of this helper data is trying all by itself as clients' social commitment to counterfeit news produces information that is huge, deficient, unstructured, and boisterous. Since the issueof phony news discovery via virtual entertainment is both testing and pertinent, we directed thisstudy to additionally work with research on the issue. In this study, we present an exhaustivesurvey of distinguishing counterfeit news via virtual entertainment, remembering counterfeit news portrayals forbrain research and social speculations, existing calculations from an information mining viewpoint, assessmentmeasurements and agent datasets. We additionally talk about related research regions, open issues, andfuture examination headings for counterfeit news discovery via virtual entertainment.

**What are the three 3 essential ways to deal with irregularity location?**

Picture resultThere are three primary classes of irregularity location procedures: solo, semi-endlessly administered.

The peculiarity of Phony news is encountering a fast and developing advancement with the development of the method for correspondence and Virtual entertainment. Counterfeit news recognition is an arising research region which is acquiring huge interest. It faces anyway a difficulties because of the restricted assets, for example, datasets and handling and breaking down strategies. In this work, we propose a framework for Counterfeit news identification that utilizations AI methods. We utilized term recurrence reverse record recurrence (TF-IDF) of pack of words and n-grams as element extraction strategy, and Backing Vector Machine (SVM) as a classifier. We propose likewise a dataset of phony and genuine news to prepare the proposed framework. Gotten results show the proficiency of the framework. In this work, we propose a framework for Counterfeit news recognition that utilizations AI procedures. We utilized term recurrence opposite report recurrence (TF-IDF) of sack of words and n-grams as element extraction procedure, and Backing Vector Machine (SVM) as a classifier. We propose likewise a dataset of phony and genuine news to prepare the proposed framework. Acquired results show the productivity of the framework.

As a rising measure of our lives is spent communicating on the web through virtual entertainment stages, an ever increasing number of individuals will generally chase out and consume news from web-based entertainment rather than customary news organizations.[1] The clarifications for this change in utilization ways of behaving are inborn inside the idea of those web-based entertainment stages: (I) it's not unexpected all the more opportune and less costly to consume news via virtual entertainment contrasted and conventional reporting , like papers or TV; and (ii) it's simpler to additional offer, examine , and talk about the news with companions or different perusers via web-based entertainment. For example, 62% of U.S. grown-ups get news via virtual entertainment in 2016, while in 2012; just 49 percent detailed seeing news via web-based entertainment [1]. It had been likewise found that virtual entertainment presently beats TV in light of the fact that the significant news source. Notwithstanding the advantages given by web-based entertainment, the norm of stories via virtual entertainment is not exactly conventional news associations. Nonetheless, on the grounds that it's reasonable to supply news on the web and far quicker and simpler to engender through virtual entertainment, huge volumes of artificial new

## II. METHOD AND MATERIALS

Import important Libraries

```python
import numpy as np
import pandas as pd
import re
from nltk.corpus import stopwords
from nltk.stem.porter import PorterStemmer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
```

displaying the English stop words

```python
# printing the stopwords in English
print(stopwords.words('english'))
```

ourselves', 'you', "you're", "you've", "you'll", "you'd", 'your', 'yours', 'yourself', 'yourselves', 'he', 'him', 'his', 'himself', 'she', "she's",

stacking the dataset to a pandas DataFrame

```
news_dataset = pd.read_csv('/content/train.csv')
```

```
news_dataset.shape
```

```
(20800, 5)
```

print the initial 5 columns of the dataframe

```
news_dataset.head()
```

| | id | title | author | text | label |
|---|---|---|---|---|---|
| 0 | 0 | House Dem Aide: We Didn't Even See Comey's Let... | Darrell Lucus | House Dem Aide: We Didn't Even See Comey's Let... | 1 |
| 1 | 1 | FLYNN: Hillary Clinton, Big Woman on Campus - ... | Daniel J. Flynn | Ever get the feeling your life circles the rou... | 0 |
| 2 | 2 | Why the Truth Might Get You Fired | Consortiumnews.com | Why the Truth Might Get You Fired October 29, ... | 1 |
| 3 | 3 | 15 Civilians Killed In Single US Airstrike Hav... | Jessica Purkiss | Videos 15 Civilians Killed In Single US Airstr... | 1 |
| 4 | 4 | Iranian woman jailed for fictional unpublished... | Howard Portnoy | Print \nAn Iranian woman has been sentenced to... | 1 |

counting the quantity of missing qualities in the dataset

```
news_dataset.isnull().sum()
```

```
id             0
title        558
author      1957
text          39
label          0
dtype: int64
```

supplanting the invalid qualities with void string and consolidating the creator name and news title:

