



ISSN: 2321-2152

IJMECE

*International Journal of modern
electronics and communication engineering*

E-Mail

editor.ijmece@gmail.com

editor@ijmece.com

www.ijmece.com

DETECTING HATE SPEECH ON TWITTER A SYSTEMATIC REVIEW OF METHODS, TAXONOMY ANALYSIS, CHALLENGES AND OPPORTUNITIES

G VISWANATH¹, DURRGA PRASAD², A DHANASEKHAR REDDY³, K SUPRIYA⁴

¹Associate Professor, Department of CSE(AIML), Sri Venkatesa Perumal College of Engineering & Technology, Puttur, Email: viswag111@gmail.com, ORCID: <https://orcid.org/0009-0001-7822-4739>

²P.G Scholar, Department of MCA, Sri Venkatesa Perumal College of Engineering & Technology, Puttur, Email: asthagiridurgaprasadnaidu@gmail.com

³Assistant Professor, Department of MCA, Sri Venkatesa Perumal College of Engineering & Technology, Puttur, Email: ghanasekhar918@gmail.com

⁴Assistant Professor, Department of CSE, Sri Venkatesa Perumal College of Engineering & Technology, Puttur, Email: srisupriya05@gmail.com

Abstract: Hate speech has soared because of virtual entertainment locales like Twitter. This broad issue causes conflicts, influences clients, and makes content control troublesome. This undertaking targets hate speech identification to diminish its effect on internet based networks and then some. The examination utilizes state of the art calculations to detect hate speech. A hearty and versatile hate speech detection utilizes ML models and DL. To ensure accuracy and reliability, model execution is completely surveyed utilizing numerous markers. Accuracy, precision, recall, and F1 score demonstrate the model's hate speech detection execution. complete separating power, showing its exhibition across limits. The thorough assessment of hate speech detection techniques yielded valuable undertaking discoveries. In spite of advances, virtual entertainment language designs stay hard to deal with. The report underlines the requirement for hate speech detection innovative work to work on satisfied control and make the web more secure. The Hate Speech identification model purposes the stacking classifier, a modern ensemble approach with 100 percent precision. The Hybrid

Approach, utilizing LSTM and BiGRU models, has 94% precision. A Flask front end with verification capacities was made to make testing simple and secure the Twitter Hate Speech Detection system. This makes assessing the model's capacity to perceive and relieve Twitter hate speech simple and dependable.

Index terms - Hate speech, classification, automatic detection, twitter, systematic review, natural language processing, social media.

1. INTRODUCTION

Twitter and other web-based entertainment have risen decisively in the earlier 10 years. As per [1], these mediums empower hate speech by advancing client namelessness and free articulation. With 300 million month to month dynamic clients, Twitter is one of the most well known person to person communication stages [2]. Twitter spreads hate speech regardless of its prominence and importance. It is one of the most famous informal organizations for mechanized hate speech ID [3], [4], [5] and harmful language study. Hate speech is ascending via web-based

entertainment. Clients become unfriendly, causing genuine showdowns and influencing organizations. Despised content is regularly taken out by online entertainment enterprises.[114]

Since English is the most by and large communicated in language and the most openly accessible information source, this study centers around Twitter virtual entertainment messages [6]. Computerizing on the web hate speech detection is required since manual screening is unbending. PC based arrangements can answer speedier than individuals, though non-mechanized positions influence it. Adding to programmed text Hate Speech discovery is critical. These realities have prodded NLP research. Hate Speech writing develops. Research people group have allotted this undertaking as supervised record classification utilizing NLP and AI [7]. Twitter was perhaps of the greatest social medium firms in 2017. It modified their protection strategy including misuse. These rules apply to tweets that advance maltreatment, provocation, self-destruction, self-hurt, viciousness, scorn, and so forth [8].

Specialists have extended their endeavors to distinguish Hate Speech on Twitter. In any case, non-English datasets are restricted. English is the most generally communicated in language. It is additionally the significant can't stand content identifier. Hate Speech is difficult to characterize since it has many structures. The term most frequently utilized for this event is Hate Speech, which is legitimate in numerous countries [7]. Numerous meanings of Hate Speech exist in writing.

In view of an examination of different depictions in the writing, reference [9] characterized Hate Speech as

language that assaults or decreases, actuates viciousness or disdain against bunches in light of explicit qualities, like actual appearance, religion, plummet, public or ethnic beginning, sexual direction, orientation personality, or others, and it can happen in unobtrusive or amusing structures. Twitter Hate Speech models: "Twitter client Pu**y a** ni**a" and "You disdain football you are a fa**ot." [10].

Many hate speech detection systems have been made lately, outperforming their procedures. In any case, the evaluations for the most part find non-disdain things as opposed to arranging threatening ones [1]. Since virtual entertainment language is developing quickly, the vast majority of these endeavors are presently battling to find an answer [9]. In this manner, a careful consciousness of the ongoing writing is required. Hate speech detection has advanced for quite a long time, yet there is no exhaustive examination assessment. SLR papers assist with finding remarkable subjects and examination holes on a particular region.

2. LITERATURE SURVEY

Web-based entertainment gives Web clients a well disposed spot to communicate their thoughts. This area offers astonishing correspondence potential yet additionally huge issues. Online hate speech epitomizes such issues. In spite of its size, virtual entertainment disdain discourse is inadequately perceived. The principal efficient enormous scope estimating examination of online virtual entertainment hate speech targets is introduced in this exploration [1]. We gather Murmur and Twitter follows for that. We then, at that point, make and test a hate speech detection algorithm for the two frameworks. Our discoveries uncover online hate speech types and

upgrade how we might interpret the peculiarities, directing anticipation and recognition [1].

A large number of individuals overall rely upon web-based entertainment. It allows individuals to pass their perspectives on to a huge crowd. Because of this openness of articulation, falsehood and despise discourse have spread generally [2]. Bigots use code (Activity Google) to sidestep web-based entertainment misuse limitations and robotized frameworks like Google's Discussion AI. In disdainful Tweets and postings, harmless expressions are utilized rather than local area references [40,54]. Clients have called African-Americans and Asians researches and Bings. The rundown of individuals who submit such satisfied allows us to explore the use example of these concentrated clients, moving past tweet grouping.

Web-based entertainment stages permit everybody to economically deliver and share material. Web-based entertainment can spread poisonous talks to specific networks. These talks incorporate harassing, offending material, and disdain discourse [4]. Numerous nations rapidly consider disdain discourse to be an extreme issue from these discussions. This work is the main precise huge scope estimating and logical examination of unequivocal disdain discourse in virtual entertainment [63, 90, 92]. We need to grasp the pervasiveness of disdain discourse via virtual entertainment, the most famous disdain articulations, the effect of namelessness, the awareness of disdain discourse, and the most hated bunches across geologies. We gather Murmur and Twitter follows to satisfy our objectives. We then make and test a disdain discourse identification calculation for the two frameworks. Our outcomes recognize disdain discourse types and uncover basic examples, growing

comprehension we might interpret online disdain discourse and giving recognition and avoidance systems.[116]

Hostility is essential to grasping human way of behaving. Individual, conduct, propensities, climate, and emotional wellness are connected. Understanding classifications of forcefulness and forceful lead can help counter online entertainment animosity [8]. An examination combination utilizing different web search tools to look for a decent job on hate speech detection, disdain, outrage, forceful conduct in virtual entertainment, and the results of these terms found that past techniques disregard discourse assortment, conceivable different classifications of disdain discourse, the relationship of discourse to human way of behaving, and negligible to no compassion toward clients. Just disdain and offending words are named disdain discourse [91,92,93,96]. Outrage ought to be incorporated. Future exploration ought to zero in on forceful lead since it joins human way of behaving to can't stand discourse.

Long range interpersonal communication causes web clients to feel appreciated. Along these lines they transparently voice their thoughts. Clients might advance unforgiving talk on the web due of its transparency. Manual recognizable proof of frightful data is tedious and may miss some [7]. Subsequently, programmed hostile substance distinguishing proof is fundamental to perceive and assess how much unsavory text in online entertainment. Parametric examination of robotized disdain message acknowledgment strategies is introduced in this study [93,105,108].

