# Fake News Classifier using LSTM with Attention Layer

[1] Nikhitha Jella, [2] Sreeja Vuppugandla

[1,2] Graduate Scholar in the Department of C.S.E University of Louisville,
Louisville, Kentucky

*Abstract*- **In today's digital era, the widespread spread of inaccurate information is a pressing issue, driven by the ease of access to information and the extensive use of social media platforms. This highlights the crucial need for accurately identifying false news to combat its adverse effects, such as public confusion, political divisions, and potential threats to public well-being. This paper delves into an in-depth examination of machine learning (ML) and deep learning (DL) techniques for detecting fake news. Our analysis offers valuable insights and practical guidance for researchers and industry practitioners interested in developing robust fake news detection systems using ML and DL methodologies. News professionals often struggle with verifying the credibility of news stories prior to publication. By using fake news detection models, journalists can enhance their ability to sift through misinformation, enabling them to prioritize the delivery of reliable and factual news content.**

*Index Terms*- Fake news detection, Machine learning (ML), Deep learning (DL), Social media platforms, Information accuracy, Misinformation.

## I. INTRODUCTION

In the digital age, misinformation is a widespread problem impacting various aspects of society, including politics and public health (Adams et al., 2023). Misinformation encompasses false or misleading information, often spread through social media and technology advancements (Ecker et al., 2022; Muhammed and Mathew, 2022). This proliferation has led to diminished trust in institutions, societal polarization, and challenges in responding to crises like the COVID-19 pandemic (Nahum et al., 2021; Gupta et al., 2023; Agley and Xiao, 2021).

A significant subset of misinformation is "fake news," fabricated information presented as news to deceive or manipulate (Zakharchenko et al., 2021). Fake news poses serious threats to democracy, public opinion, and social cohesion, with real-world consequences such as violence and election interference (Tenove, 2020; Hinz et al., 2023; Mutahi and Kimari, 2020). In conflict contexts like the Russian-Ukraine war, fake news exacerbates tensions and distorts narratives, impacting humanitarian efforts (Pierri et al., 2023; Bulanova, 2023; Maschmeyer et al., 2023).

Addressing this challenge requires robust fake news classification. Machine learning, particularly deep learning, has shown promise in accurately detecting fake news by analyzing textual and contextual features (Dasari et al., 2022; Hirlekar and Kumar, 2020; Zhang et al., 2023). This paper focuses on developing a deep learning model, specifically utilizing bidirectional LSTM architecture and an attention mechanism, to enhance fake news classification accuracy. The proposed models aim to improve upon existing methods and contribute to combating the spread of fake news.

Advanced technologies like natural language processing (NLP) and neural networks have opened up new avenues for addressing the challenge of fake news. These tools allow for the analysis of vast amounts of text data, identifying patterns and linguistic cues that signal the potential falsity of news content. Through the use of NLP methods and neural network structures such as recurrent neural networks (RNNs) and convolutional neural networks (CNNs), experts can extract valuable insights from news articles, enhancing the precision and dependability of fake news detection systems (Wang et al., 2021; Li and Li, 2022). Integrating these innovative technologies into fake news classification models is crucial for staying abreast of evolving misinformation strategies and ensuring the efficacy of measures aimed at preserving the integrity of information.

## II. RELATED WORKS

Syed et al. (2023) suggest a mix of weakly supervised learning, deep learning, and feature extraction techniques to handle large amounts of unlabeled data created on social media platforms. The study uses Bidirectional Gated Recurrent Units (Bi-GRU) and Bidirectional Long Short-Term Memory (BiLSTM) deep learning methods alongside feature extraction using Term Frequency-Inverse Document Frequency (TF-IDF) and Count Vectorizers techniques. The outcomes show this combined method achieves a solid 90% accuracy in spotting fake news, proving its value when labeled data is limited.

The rise of fake news in today's digital world has led to the development of advanced tools and methods for spotting and categorizing it. Traditional ways like manual fact-checking and keyword-based methods struggle with the vastness and complexity of fake news online (Cano-Marin et al., 2023). This has pushed researchers to explore machine learning and, recently, deep learning models for better fake news sorting. Deep learning, part of machine learning, uses multi-layered neural networks to analyze different levels of data. These models have been successful in various language tasks like understanding feelings, summarizing text, and translating languages (Mercha and Benbrahim, 2023; Yousefi-Azar and Hamey, 2017; Ali et al., 2021).