**Copyright to IJARSCT**
**www.ijarsct.co.in**

**DOI: 10.48175/IJARSCT-9117**

ISSN
2581-9429
IJARSCT

197

```
news_dataset = news_dataset.fillna('')
```

```
# merging the author name and news title
news_dataset['content'] = news_dataset['author']+' '+news_dataset['title']
```

```
print(news_dataset['content'])
```

```
0         Darrell Lucus House Dem Aide: We Didn't Even S...
1         Daniel J. Flynn FLYNN: Hillary Clinton, Big Wo...
2         Consortiumnews.com Why the Truth Might Get You...
3         Jessica Purkiss 15 Civilians Killed In Single ...
4         Howard Portnoy Iranian woman jailed for fictio...
                             ...
20795     Jerome Hudson Rapper T.I.: Trump a 'Poster Chi...
20796     Benjamin Hoffman N.F.L. Playoffs: Schedule, Ma...
20797     Michael J. de la Merced and Rachel Abrams Macy...
20798     Alex Ansary NATO, Russia To Hold Parallel Exer...
20799               David Swanson What Keeps the F-35 Alive
Name: content, Length: 20800, dtype: object
```

isolating the information and mark

```
X = news_dataset.drop(columns='label', axis=1)

Y = news_dataset['label']
```

**Copyright to IJARSCT**
**www.ijarsct.co.in**

**DOI: 10.48175/IJARSCT-9117**

ISSN
2581-9429
IJARSCT

198

```
print(X)
print(Y)
```

```
           id ...                                         content
0           0 ...  Darrell Lucus House Dem Aide: We Didn't Even S...
1           1 ...  Daniel J. Flynn FLYNN: Hillary Clinton, Big Wo...
2           2 ...  Consortiumnews.com Why the Truth Might Get You...
3           3 ...  Jessica Purkiss 15 Civilians Killed In Single ...
4           4 ...  Howard Portnoy Iranian woman jailed for fictio...
...       ... ...                                             ...
20795   20795 ...  Jerome Hudson Rapper T.I.: Trump a 'Poster Chi...
20796   20796 ...  Benjamin Hoffman N.F.L. Playoffs: Schedule, Ma...
20797   20797 ...  Michael J. de la Merced and Rachel Abrams Macy...
20798   20798 ...  Alex Ansary NATO, Russia To Hold Parallel Exer...
20799   20799 ...          David Swanson What Keeps the F-35 Alive

[20800 rows x 5 columns]
0           1
1           0
2           1
3           1
4           1
           ..
20795       0
20796       0
20797       0
20798       1
20799       1
Name: label, Length: 20800, dtype: int64
```

Stemming:

Stemming is the most common way of decreasing a word to its Root word model: entertainer, entertainer, acting - - > act

```python
port_stem = PorterStemmer()

def stemming(content):
    stemmed_content = re.sub('[^a-zA-Z]',' ',content)
    stemmed_content = stemmed_content.lower()
    stemmed_content = stemmed_content.split()
    stemmed_content = [port_stem.stem(word) for word in stemmed_content if not word in stopwords.words('english')]
    stemmed_content = ' '.join(stemmed_content)
    return stemmed_content

news_dataset['content'] = news_dataset['content'].apply(stemming)
```

**Copyright to IJARSCT**
**www.ijarsct.co.in**

**DOI: 10.48175/IJARSCT-9117**

ISSN
2581-9429
IJARSCT

199

isolating the information and name

```
X = news_dataset['content'].values
Y = news_dataset['label'].values
```

```
print(X)
```

```
['darrel lucu hous dem aid even see comey letter jason chaffetz tweet'
 'daniel j flynn flynn hillari clinton big woman campu breitbart'
 'consortiumnew com truth might get fire' ...
 'michael j de la merc rachel abram maci said receiv takeov approach hudson bay new york time'
 'alex ansari nato russia hold parallel exercis balkan'
 'david swanson keep f aliv']
```

```
print(Y)
```

```
[1 0 1 ... 0 1 1]
```

```
Y.shape
```

```
(20800,)
```

```
print(Y)
```

```
[1 0 1 ... 0 1 1]
```

```
Y.shape
```

```
(20800,)
```