3. METHODOLOGY

i) Proposed Work:

The proposed framework propels hate speech detection through cutting edge NLP, ML, DL, and gathering models (e.g., stacking classifier, voting classifier) [30, 31]. This framework will be prepared on different datasets. Semantic and sentiment analysis will further develop hate speech recognizable proof by extending setting information. Continuous programmed ID will accelerate web-based entertainment disdain discourse sifting. The hate speech detection model purposes the stacking classifier, a refined troupe approach with 100 percent accuracy. The Hybrid Methodology, utilizing LSTM and BiGRU models, has 94% accuracy. A Flask front end with verification capacities was made to make testing simple and secure the Twitter Hate Speech Detection system. This makes assessing the model's capacity to perceive and alleviate Twitter disdain discourse simple and reliable.[118]

ii) System Architecture:

Import the Stock Tweets Dataset, Single Stock Information, and Multi-Source Information. These databases support sentiment analysis and stock price prediction. Stock Tweets Dataset text is cleaned of accentuations, HTML components, URLs, and emojis. This plans message for opinion examination [29]. Handled Single Stock Information and Multi-Source Information wipe out copies, oversee invalid qualities, and scale. This gives monetary information to stock cost conjecture. For feeling order, MLP, CNN, LSTM, MS-LSTM, MS-SSA-LSTM, Voting Classifier, and LSTM + GRU are prepared. Market feeling is determined utilizing scrubbed tweet information. For stock cost expectation, MLP, CNN, LSTM, MS-

LSTM, MS-SSA-LSTM [63,65,94], and expansion Voting Regression are prepared. Monetary information is utilized to anticipate stock costs. Models conjecture in the wake of preparing. Market feeling is shown through figures in opinion examination. Stock cost forecast techniques gauge future costs. Opinion examination and stock cost models assist financial backers and brokers with making decisions. The consolidated outcomes help clients explore the convoluted securities exchange, decline chances, and amplify rewards.



Fig 1 Proposed architecture

iii) Dataset collection:

The twitter dataset should be stacked and investigated for this venture. Investigating the dataset's design, missing qualities, and class dispersion (hate speech vs. non-hate speech) is finished. Data on dataset attributes is likewise procured.



id	text	label
1	1. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
2	2. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
3	3. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
4	4. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
5	5. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
6	6. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
7	7. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
8	8. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
9	9. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
10	10. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
11	11. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
12	12. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
13	13. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
14	14. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
15	15. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
16	16. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
17	17. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
18	18. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
19	19. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
20	20. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
21	21. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
22	22. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
23	23. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
24	24. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
25	25. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
26	26. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
27	27. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
28	28. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
29	29. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
30	30. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
31	31. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
32	32. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
33	33. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
34	34. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
35	35. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
36	36. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
37	37. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
38	38. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
39	39. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
40	40. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
41	41. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
42	42. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
43	43. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
44	44. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
45	45. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
46	46. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
47	47. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
48	48. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
49	49. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
50	50. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
51	51. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
52	52. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
53	53. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
54	54. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
55	55. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
56	56. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
57	57. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
58	58. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
59	59. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
60	60. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
61	61. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
62	62. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
63	63. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
64	64. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
65	65. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
66	66. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
67	67. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
68	68. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
69	69. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
70	70. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
71	71. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
72	72. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
73	73. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
74	74. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
75	75. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
76	76. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
77	77. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
78	78. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
79	79. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
80	80. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
81	81. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
82	82. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
83	83. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
84	84. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
85	85. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
86	86. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
87	87. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
88	88. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
89	89. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
90	90. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
91	91. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
92	92. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
93	93. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
94	94. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
95	95. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
96	96. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
97	97. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
98	98. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
99	99. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0
100	100. @realDonaldTrump: A failure in Afghanistan and a... (truncated)	0

Fig 2 Tweets hate dataset

iv) Data Processing:

Data processing transforms crude information into business-helpful data. Information researchers accumulate, sort out, clean, confirm, investigate, and organize information into charts or papers. Information can be handled physically, precisely, or electronically. Data ought to be more important and decision-production simpler. Organizations might upgrade activities and settle on basic decisions quicker. PC programming improvement and other robotized data processing innovations add to this. Large information can be transformed into pertinent bits of knowledge for quality administration and independent direction.

v) Feature selection:

Feature selection chooses the most steady, non-repetitive, and pertinent elements for model turn of events. As data sets extend in amount and assortment, purposefully bringing down their size is significant. The fundamental reason for feature selection is to increment prescient model execution and limit processing cost.

One of the vital pieces of feature engineering is picking the main attributes for machine learning algorithms. To diminish input factors, feature selection methodologies take out copy or superfluous elements and limit the assortment to those generally critical to the ML model. Rather than permitting the ML model pick the main qualities, feature selection ahead of time enjoys a few benefits.[120]

vi) Algorithms:

BERT (Bidirectional Encoder Representations from Transformers) utilizes a transformer-based neural network to perceive and make human-like language.

BERT is encoder-as it were. The first Transformer configuration has encoder and decoder parts. BERT's encoder-just engineering accentuates fathoming input groupings over making yield successions. Conventional language models dissect text left-to-right or right-to-left. This system confines the model to the setting before the objective word [50,56].

```
def create_model(bert_model):
    input_ids = tf.keras.Input(shape=(80,), dtype='int32')
    attention_masks = tf.keras.Input(shape=(80,1), dtype='int32')

    output = bert_model([input_ids, attention_masks])
    output = output[1]

    output = tf.keras.layers.Dense(12, activation='relu')(output)
    output = tf.keras.layers.Dropout(0.2)(output)

    output = tf.keras.layers.Dense(1, activation='sigmoid')(output)
    model = tf.keras.models.Model(inputs = [input_ids, attention_masks], outputs = output)
    model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy', 'f1_score', 'precision', 'recall'])
    return model

from transformers import BertTokenizer
bert_tokenizer = BertTokenizer.from_pretrained('bert-base-uncased')
```

Fig 3 BERT

Bidirectional LSTM or a sequence model with two LSTM layers, one for forward handling and one for in reverse handling, is called BiLSTM. Generally utilized for NLP. This strategy works by handling input in the two bearings to assist the model handle with sequencing connections (e.g., grasping the following and past words in an expression). A bidirectional LSTM has two unidirectional LSTMs that cycle the grouping forward and in reverse [64].

```
Bi-LSTM

# Build the model
from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import LSTM, Dense, Dropout
from tensorflow.keras.layers import SpatialDropout2D
from tensorflow.keras.layers import Embedding

embedding_vector_length = 128

model = Sequential([
    layers.Embedding(vocab.get_vocab('tokens'), embedding_vector_length),
    layers.Bidirectional(layers.LSTM(64, return_sequences=True, recurrent_dropout=0.4)),
    layers.GlobalMaxPooling1D(),
    layers.Dropout(0.5),
    layers.Dense(128, activation='relu'),
    layers.Dropout(0.5),
    layers.Dense(num_classes, activation='softmax')
])
```

Fig 4 BiLSTM

GRU (Gated Recurrent Unit): A recurrent neural network. Contrasted with LSTM networks, it is more

straightforward. GRU processes successive text, voice, and time-series information like LSTM. GRU refreshes the organization's secret state specifically at each time step utilizing gating techniques. Data enters and leaves the organization through gating instruments. The reset and update entryways are GRU gating techniques. The reset entryway chooses the amount of the earlier disguised state to neglect, though the update door concludes how much new contribution to use. GRU yield is reliant upon refreshed secret state. GRU [35] handles consecutive information all the more productively in this review, making hate speech detection more strong and compelling.

GRU

```
model = Sequential([
    layers.Embedding(5000, 100, input_length=100),
    layers.GRU(64, return_sequences=True, recurrent_dropout=0.4),
    layers.GlobalAveragePooling1D(), # or layers.Flatten()
    layers.Dense(64, activation='relu'),
    layers.Dropout(0.4),
    layers.Dense(num_classes, activation='softmax')
])

WARNING:tensorflow:Layer gru will not use cuDNN kernels since it doesn't
fallback when running on GPU.

model.compile(loss=tf.keras.losses.SparseCategoricalCrossentropy(),
              optimizer=tf.keras.optimizers.Adam(learning_rate=0.001),
              metrics=['accuracy', 'f1_score', 'precision_score', 'recall_score'])

trained1 = model.fit(X_train, y_train,
                    epochs=5,
                    steps_per_epoch = 20,
                    validation_steps = 20,
                    validation_data=(X_val, y_val),
                    shuffle=True)
```

Fig 5 GRU

CNNs are class of deep neural networks that decipher pictures and spatial information well. CNNs use channels to catch nearby examples in text as a picture in natural language processing. This study utilizes CNNs to find neighborhood qualities and patterns in printed information to detect hate speech by detecting frightful language structures.

```
CNN

from keras.layers import Input, Embedding, Conv1D, MaxPooling1D, Flatten
def build_cnn_model():
    model = Sequential()

    model.add(Embedding(5000, 100, input_length=100))

    model.add(Conv1D(16, kernel_size=3, padding='same', activation='relu'))
    model.add(MaxPooling1D())
    model.add(Conv1D(32, kernel_size=3, padding='same', activation='relu'))
    model.add(MaxPooling1D())
    model.add(Conv1D(64, kernel_size=3, padding='same', activation='relu'))
    model.add(MaxPooling1D())
    model.add(Flatten())

    optimizer = Adam(lr=0.0001, amsgrad=True)
    model.compile(loss=tf.keras.losses.SparseCategoricalCrossentropy,
                  optimizer=optimizer, metrics=['accuracy', 'f1_score', 'precision_score', 'recall_score'])

    return model

cnn_model = build_cnn_model()

cnn_history = cnn_model.fit(x_train, y_train,
                           epochs=5,
                           steps_per_epoch = 20,
                           validation_steps = 20,
                           validation_data=(x_val, y_val),
                           shuffle=True)
```

Fig 6 CNN

CNN + LSTM (Convolutional Neural Network with Long Short-Term Memory), this hybrid architecture utilizes CNNs' neighborhood include catch and LSTMs' successive learning. The CNN layer catches spatial examples in input information, while the LSTM layer models long-range connections. CNN + LSTM is utilized to utilize neighborhood and successive data to further develop the models hate speech recognition and setting understanding.

```
CNN + LSTM

from keras.models import Sequential
from keras.layers import Input, Embedding, Conv1D, MaxPooling1D, LSTM, Flatten
from keras.layers import LSTM, Dropout, GRU, Bidirectional
from keras.preprocessing import text, sequence
from keras.callbacks import TensorBoard

model = Sequential()
model.add(Embedding(5000, 100, input_length=100))
model.add(Conv1D(16, kernel_size=3, padding='same', activation='relu'))
model.add(MaxPooling1D())
model.add(Conv1D(32, kernel_size=3, padding='same', activation='relu'))
model.add(MaxPooling1D())
model.add(LSTM(64, return_sequences=True))
model.compile(loss=tf.keras.losses.SparseCategoricalCrossentropy,
              optimizer=tf.keras.optimizers.Adam(learning_rate=0.001),
              metrics=['accuracy', 'f1_score', 'precision_score', 'recall_score'])

trainer = TensorBoard(log_dir='./logs')
callbacks = [trainer]
model.fit(x_train, y_train,
        epochs=5,
        steps_per_epoch = 20,
        validation_steps = 20,
        validation_data=(x_val, y_val),
        callbacks=callbacks)
```

Fig 7 CNN + LSTM

CNN + BiLSTM (Convolutional Neural Network with Bidirectional Long Short-Term Memory), CNN + BiLSTM joins CNNs' nearby component catch with BiLSTM's bidirectional consecutive learning, as CNN + LSTM. The model might catch

association in the two bearings by thinking about past and future setting. This hybrid configuration catches unpretentious setting and worldly examples in language, improving hate speech detection execution. These plans are picked for their capacity to catch various components of text based material, empowering more complete hate speech detection investigation.

CNN + BiLSTM

```
model = Sequential()
model.add(Embedding(1000, 100, input_length=100))
model.add(Conv1D(filters=10, kernel_size=3, padding='same', activation='relu'))
model.add(MaxPooling1D(pool_size=2))
model.add(Bidirectional(LSTM(100)))
model.add(Dense(units=100, activation='sigmoid'))
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy', 'f1_score', 'recall'])
print(model.summary())
filepath = 'weights_cnn_bi_lstm.h5'
checkpoint = ModelCheckpoint(filepath, monitor='val_acc', verbose=1, save_best_only=True, mode='max', save_weights_only=True)
callbacks_list = [checkpoint]
```

Fig 8 CNN + BiLSTM

CNN + GRU, This hybrid engineering joins CNN spatial example acknowledgment with GRU successive learning and productivity. The CNN layer assembles neighborhood qualities, and the GRU layer handles successive information. This blend is possible used in the venture to oversee neighborhood and long-range conditions, further developing hate speech identification.

CNN + GRU

```
model = Sequential()
model.add(Embedding(1000, 100, input_length=100))
model.add(Conv1D(filters=10, kernel_size=3, padding='same', activation='relu'))
model.add(MaxPooling1D(pool_size=2))
model.add(Bidirectional(GRU(100)))
model.add(Dense(units=100, activation='sigmoid'))
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy', 'f1_score', 'recall'])
print(model.summary())
filepath = 'weights_cnn_gru.h5'
checkpoint = ModelCheckpoint(filepath, monitor='val_acc', verbose=1, save_best_only=True, mode='max', save_weights_only=True)
callbacks_list = [checkpoint]
```

Fig 9 CNN + GRU

LSTM, an overhauled type of RNN that catches long haul connections and is ideal for succession

expectation. Applied to time series investigation, machine interpretation, and voice acknowledgment. LSTM memory cells have input, neglect, and result entryways, in contrast to RNNs. Data is specifically held or disposed of by these doors. LSTMs might be combined with CNNs for picture and video investigation in light of the fact that to their novel potential. LSTM [65] is possible utilized in the review to display relevant data across broadened arrangements and decipher hate speech's complex etymological examples.[121]

```
embed_dim = 128 # dimension of the word embedding vector for each word in a sequence
lstm_out = 100 # no. of lstm layers
lstm_model = Sequential()
lstm_model.add(Embedding(num_words, embed_dim, input_length=X_train.shape[1]))
lstm_model.add(LSTM(lstm_out, recurrent_dropout=0.2))
lstm_model.add(LSTM(lstm_out, recurrent_dropout=0.2))
lstm_model.add(Dense(1, kernel_regularizer=regularizers.L2(0.001), activation='sigmoid'))
lstm_model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy', 'f1_score', 'recall'])
print(lstm_model.summary())
```

Fig 10 LSTM

LSTM + GRU (Long Short-Term Memory with Gated Recurrent Unit), LSTM's long-term conditions and GRUs' computational effectiveness are joined in this hybrid architecture. LSTMs catch far off context oriented data well, and GRUs train quicker and handle transient conditions. For its decent way to deal with demonstrating short and long haul successive conditions, LSTM + GRU might be considered for the venture to all the more likely figure out hate speech articulations.

LSTM + GRU

```
model = Sequential()
model.add(Embedding(num_words, embed_dim, input_length=X_train.shape[1]))
model.add(LSTM(lstm_out, recurrent_dropout=0.2, return_sequences=True))
model.add(GRU(gru_out, recurrent_dropout=0.2, return_sequences=False))
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy', 'f1_score', 'recall'])
print(model.summary())
filepath = 'weights_lstm_gru.h5'
checkpoint = ModelCheckpoint(filepath, monitor='val_acc', verbose=1, save_best_only=True, mode='max', save_weights_only=True)
callbacks_list = [checkpoint]
```

Fig 11 LSTM + GRU

LSTM + BiGRU (Long Short-Term Memory + Bidirectional Gated Recurrent Unit), This hybrid architecture utilizes LSTM sequential learning and Gated Recurrent Unit bidirectional processing. LSTM catches long-range conditions, while BiGRU processes data forward and in reverse. We picked this mix since it handles consecutive information with refined worldly examples well, expanding the model's hate speech recognition.

LSTM + BiGRU

```
from tensorflow.keras.layers import Input, LSTM, Dense, Concatenate, Activation, BatchNormalization, Dropout, Bidirectional

model = Sequential()
model.add(LSTM(input_shape=(input_length, input_dim), return_sequences=True))
model.add(Bidirectional(LSTM(input_shape=(input_length, input_dim), return_sequences=True)))
model.add(Dense(input_dim))
model.add(Activation('softmax'))
model.compile(optimizer='adam', metrics=['accuracy', 'f1_score', 'precision', 'recall'])
print(model.summary())
filepath = 'weights_lstm_bi-gru.h5'
checkpoint = ModelCheckpoint(filepath, monitor='val_acc', verbose=1, save_best_only=True, mode='max', save_weights_only=True)
callbacks_list = [checkpoint]
```

Fig 12 LSTM + BiGRU

Naïve Bayes classifiers utilize Bayes' Hypothesis to arrange. This group of calculations all offer the possibility that each sets of attributes being classed is autonomous. To start with, consider a dataset. The basic and viable Naïve Bayes classifier empowers speedy formation of ML models with expectation abilities. The Naïve Bayes classifier's name alludes to its working on suspicions. The classifier expects that perception attributes are restrictively free given the class name. "Bayes" alludes to Reverend Thomas Bayes. For hate speech identification, Naïve Bayes' straightforwardness and quick preparation time can act as a pattern model for further developed calculations.