changing the text based information over completely to mathematical information

```
vectorizer = TfidfVectorizer()
vectorizer.fit(X)

X = vectorizer.transform(X)
```

```
print(X)
```

```
  (0, 15686)    0.28485063562728646
  (0, 13473)    0.2565896679337957
  (0, 8909)     0.3635963806326075
  (0, 8630)     0.29212514087043684
  (0, 7692)     0.24785219520671603
  (0, 7005)     0.21874169089359144
  (0, 4973)     0.233316966909351
  (0, 3792)     0.2705332480845492
  (0, 3600)     0.3598939188262559
  (0, 2959)     0.2468450128533713
  (0, 2483)     0.3676519686797209
  (0, 267)      0.27010124977708766
  (1, 16799)    0.30071745655510157
  (1, 6816)     0.1904660198296849
  (1, 5503)     0.7143299355715573
  (1, 3568)     0.26373768806048464
  (1, 2813)     0.19094574062359204
  (1, 2223)     0.3827320386859759
  (1, 1894)     0.15521974226349364
  (1, 1497)     0.2939891562094648
  (2, 15611)    0.41544962664721613
  (2, 9620)     0.49351492943649944
  (2, 5968)     0.3474613386728292
  (2, 5389)     0.3866530551182615
  (2, 3103)     0.46097489583229645
   :      :
  (20797, 13122)    0.2482526352197606
  (20797, 12344)    0.27263457663336677
  (2, 15611)    0.41544962664721613
  (2, 9620)     0.49351492943649944
  (2, 5968)     0.3474613386728292
  (2, 5389)     0.3866530551182615
  (2, 3103)     0.46097489583229645
   :      :
  (20797, 13122)    0.2482526352197606
  (20797, 12344)    0.27263457663336677
  (20797, 12138)    0.24778257724396507
  (20797, 10306)    0.08038079000566466
  (20797, 9588)  0.174553480255222
  (20797, 9518)  0.2954204003420313
  (20797, 8988)  0.36160868928090795
  (20797, 8364)  0.22322585870464118
  (20797, 7042)  0.21799048897828688
  (20797, 3643)  0.21155500613623743
  (20797, 1287)  0.33538056804139865
  (20797, 699)   0.306858460799762347
  (20797, 43)    0.29710241860700626
  (20798, 13046)    0.22363267488270608
  (20798, 11052)    0.4460515589182236
  (20798, 10177)    0.3192496370187028
  (20798, 6889)  0.32496285694299426
  (20798, 5032)  0.4083701450239529
  (20798, 1125)  0.4460515589182236
  (20798, 588)   0.3112141524638974
  (20798, 350)   0.28446937819072576
  (20799, 14852)    0.5677577267055112
  (20799, 8036)  0.45983893273780013
  (20799, 3623)  0.37927626273066584
  (20799, 377)   0.5677577267055112
```

Parting the dataset to preparing and test information:-

```
X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size = 0.2, stratify=Y, random_state=2)
```

Training the Model: Logistic Regression

```
model = LogisticRegression()
```

```
model.fit(X_train, Y_train)
```

```
LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
                   intercept_scaling=1, l1_ratio=None, max_iter=100,
                   multi_class='auto', n_jobs=None, penalty='l2',
                   random_state=None, solver='lbfgs', tol=0.0001, verbose=0,
                   warm_start=False)
```

**Assessment and exactness score:-**

exactness score on the preparation information:-

```
X_train_prediction = model.predict(X_train)
training_data_accuracy = accuracy_score(X_train_prediction, Y_train)
```

```
print('Accuracy score of the training data : ', training_data_accuracy)
```

```
Accuracy score of the training data :  0.9865985576923076
```

exactness score on the test information :-

```
X_test_prediction = model.predict(X_test)
test_data_accuracy = accuracy_score(X_test_prediction, Y_test)
```

```
print('Accuracy score of the test data : ', test_data_accuracy)
```

```
Accuracy score of the test data :  0.9790865384615385
```

Making a Prescient Framework

```
] X_new = X_test[3]

  prediction = model.predict(X_new)
  print(prediction)

  if (prediction[0]==0):
    print('The news is Real')
  else:
    print('The news is Fake')

  [0]
  The news is Real

] print(Y_test[3])

  0
```

## III. CONCLUSION

In the research of Fake news Detection, we do many things, first of all, Importing the Dependencies , display the English stop words, print the initial 5 columns of the data frame, counting the number of missing qualities in the dataset, supplanting the invalid qualities with void string and consolidating the creator name and news title: isolating the information and mark, isolating the information and name, Parting the dataset to preparing and test information, Training the Model: Logistic Regression, Assessment, and exactness score, exactness score on the preparation information, exactness score on the test information, Making a Prescient Framework. From this code and project we can able to detect the news is fake is real.