Naive Bayes

```
from sklearn.naive_bayes import MultinomialNB
clf = MultinomialNB()
clf.fit(X_train, y_train)
y_pred = clf.predict(X_test)

lr_acc = accuracy_score(y_test, y_pred)
lr_prec = precision_score(y_test, y_pred, average='weighted')
lr_rec = recall_score(y_test, y_pred, average='weighted')
lr_f1 = f1_score(y_test, y_pred, average='weighted')

storeResults('Naive Bayes', lr_acc, lr_prec, lr_rec, lr_f1)
```

Fig 13 Naïve bayes

Random Forest is a typical supervised ML strategy. It can address ML order and relapse issues. Ensemble learning utilizes a few classifiers to deal with convoluted issues and improve model execution. As the name says, "Random Forest is a classifier that contains various decision trees on different subsets of the given dataset and takes the normal to work on the prescient accuracy of that dataset." Rather than utilizing one decision tree, the random forest conjectures a definitive result in light of the greater part votes of each tree.[122]

Random Forest

```
from sklearn.ensemble import RandomForestClassifier
RandomForest = RandomForestClassifier(n_estimators=10, random_state=0)

RandomForest.fit(X_train, y_train)
y_pred = RandomForest.predict(X_test)

rf_acc = accuracy_score(y_test, y_pred)
rf_prec = precision_score(y_test, y_pred, average='weighted')
rf_rec = recall_score(y_test, y_pred, average='weighted')
rf_f1 = f1_score(y_test, y_pred, average='weighted')

storeResults('Random Forest', rf_acc, rf_prec, rf_rec, rf_f1)
```

Fig 14 Random forest

LinearSVC (Linear Support Vector Classifier): LinearSVC is a Support Vector Machine (SVM) strategy that spotlights on linear classification. SVMs are great in making hyperplanes that partition different

classes in a high-layered space. LinearSVC can be valuable in hate speech recognizable proof in light of its ability to deal with non-direct choice cutoff points and precisely arrange occurrences of can't stand discourse.

LinearSVC

```
from sklearn.svm import LinearSVC
svm = LinearSVC()
svm.fit(X_train, y_train)
y_pred = svm.predict(X_test)

svm_acc = accuracy_score(y_test, y_pred)
svm_prec = precision_score(y_test, y_pred, average='weighted')
svm_rec = recall_score(y_test, y_pred, average='weighted')
svm_f1 = f1_score(y_test, y_pred, average='weighted')

storeResults('LinearSVC', svm_acc, svm_prec, svm_rec, svm_f1)
```

Fig 15 LinearSVC

RF + SVM + NB (Random Forest + Support Vector Machine + Naive Bayes), this group strategy joins the benefits of the Random Forest, Support Vector Machine (SVM), and Naive Bayes (NB) calculations. Random Forest gives strength through decision tree ensembles, SVM succeeds at building compelling hyperplanes, and Gullible Bayes offers probabilistic classification. This gathering is probably utilized as a result of its ability to gather a large number of the information, which further develops generally speaking hate speech detection accuracy.

RF + SVM + NB

```
from sklearn.ensemble import VotingClassifier
from sklearn.svm import SVC
from sklearn.ensemble import RandomForestClassifier

estimator = []
estimator.append(('SVM', SVC(probability=True)))
estimator.append(('RFC', RandomForestClassifier()))
estimator.append(('NB', MultinomialNB()))
vot_hard = VotingClassifier(estimators = estimator, voting = 'soft')
vot_hard.fit(X_train, y_train)
y_pred = vot_hard.predict(X_test)

vot_acc = accuracy_score(y_test, y_pred)
vot_prec = precision_score(y_test, y_pred, average='weighted')
vot_rec = recall_score(y_test, y_pred, average='weighted')
vot_f1 = f1_score(y_test, y_pred, average='weighted')

storeResults('RF + SVM + NB', vot_acc, vot_prec, vot_rec, vot_f1)
```

Fig 16 RF+SVM+NB

Stacking Classifier, Stacking is an ensemble learning system in which various models are prepared to foresee a similar result, and afterward a meta-model is prepared to make expectations in view of the singular models' outcomes. With regards to hate speech detection, a Stacking Classifier is most often used to join the capacities of many base models, bringing about a stronger and exact by and large hate speech detection system.

```
from sklearn.ensemble import RandomForestClassifier
from sklearn.neural_network import MLPClassifier
from lightgbm import LGBMClassifier
from xgboost import XGBClassifier
from sklearn.ensemble import StackingClassifier

estimators = [('rf', RandomForestClassifier(n_estimators=100)), ('mlp', MLPClassifier(random_state=1, max_iter=100))]
clf = StackingClassifier(estimators=estimators, final_estimator=MLPClassifier())
clf.fit(X_train, y_train)

y_pred = clf.predict(X_test)

stat_acc = accuracy_score(y_test, y_pred)
stat_prec = precision_score(y_test, y_pred, average='weighted')
stat_rec = recall_score(y_test, y_pred, average='weighted')
stat_f1 = f1_score(y_test, y_pred, average='weighted')

storeResults('Stacking Classifier', stat_acc, stat_prec, stat_rec, stat_f1)
```

Fig 17 Stacking classifier

4. EXPERIMENTAL RESULTS

Precision: Precision estimates the level of positive cases or tests precisely sorted. Precision is determined utilizing the recipe:

$$\text{Precision} = \frac{\text{True positives}}{(\text{True positives} + \text{False positives})} = \frac{TP}{(TP + FP)}$$

$$\text{Precision} = \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$

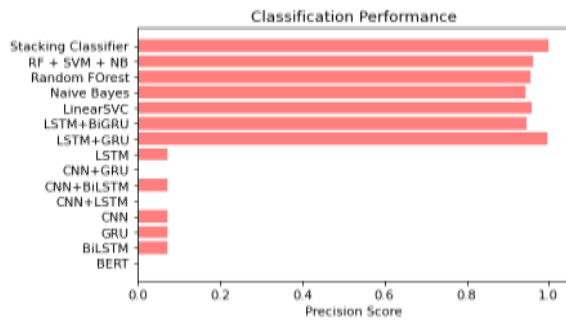


Fig 18 Precision comparison graph

Recall: Machine learning recall assesses a model's ability to perceive all significant examples of a class. It shows a model's culmination in catching occasions of a class by contrasting accurately anticipated positive perceptions with complete positives.

$$\text{Recall} = \frac{TP}{TP + FN}$$

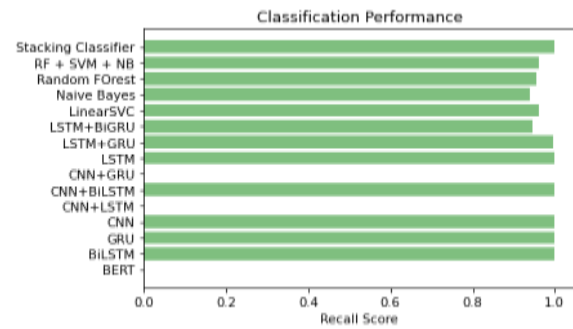


Fig 19 Recall comparison graph

Accuracy: A test's accuracy is its ability to recognize debilitated from sound cases. To quantify test accuracy, figure the small part of true positive and true negative in completely broke down cases. Numerically, this is:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN}$$

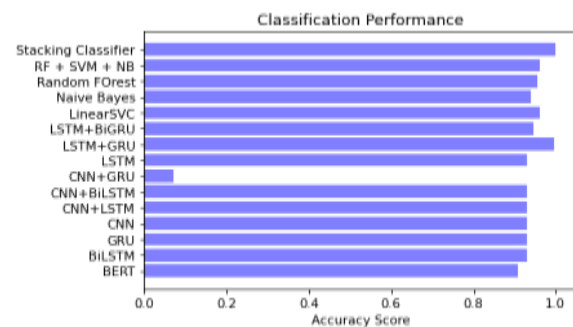


Fig 20 Accuracy graph

F1 Score: Machine learning model accuracy is estimated by F1 score. Consolidating model precision

and recall scores. The accuracy measurement estimates how frequently a model anticipated accurately all through the dataset.

$$F1 \text{ Score} = 2 * \frac{\text{Recall} \times \text{Precision}}{\text{Recall} + \text{Precision}} * 100$$

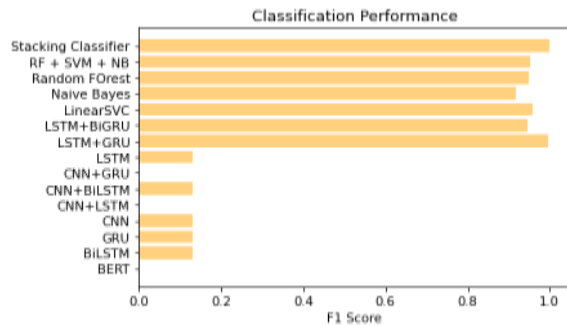


Fig 21 F1Score

ML Model	Accuracy	Precision	Recall	F1 - score
BKRT	0.908	0.000	0.000	0.000
BiLSTM	0.930	0.070	1.000	0.130
GRU	0.930	0.070	1.000	0.130
CNN	0.930	0.070	1.000	0.130
CNN+LSTM	0.930	0.000	0.000	0.000
CNN+BiLSTM	0.930	0.070	1.000	0.130
CNN+GRU	0.970	0.000	0.000	0.000
LSTM	0.630	0.070	1.000	0.130
LSTM+GRU	0.966	0.966	0.966	0.966
LSTM+BiGRU	0.942	0.942	0.942	0.942
LinearSVC	0.961	0.960	0.961	0.957
Naive Bayes	0.939	0.948	0.939	0.917
Random Forest	0.956	0.954	0.956	0.949
RF + SVM + NB	0.960	0.960	0.960	0.953
Stacking Classifier	1.000	1.000	1.000	1.000

Fig 22 Performance Evaluation



Fig 23 Home page

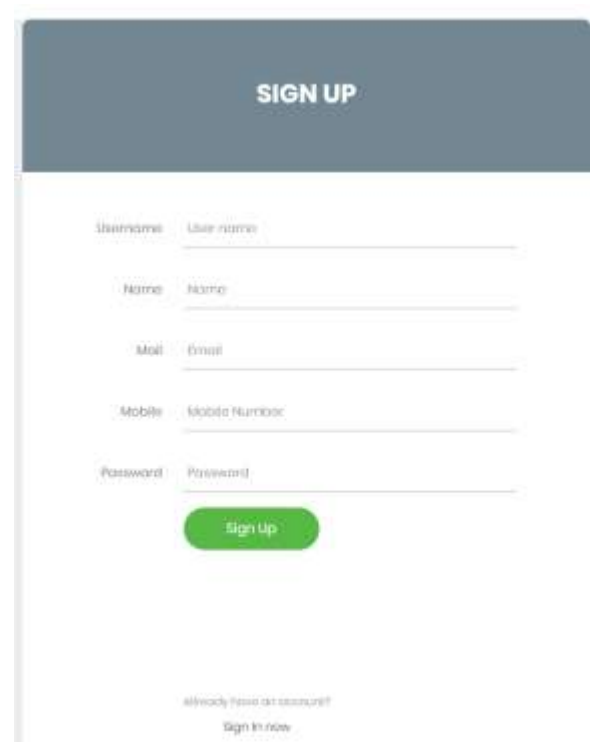


Fig 24 Signin page



Fig 25 Login page



Fig 26 User input



Fig 27 Prediction result

5. CONCLUSION

Web-based entertainment clients who persevere through internet based misuse benefit most from the undertaking. The work makes the web more secure and more sure by perceiving and diminishing hate speech. Diminished provocation makes a more comprehensive internet based local area. The drive detects hate speech for controllers and stage administrators [32]. The innovation upholds advanced disdain discourse guidelines. Administrative associations can make proactive moves to safeguard a sound web environment. A strong ensemble stacking classifier accomplishes 100 percent accuracy. Front-end testing with validation showed the model's capacity to perceive and address Twitter hate speech. Ensemble approaches increment estimating accuracy by consolidating various models. Flask with SQLite for client information exchange and signin guarantee security and confirmation. This safeguards client protection and makes the hate speech detection system dependable.

6. FUTURE SCOPE

Creating calculations that can detect hate speech in dialects other than English would guarantee an overall

impact and advance consideration. Disdain discourse recognition strategies will be refreshed to answer web language patterns and ongoing learning and change. To recognize blameless expressions and hate speech, complex NLP approaches [64, 8, 87] including opinion investigation and mockery acknowledgment could work on the model's context oriented appreciation. Future advances might engage clients with adjustable channels and content decisions to alter disdain discourse identification force, making the experience more easy to use and versatile.

REFERENCES

- [1] Z. Zhang and L. Luo, "Hate speech detection: A solved problem? The challenging case of long tail on Twitter," *Semantic Web*, vol. 10, no. 5, pp. 925–945, Sep. 2019.
- [2] L. Silva, M. Mondal, D. Correa, F. Benevenuto, and I. Weber, "Analyzing the targets of hate in online social media," in *Proc. Int. AAAI Conf. Web Social Media*, 2016, vol. 10, no. 1, pp. 1–4.
- [3] R. Magu, K. Joshi, and J. Luo, "Detecting the hate code on social media," in *Proc. Int. AAAI Conf. Web Social Media*, 2017, vol. 11, no. 1, pp. 608–611.
- [4] M. Mondal, L. A. Silva, D. Correa, and F. Benevenuto, "Characterizing usage of explicit hate expressions in social media," *New Rev. Hypermedia Multimedia*, vol. 24, no. 2, pp. 110–130, Apr. 2018.
- [5] A. M. Founta, D. Chatzakou, N. Kourtellis, J. Blackburn, A. Vakali, and I. Leontiadis, "A unified deep learning architecture for abuse detection," in *Proc. 10th ACM Conf. Web Sci.*, Jun. 2019, pp. 105–114.

- [6] A. Arango, J. Pérez, and B. Poblete, “Hate speech detection is not as easy as you may think: A closer look at model validation (extended version),” *Inf. Syst.*, vol. 105, Mar. 2022, Art. no. 101584.
- [7] A. Schmidt and M. Wiegand, “A survey on hate speech detection using natural language processing,” in *Proc. 5th Int. Workshop Natural Lang. Process. Social Media*, 2017, pp. 1–10.
- [8] J. Langham and K. Gosha, “The classification of aggressive dialogue in social media platforms,” in *Proc. ACM SIGMIS Conf. Comput. People Res.*, Jun. 2018, pp. 60–63.
- [9] P. Fortuna and S. Nunes, “A survey on automatic detection of hate speech in text,” *ACM Comput. Surv.*, vol. 51, no. 4, pp. 1–30, 2018.
- [10] W. Dorris, R. Hu, N. Vishwamitra, F. Luo, and M. Costello, “Towards automatic detection and explanation of hate speech and offensive language,” in *Proc. 6th Int. Workshop Secur. Privacy Anal.*, Mar. 2020, pp. 23–29.
- [11] A. Alrehili, “Automatic hate speech detection on social media: A brief survey,” in *Proc. IEEE/ACS 16th Int. Conf. Comput. Syst. Appl. (AICCSA)*, Nov. 2019, pp. 1–6.
- [12] S. Modi, “AHTDT—Automatic hate text detection techniques in social media,” in *Proc. Int. Conf. Circuits Syst. Digit. Enterprise Technol. (ICCSDET)*, Dec. 2018, pp. 1–3.
- [13] F. E. Ayo, O. Folorunso, F. T. Ibharalu, and I. A. Osinuga, “Machine learning techniques for hate speech classification of Twitter data: Stateofthe-art, future challenges and research directions,” *Comput. Sci. Rev.*, vol. 38, Nov. 2020, Art. no. 100311.
- [14] E. Shushkevich and J. Cardiff, “Automatic misogyny detection in social media: A survey,” *Computación Y Sistemas*, vol. 23, no. 4, pp. 1159–1164, Dec. 2019.
- [15] F. Poletto, V. Basile, M. Sanguinetti, C. Bosco, and V. Patti, “Resources and benchmark corpora for hate speech detection: A systematic review,” *Lang. Resour. Eval.*, vol. 55, pp. 477–523, Jun. 2020.
- [16] T. X. Moy, M. Raheem, and R. Logeswaran, “Hate speech detection in English and non-English languages: A review of techniques and challenges,” *Webology*, vol. 18, no. 5, pp. 929–938, Oct. 2021, doi: 10.14704/WEB/V18SI05/WEB18272.
- [17] W. Yin and A. Zubiaga, “Towards generalisable hate speech detection: A review on obstacles and solutions,” *PeerJ Comput. Sci.*, vol. 7, pp. 1–38, Jun. 2021, doi: 10.7717/PEERJ-CS.598.
- [18] N. S. Mullah and W. M. N. W. Zainon, “Advances in machine learning algorithms for hate speech detection in social media: A review,” *IEEE Access*, vol. 9, pp. 88364–88376, 2021, doi: 10.1109/ACCESS.2021.3089515.
- [19] O. Istaiteh, R. Al-Omoush, and S. Tedmori, “Racist and sexist hate speech detection: Literature review,” in *Proc. Int. Conf. Intell. Data Sci. Technol. Appl. (IDSTA)*, Oct. 2020, pp. 95–99, doi: 10.1109/IDSTA50958.2020.9264052.
- [20] R. Rini, E. Utami, and A. D. Hartanto, “Systematic literature review of hate speech detection with text mining,” in *Proc. 2nd Int. Conf. Cybern. Intell. Syst. (ICORIS)*, Oct. 2020, pp. 1–6, doi: 10.1109/ICORIS50180.2020.9320755.