## REFERENCES

[1]. Harshwardhansinh K. Chauhan, Dr. Sheshang Degadwala, "Project Base Prediction Using Machine Learning and Deep Learning", International Journal of Scientific Research in Computer Science, Engineering and Information Technology (IJSRCSEIT), ISSN : 2456-3307, Volume 9, Issue 2, pp.22-29, March-April2023. Available at doi : https://doi.org/10.32628/CSEIT2390150 Journal URL : https://ijsrcseit.com/CSEIT2390150

[2]. Harshwardhansinh K. Chauhan, Dr. SheshangDegadwala, "Sensitivity Analysis of Project using Machine Learning ", International Journal of Scientific Research in Science, Engineering and Technology (IJSRSET), Online ISSN : 2394-4099, Print ISSN : 2395- 1990, Volume 10 Issue 2, pp. 197-202, March-April 2023. Available at doi : https://doi.org/10.32628/IJSRSET2310128 Journal URL : https://ijsrset.com/IJSRSET2310128

[3]. Harshwardhansinh K. Chauhan, Miss Kamini R. Parmar, Dr. SheshangDegadwala, " Financial Risk Analysis using Machine Learning", International Journal of Scientific Research in Science and Technology(IJSRST), Print ISSN : 2395-6011, Online ISSN : 2395-602X, Volume 10, Issue 2, pp.352-358, March-April-2023. Available at doi : https://doi.org/10.32628/IJSRST2310117 Journal URL : https://ijsrst.com/IJSRST2310117

[4]. 4) Proceedings of the Sixth International Conference on I-SMAC (IoT in Social, Mobile, Analytics and Cloud) (I-SMAC-2022). IEEE Xplore Part Number: CFP22OSV-ART; ISBN: 978-1-6654-6941-8- Lung Respiratory Audio Prediction using Transfer Learning Models Arohi Patel Assistant professor Sigma institute of engineering Vadodara, Gujarat, India arohipatel3010@gmail.com, SheshangDegadwala Associate professor Sigma institute of engineering Vadodara, Gujarat, India sheshang13@gmail.com, SheshangDegadwala Associate professor Sigma institute of engineering Vadodara, Gujarat, India

[5]. Proceedings of the Sixth International Conference on Electronics, Communication and Aerospace Technology (ICECA 2022) IEEE Xplore Part Number: CFP22J88-ART; ISBN: 978-1-6654-8271-4 Mihir Prajapati Mitul Nakrani Dr. Tarjni Vyas Computer Science Engineering Computer Science Engineering Assistant Professor Nirma University Nirma University Nirma University Ahmedabad, India Ahmedabad, India Ahmedabad, India 19bce128@nirmauni.ac.in

[6]. Proceedings of the Sixth International Conference on I-SMAC (IoT in Social, Mobile, Analytics and Cloud) (I-SMAC-2022). IEEE Xplore Part Number: CFP22OSV-ART; ISBN: 978-1-6654-6941-8 - Crop Prediction System based on Soil and Weather Characteristics JayashriMahale Assistant professor Sigma institute of engineering Vadodara, Gujarat, India

[7]. Siyu Chen *, Guohua Fang *, Xianfeng Huang and Yuhong Zhang College of Water Conservancy and Hydropower Engineering, Hohai University, Nanjing 210098, China

[8]. Gonzalez, J. R. C., Romero, J. J. F., Guerrero, M. G., and Calderon, F. (2015). Multi-class multi-tag classifier system for StackOverflow questions. 2015 IEEE International Autumn Meeting on Power, Electronics and Computing (ROPEC).

[9]. B. Sun, Y. Zhu, Y. Xiao, R. Xiao and Y. Wei, "Automatic Question Tagging with Deep Neural Networks," in IEEE Transactions on Learning Technologies, vol. 12, no. 1, pp. 29-43, 1 Jan.-March 2019, doi: 10.1109/TLT.2018.2808187.

[10]. P. Devine and K. Blincoe, "Unsupervised Extreme Multi Label Classification of Stack Overflow Posts," 2022 IEEE/ACM 1st International Workshop on Natural Language-Based Software Engineering (NLBSE), 2022, pp. 1-8.