- [21] S. MacAvaney, H. R. Yao, E. Yang, K. Russell, N. Goharian, and O. Frieder, "Hate speech detection: Challenges and solutions," *PLoS ONE*, vol. 14, no. 8, pp. 1–16, 2019, doi: 10.1371/journal.pone.0221152.
- [22] F. Alkomah and X. Ma, "A literature review of textual hate speech detection methods and datasets," *Information*, vol. 13, no. 6, p. 122, 2022, doi: 10.3390/info13060273.
- [23] B. Kitchenham and S. Charters, "Guidelines for performing systematic literature reviews in software engineering," *Tech. Rep.*, 2007.
- [24] D. Moher, "Preferred reporting items for systematic reviews and metaanalyses: The PRISMA statement," *Ann. Internal Med.*, vol. 151, no. 4, p. 264, Aug. 2009.
- [25] Z. Al-Makhadmeh and A. Tolba, "Automatic hate speech detection using killer natural language processing optimizing ensemble deep learning approach," *Computing*, vol. 102, no. 2, pp. 501–522, Feb. 2020.
- [26] M. Suhaidi, R. A. Kadir, and S. Tiun, "A review of feature extraction methods on machine learning," *J. Inf. Syst. Technol. Manag.*, vol. 6, no. 22, pp. 51–59, Sep. 2021, doi: 10.35631/jistm.622005.
- [27] M. D. Oskouei and S. N. Razavi, "An ensemble feature selection method to detect web spam," *Asia–Pacific J. Inf. Technol. Multimedia*, vol. 7, no. 2, pp. 99–113, Dec. 2018, doi: 10.17576/apjitm-2018-0702-08.
- [28] P. Badjatiya, S. Gupta, M. Gupta, and V. Varma, "Deep learning for hate speech detection in tweets," in *Proc. 26th Int. Conf. World Wide Web Companion*, 2017, pp. 759–760.
- [29] A. F. Naswir, L. Q. Zakaria, and S. Saad, "The effectiveness of URL features on phishing emails classification using machine learning approach," *Asia–Pacific J. Inf. Technol. Multimedia J. Teknol. Mklm. Dan Multimedia Asia-Pasifik*, vol. 7, no. 2, pp. 61–69, 2022, doi: 10.17576/apjitm-2022-1102-04.
- [30] Y. Yadav, P. Bajaj, R. K. Gupta, and R. Sinha, "A comparative study of deep learning methods for hate speech and offensive language detection in textual data," in *Proc. IEEE 18th India Council Int. Conf. (INDICON)*, Dec. 2021, pp. 1–6, doi: 10.1109/INDICON52576.2021.9691704.
- [31] M. Li, S. Liao, E. Okpala, M. Tong, M. Costello, L. Cheng, H. Hu, and F. Luo, "COVID-HateBERT: A pre-trained language model for COVID-19 related hate speech detection," in *Proc. 20th IEEE Int. Conf. Mach. Learn. Appl. (ICMLA)*, Dec. 2021, pp. 233–238, doi: 10.1109/ICMLA52953.2021.00043.
- [32] M. K. A. Aljero and N. Dimililer, "Genetic programming approach to detect hate speech in social media," *IEEE Access*, vol. 9, pp. 115115–115125, 2021, doi: 10.1109/ACCESS.2021.3104535.
- [33] Y. Saini, V. Bachchas, Y. Kumar, and S. Kumar, "Abusive text examination using latent Dirichlet allocation, self organizing maps and K means clustering," in *Proc. 4th Int. Conf. Intell. Comput. Control Syst. (ICICCS)*, May 2020, pp. 1233–1238.
- [34] M. Moh, T.-S. Moh, and B. Khieu, "No 'Love' lost: Defending hate speech detection models against

adversaries,” in Proc. 14th Int. Conf. Ubiquitous Inf. Manage. Commun. (IMCOM), Jan. 2020, pp. 1–6.

[35] Z. Zhang, D. Robinson, and J. Tepper, “Detecting hate speech on Twitter using a convolution-GRU based deep neural network,” in Proc. Eur. Semantic Web Conf. Cham, Switzerland: Springer, 2018, pp. 745–760.

[36] D. Robinson, Z. Zhang, and J. Tepper, “Hate speech detection on Twitter: Feature engineering vs feature selection,” in Proc. Eur. Semantic Web Conf. Cham, Switzerland: Springer, 2018, pp. 46–49.

[37] R. Cao, R. K.-W. Lee, and T.-A. Hoang, “DeepHate: Hate speech detection via multi-faceted text representations,” in Proc. 12th ACM Conf. Web Sci., Jul. 2020, pp. 11–20.

[38] T. Chakrabarty, K. Gupta, and S. Muresan, “Pay ‘attention’ to your context when classifying abusive language,” in Proc. 3rd Workshop Abusive Lang. Online, 2019, pp. 70–79.

[39] G. K. Pitsilis, H. Ramampiaro, and H. Langseth, “Effective hate-speech detection in Twitter data using recurrent neural networks,” Appl. Intell., vol. 48, no. 12, pp. 4730–4742, 2018.

[40] W. Alorainy, P. Burnap, H. Liu, A. Javed, and M. L. Williams, “Suspended accounts: A source of tweets with disgust and anger emotions for augmenting hate speech data sample,” in Proc. Int. Conf. Mach. Learn. Cybern. (ICMLC), Jul. 2018, pp. 581–586.

[41] W. Alorainy, P. Burnap, H. Liu, and M. L. Williams, “‘The enemy among us’: Detecting cyber hate speech with threats-based othering language

embeddings,” ACM Trans. Web, vol. 13, no. 3, pp. 1–26, Aug. 2019.

[42] H. Liu, P. Burnap, W. Alorainy, and M. L. Williams, “Fuzzy multi-task learning for hate speech type detection,” in Proc. World Wide Web Conf., May 2019, pp. 3006–3012.

[43] P. Burnap and M. L. Williams, “Us and them: Identifying cyber hate on Twitter across multiple protected characteristics,” EPJ Data Sci., vol. 5, pp. 1–15, Oct. 2016.

[44] S. Malmasi and M. Zampieri, “Detecting hate speech in social media,” 2017, arXiv:1712.06427.

[45] S. Malmasi and M. Zampieri, “Challenges in discriminating profanity from hate speech,” J. Experim. Theor. Artif. Intell., vol. 30, no. 2, pp. 187–202, Mar. 2018.

[46] M. Anzovino, E. Fersini, and P. Rosso, “Automatic detection and classification of misogynistic language on Twitter,” in Proc. Int. Conf. Appl. Natural Lang. Inf. Syst. Cham, Switzerland: Springer, 2018, pp. 57–64.

[47] S. Frenda, B. Ghanem, M. Montes-y-Gómez, and P. Rosso, “Online hate speech against women: Automatic detection of misogyny and sexism on Twitter,” J. Intell. Fuzzy Syst., vol. 36, no. 5, pp. 4743–4752, May 2019.

[48] M. A. Bashar, R. Nayak, N. Suzor, and B. Weir, “Misogynistic tweet detection: Modelling CNN with small datasets,” in Proc. Australas. Conf. Data Mining. Cham, Switzerland: Springer, 2018, pp. 3–16.

- [49] A. Arango, J. Pérez, and B. Poblete, “Hate speech detection is not as easy as you may think: A closer look at model validation,” in Proc. 42nd Int. ACM SIGIR Conf. Res. Develop. Inf. Retr., 2019, pp. 45–54.
- [50] J. S. Meyer and B. Gambäck, “A platform agnostic dual-strand hate speech detector,” in Proc. 3rd Workshop Abusive Lang. Online, 2019, pp. 146–156.
- [51] B. Vidgen and T. Yasseri, “Detecting weak and strong islamophobic hate speech on social media,” J. Inf. Technol. Politics, vol. 17, no. 1, pp. 66–78, Jan. 2020.
- [52] R. Martins, M. Gomes, J. J. Almeida, P. Novais, and P. Henriques, “Hate speech classification in social media using emotional analysis,” in Proc. 7th Brazilian Conf. Intell. Syst. (BRACIS), Oct. 2018, pp. 61–66.
- [53] M. Sajjad, F. Zulifqar, M. U. G. Khan, and M. Azeem, “Hate speech detection using fusion approach,” in Proc. Int. Conf. Appl. Eng. Math. (ICAEM), Aug. 2019, pp. 251–255.
- [54] H. Watanabe, M. Bouazizi, and T. Ohtsuki, “Hate speech on Twitter: A pragmatic approach to collect hateful and offensive expressions and perform hate speech detection,” IEEE Access, vol. 6, pp. 13825–13835, 2018.
- [55] M. Sap, D. Card, S. Gabriel, Y. Choi, and N. A. Smith, “The risk of racial bias in hate speech detection,” in Proc. 57th Annu. Meeting Assoc. Comput. Linguistics, 2019, pp. 1668–1678.
- [56] M. Mozafari, R. Farahbakhsh, and N. Crespi, “A bert-based transfer learning approach for hate speech detection in online social media,” in Proc. Int. Conf. Complex Netw. Their Appl. Cham, Switzerland: Springer, 2019, pp. 928–940.
- [57] G. Rizos, K. Hemker, and B. Schuller, “Augment to prevent: Short-text data augmentation in deep learning for hate-speech classification,” in Proc. 28th ACM Int. Conf. Inf. Knowl. Manag., Nov. 2019, pp. 991–1000.
- [58] A. G. D’Sa, I. Illina, and D. Fohr, “BERT and fastText embeddings for automatic detection of toxic speech,” in Proc. Int. Multi-Conf. Org. Knowl. Adv. Technologie (OCTA), Feb. 2020, pp. 1–5.
- [59] G. Koushik, K. Rajeswari, and S. K. Muthusamy, “Automated hate speech detection on Twitter,” in Proc. 5th Int. Conf. Comput., Commun., Control Autom. (ICCUBE), Sep. 2019, pp. 1–4.
- [60] L. Jiang and Y. Suzuki, “Detecting hate speech from tweets for sentiment analysis,” in Proc. 6th Int. Conf. Syst. Informat. (ICSAI), Nov. 2019, pp. 671–676.
- [61] B. Mathew, N. Kumar, P. Goyal, and A. Mukherjee, “Interaction dynamics between hate and counter users on Twitter,” in Proc. 7th ACM IKDD CoDS 25th COMAD, Jan. 2020, pp. 116–124.
- [62] H. Liu, P. Burnap, W. Alorainy, and M. L. Williams, “A fuzzy approach to text classification with two-stage training for ambiguous instances,” IEEE Trans. Comput. Soc. Syst., vol. 6, no. 2, pp. 227–240, Apr. 2019.

- [63] R. Hu, W. Dorris, N. Vishwamitra, F. Luo, and M. Costello, “On the impact of word representation in hate speech and offensive language detection and explanation,” in Proc. 10th ACM Conf. Data Appl. Secur. Privacy, Mar. 2020, pp. 171–173.
- [64] J. Qian, M. ElSherief, E. M. Belding, and W. Yang Wang, “Leveraging intra-user and inter-user representation learning for automated hate speech detection,” 2018, arXiv:1804.03124.
- [65] A. Bisht, “Detection of hate speech and offensive language in Twitter data using LSTM model,” in Recent Trends in Image and Signal Processing in Computer Vision. Berlin, Germany: Springer, 2020, pp. 243–264.
- [66] R. Ahluwalia, E. Shcherbinina, E. Callow, A. C. Nascimento, and M. De Cock, “Detecting misogynous tweets,” in Proc. IberEval SEPLN, 2018, pp. 242–248.
- [67] D. Nozza, C. Volpetti, and E. Fersini, “Unintended bias in misogyny detection,” in Proc. IEEE/WIC/ACM Int. Conf. Web Intell., Oct. 2019, pp. 149–155.
- [68] T. Gröndahl, L. Pajola, M. Juuti, M. Conti, and N. Asokan, “All you need is ‘Love’: Evading hate speech detection,” in Proc. 11th ACM Workshop Artif. Intell. Secur., Jan. 2018, pp. 2–12.
- [69] R. Oak, “Poster: Adversarial examples for hate speech classifiers,” in Proc. ACM SIGSAC Conf. Comput. Commun. Secur., Nov. 2019, pp. 2621–2623.
- [70] T. Wullach, A. Adler, and E. Minkov, “Towards hate speech detection at large via deep generative modeling,” IEEE Internet Comput., vol. 25, no. 2, pp. 48–57, Mar. 2021.
- [71] R. I. Rasel, N. Sultana, S. Akhter, and P. Meesad, “Detection of cyberaggressive comments on social media networks: A machine learning and text mining approach,” in Proc. 2nd Int. Conf. Natural Lang. Process. Inf. Retr., Sep. 2018, pp. 37–41.
- [72] B. Gupta, N. Goel, D. Jain, and N. Gupta, “A novel IN-Gram technique for improving the hate speech detection for larger datasets,” in Micro-Electronics and Telecommunication Engineering. Berlin, Germany: Springer, 2020, pp. 611–620.
- [73] M. M. Al-Ani, N. Omar, and A. A. Nafea, “A hybrid method of long short-term memory and auto-encoder architectures for sarcasm detection,” J. Comput. Sci., vol. 17, no. 11, pp. 1093–1098, Nov. 2021, doi: 10.3844/JCSSP.2021.1093.1098.
- [74] R. T. Mutanga and N. Naicker, “Hate speech detection in Twitter using transformer methods,” Int. J. Adv. Comput. Sci. Appl., vol. 11, no. 9, pp. 1–7, 2020.
- [75] R. Joshua. (May 2009). Tweepy Documentation. [Online]. Available: <http://tweepy.readthedocs.io/en/v3>
- [76] W. Alorainy, P. Burnap, H. Liu, and M. L. Williams, “The enemy among us’: Detecting cyber hate speech with threats-based othering language embeddings,” ACM Trans. Web, vol. 13, no. 3, pp. 1–26, Aug. 2019.
- [77] Z. Waseem and D. Hovy, “Hateful symbols or hateful people? Predictive features for hate speech

detection on Twitter,” in Proc. NAACL Student Res. Workshop, 2016, pp. 88–93.

[78] V. Basile, C. Bosco, E. Fersini, D. Nozza, V. Patti, F. M. R. Pardo, P. Rosso, and M. Sanguinetti, “SemEval-2019 task 5: Multilingual detection of hate speech against immigrants and women in Twitter,” in Proc. 13th Int. Workshop Semantic Eval., 2019, pp. 54–63.

[79] E. Fersini, P. Rosso, and M. Anzovino, “Overview of the task on automatic misogyny detection at IberEval 2018,” in Proc. IberEval SEPLN, vol. 2150, Sep. 2018, pp. 214–228.

[80] E. Fersini, D. Nozza, and P. Rosso, “Overview of the evalita 2018 task on automatic misogyny detection (AMI),” in Proc. EVALITA Eval. NLP Speech Tools Italian, vol. 12, 2018, p. 59.

[81] M. Mondal, L. A. Silva, and F. Benevenuto, “A measurement study of hate speech in social media,” in Proc. 28th ACM Conf. Hypertext Social Media, Jul. 2017, pp. 85–94.

[82] A. Founta, C. Djouvas, D. Chatzakou, I. Leontiadis, J. Blackburn, G. Stringhini, A. Vakali, M. Sirivianos, and N. Kourtellis, “Large scale crowdsourcing and characterization of Twitter abusive behavior,” in Proc. Int. AAAI Conf. Web Social Media, 2018, vol. 12, no. 1, pp. 1–10.

[83] S. Abro, S. Shaikh, Z. Hussain, Z. Ali, S. Khan, and G. Mujtaba, “Automatic hate speech detection using machine learning: A comparative study,” Int. J. Adv. Comput. Sci. Appl., vol. 11, no. 8, pp. 1–8, 2020.

[84] A. H. Wani, N. S. Molvi, and S. I. Ashraf, “Detection of hate and offensive speech in text,” in Proc. Int. Conf. Intell. Hum. Comput. Interact. Cham, Switzerland: Springer, 2019, pp. 87–93.

[85] K. J. Madukwe and X. Gao, “The thin line between hate and profanity,” in Proc. Australas. Joint Conf. Artif. Intell. Cham, Switzerland: Springer, 2019, pp. 344–356.

[86] M. ElSherief, V. Kulkarni, D. Nguyen, W. Y. Wang, and E. Belding, “Hate lingo: A target-based linguistic analysis of hate speech in social media,” in Proc. Int. AAAI Conf. Web Social Media, 2018, vol. 12, no. 1, pp. 1–10.

[87] Z. Waseem, “Are you a racist or am i seeing things? Annotator influence on hate speech detection on Twitter,” in Proc. 1st Workshop NLP Comput. Social Sci., 2016, pp. 138–142.

[88] S. Agrawal and A. Awekar, “Deep learning for detecting cyberbullying across multiple social media platforms,” in Proc. Eur. Conf. Inf. Retr. Cham, Switzerland: Springer, 2018, pp. 141–153.

[89] P. Chiril, E. W. Pamungkas, F. Benamara, V. Moriceau, and V. Patti, “Emotionally informed hate speech detection: A multi-target perspective,” Cognit. Comput., vol. 14, no. 1, pp. 322–352, Jan. 2022, doi: 10.1007/s12559-021-09862-5.

[90] N. Rai, P. Meena, and C. Agrawal, “Improving the hate speech analysis through dimensionality reduction approach,” in Proc. 6th Int. Conf. Adv. Comput. Commun. Syst. (ICACCS), Mar. 2020, pp. 321–325.

- [91] R. Ahluwalia, H. Soni, E. Callow, A. Nascimento, and M. De Cock, “Detecting hate speech against women in English tweets,” *EVALITA Eval. NLP Speech Tools Italian*, vol. 12, p. 194, Dec. 2018.
- [92] J. Dhillon, V. Gupta, R. Govil, B. Varshney, and A. Sinha, “Crowdsourcing of hate speech for detecting abusive behavior on social media,” in *Proc. Int. Conf. Signal Process. Commun. (ICSC)*, Mar. 2019, pp. 41–46.
- [93] T. Davidson, D. Warmley, M. Macy, and I. Weber, “Automated hate speech detection and the problem of offensive language,” in *Proc. Int. AAAI Conf. Web Social Media*, 2017, vol. 11, no. 1, pp. 512–515.
- [94] M. Ribeiro, P. Calais, Y. Santos, V. Almeida, and W. Meira Jr., “Characterizing and detecting hateful users on Twitter,” in *Proc. Int. AAAI Conf. Web Social Media*, 2018, vol. 12, no. 1, pp. 1–10.
- [95] Y. Senarath and H. Purohit, “Evaluating semantic feature representations to efficiently detect hate intent on social media,” in *Proc. IEEE 14th Int. Conf. Semantic Comput. (ICSC)*, Feb. 2020, pp. 199–202.
- [96] H. Rathpisey and T. B. Adjii, “Handling imbalance issue in hate speech classification using sampling-based methods,” in *Proc. 5th Int. Conf. Sci. Inf. Technol. (ICSITech)*, Oct. 2019, pp. 193–198.
- [97] M. A. Bashar and R. Nayak, “Active learning for effectively fine-tuning transfer learning to downstream task,” *ACM Trans. Intell. Syst. Technol.*, vol. 12, no. 2, pp. 1–24, Apr. 2021, doi: 10.1145/3446343.
- [98] B. Vidgen, “Detecting East Asian prejudice on social media,” in *Proc. Social Inf. Netw.*, 2020, pp. 162–172, doi: 10.18653/v1/2020.alw-1.19.
- [99] F. Barbieri, J. Camacho-Collados, L. Neves, and L. Espinosa-Anke, “TWEETEVAL: Unified benchmark and comparative evaluation for tweet classification,” in *Proc. Find. Assoc. Comput. Linguist. Find. ACL (EMNLP)*, 2020, pp. 1644–1650, doi: 10.18653/v1/2020.findingsemnlp.148.
- [100] S. M. Al-Ghuribi, S. A. M. Noah, and S. Tiun, “Unsupervised semantic approach of aspect-based sentiment analysis for large-scale user reviews,” *IEEE Access*, vol. 8, pp. 218592–218613, 2020, doi: 10.1109/ACCESS.2020.3042312.
- [101] D. G. Kyrollos and J. R. Green, “MetaHate: A meta-model for hate speech detection,” in *Proc. IEEE Int. Conf. Big Data (Big Data)*, Dec. 2021, pp. 2496–2502, doi: 10.1109/BigData52589.2021.9672023.
- [102] K. J. Madukwe, X. Gao, and B. Xue, “Token replacement-based data augmentation methods for hate speech detection,” *World Wide Web*, vol. 25, no. 3, pp. 1129–1150, May 2022, doi: 10.1007/s11280-022-01025-2.
- [103] R. M. O. Cruz, W. V. De Sousa, and G. D. C. Cavalcanti, “Selecting and combining complementary feature representations and classifiers for hate speech detection,” *Online Social Netw. Media*, vol. 28, Mar. 2022, Art. no. 100194, doi: 10.1016/j.osnem.2021.100194.
- [104] M. Mastromattei, L. Ranaldi, F. Fallucchi, and F. M. Zanzotto, “Syntax and prejudice: Ethically-charged biases of a syntax-based hate speech

recognizer unveiled,” PeerJ Comput. Sci., vol. 8, pp. 1–19, Feb. 2022, doi: 10.7717/peerj-cs.859.

[105] V. Mercan, A. Jamil, A. A. Hameed, I. A. Magsi, S. Bazai, and S. A. Shah, “Hate speech and offensive language detection from social media,” in Proc. Int. Conf. Comput., Electron. Electr. Eng. (ICE Cube), Oct. 2021, pp. 1–5, doi: 10.1109/ICECube53880.2021.9628255.

[106] A. Razdan and S. Shridev, “Hate speech detection using ML algorithms,” in Proc. Int. Conf. Artif. Intell. Mach. Vis. (AIMV), Sep. 2021, pp. 1–6, doi: 10.1109/AIMV53313.2021.9670987.

[107] S. Dascálu and F. Hristea, “Towards a benchmarking system for comparing automatic hate speech detection with an intelligent baseline proposal,” Mathematics, vol. 10, no. 6, p. 945, Mar. 2022, doi: 10.3390/math10060945.

[108] G. H. Panchala, V. V. S. Sasank, D. R. H. Adidela, P. Yellamma, K. Ashesh, and C. Prasad, “Hate speech & offensive language detection using ML & NLP,” in Proc. 4th Int. Conf. Smart Syst. Inventive Technol. (ICSSIT), Jan. 2022, pp. 1262–1268, doi: 10.1109/ICSSIT53264.2022.9716417.

[109] R. T. Mutanga, N. Naicker, and O. O. Olugbara, “Detecting hate speech on Twitter network using ensemble machine learning,” Int. J. Adv. Comput. Sci. Appl., vol. 13, no. 3, pp. 331–339, 2022, doi: 10.14569/IJACSA.2022.0130341.

[110] A. Kumar, V. Tyagi, and S. Das, “Deep learning for hate speech detection in social media,” in Proc. IEEE 4th Int. Conf. Comput., Power Commun.

Technol. (GUCON), Sep. 2021, pp. 1–4, doi: 10.1109/GUCON50781.2021.9573687.

[111] B. Pariyani, K. Shah, M. Shah, T. Vyas, and S. Degadwala, “Hate speech detection in Twitter using natural language processing,” in Proc. 3rd Int. Conf. Intell. Commun. Technol. Virtual Mobile Netw. (ICICV), Feb. 2021, pp. 1146–1152, doi: 10.1109/ICICV50876.2021.9388496.

[112] Z. Zhang, X. Han, Z. Liu, X. Jiang, M. Sun, and Q. Liu, “ERNIE: Enhanced language representation with informative entities,” in Proc. 57th Annu. Meeting Assoc. Comput. Linguistics, 2019, pp. 1441–1451, doi: 10.18653/v1/p19-1139.

[113] C. Ziems, B. He, S. Soni, and S. Kumar, “Racism is a virus: Anti-Asian hate and counter hate in social media during the COVID-19 crisis,” 2020, arXiv:2005.12423. [Online]. Available: <https://claws.cc.gatech.edu/covid/#dataset>

[114] G. Viswanath, “Hybrid encryption framework for securing big data storage in multi-cloud environment”, Evolutionary intelligence, vol.14, 2021, pp.691-698.

[115] Viswanath Gudditi, “Adaptive Light Weight Encryption Algorithm for Securing Multi-Cloud Storage”, Turkish Journal of Computer and Mathematics Education (TURCOMAT), vol.12, 2021, pp.545-552.

[116] Viswanath Gudditi, “A Smart Recommendation System for Medicine using Intelligent NLP Techniques”, 2022 International Conference on Automation, Computing and Renewable Systems (ICACRS), 2022, pp.1081-1084.

[117] G.Viswanath, “Enhancing power unbiased cooperative media access control protocol in manets”, International Journal of Engineering Inventions, 2014, vol.4, pp.8-12.

[118] Viswanath G, “A Hybrid Particle Swarm Optimization and C4.5 for Network Intrusion Detection and Prevention System”, 2024, International Journal of Computing, DOI: <https://doi.org/10.47839/ijc.23.1.3442>, vol.23, 2024, pp.109-115.

[119] G.Viswanath, “A Real Time online Food Ordering application based DJANGO Restfull Framework”, Juni Khyat, vol.13, 2023, pp.154-162.

[120] Gudditi Viswanath, “Distributed Utility-Based Energy Efficient Cooperative Medium Access Control in MANETS”, 2014, International Journal of Engineering Inventions, vol.4, pp.08-12.

[121] G.Viswanath,“ A Real-Time Video Based Vehicle Classification, Detection And Counting System”, 2023, Industrial Engineering Journal, vol.52, pp.474-480.

[122] G.Viswanath, “A Real- Time Case Scenario Based On Url Phishing Detection Through Login Urls ”, 2023, Material Science Technology, vol.22, pp.103-108.

[123] Manmohan Singh,Susheel Kumar Tiwari, G. Swapna, Kirti Verma, Vikas Prasad, Vinod Patidar, Dharmendra Sharma and Hemant Mewada, “A Drug-Target Interaction Prediction Based on Supervised Probabilistic Classification” published in Journal of Computer Science, Available at: <https://pdfs.semanticscholar.org/69ac/f07f2e756b79181e4f1e75f9e0f275a56b8e.pdf>