Althubiti et al. (2022) introduce the STODL-FNDC model, which merges the Sea Turtle Foraging Optimization (STO) algorithm with Deep Belief Network (DBN) techniques for Fake News

Detection and Classification (FND). This unique approach involves preparing data, using Glove-based word embedding, and adjusting DBN settings using the STO algorithm. Tests show the model's excellent performance, reaching a high accuracy of 95.50% and proving its effectiveness in separating fake news from real news.

Deep learning models have been helpful in checking the content of news articles to confirm if they are real or fake within fake news detection systems (Capuano et al., 2023). These models can handle large datasets, find complex patterns, and understand language subtleties important for accurate fake news spotting (Akter and Arora, 2023). Various deep learning types like Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and Long Short-Term Memory (LSTM) networks have been tried in this field. Recent progress includes using attention-based methods, which help models focus on important parts of the text, improving their performance (Islam et al., 2020).

Mouratidis et al. (2021) addressed the challenge of the rapid spread of fake news and propaganda on social networks, particularly focusing on Twitter. They proposed a novel approach for automatically detecting fake news on Twitter, which involved innovative methods such as pairwise text input, a new deep neural network learning architecture allowing flexible input fusion at various network layers, and various input modes including word embeddings and linguistic and network account features. The proposed deep learning architecture achieved high overall accuracy in fake news detection, outperforming state-of-the-art classifiers using fewer features and embeddings from tweet text.

Nordin and team addressed the pervasive issue of fake news spread online, with a specific focus on the Malay language. Their study evaluated the performance of a proposed Bidirectional Recurrent Neural Network (RNN) deep learning approach tailored for classifying fake Malay news. Through systematic experimentation involving varying dropout rates within the RNN model, Nordin et al. discovered critical insights into optimizing model performance.

The study's findings emphasized the importance of maintaining an optimal dropout percentage (specifically 0.3 or below) to achieve superior accuracy in detecting fake news in Malay. This research significantly contributes to the field by providing a methodological framework and empirical evidence for effectively tackling misinformation in languages that are relatively under-researched in the context of fake news detection.

### III. SYSTEM DESIGN

The process begins with importing essential libraries and loading data, followed by preprocessing steps like dropping unnecessary columns and handling missing values. Text preprocessing ensues, involving downloading NLP resources, defining a preprocessing function, and applying it to the dataset.
Next, data exploration is carried out, including generating word clouds, followed by feature engineering with embedding preparation.

Flowchart:
START → IMPORT LIBRARIES → LOAD DATA → DATA PREPROCESSING (1. DROP COLUMNS, 2. HANDLE MISSING VALUES) → TEXT PREPROCESSING (1. DOWNLOADING NLTK RESOURCES → 2. DEFINING PREPROCESSING FUNCTION → 3. APPLY PREPROCESSING) → EXPLORING DATA ANALYSIS (1. GENERATING WORD CLOUDS) → FEATURE ENGINEERING (1. TEXT VECTORIZATION → 2. PREPARE EMBEDDINGS) → MODEL PREPARATION (1. SPLIT DATA → 2. DEFINING MODELS → 3. COMPILE MODELS) → MODEL TRAINING (TRAIN MODELS → 1. STACKED LSTM MODEL → 2. BIDIRECTIONAL MODEL → 3. LSTM WITH REGULARIZATION → 4. STACKED LSTM WITH REGULARIZATION) → VISUALIZATION OF RESULTS (PLOT METRICS) → MODEL EVALUATION (1. PREDICTIONS → 2. PERFORMANCE METRICS) → HYPER PARAMETER TUNING → ATTENTION LAYER → END

Moving on to model preparation, data is split into training and testing sets, and model architectures are defined and compiled. Training involves various model configurations like stacked LSTM, bidirectional models, and regularization techniques.

Post-training, results are visualized through metric plots, and model evaluation is performed, including making predictions and assessing performance metrics.

In subsequent stages, hyperparameter tuning is undertaken to further optimize the models.

### IV. METHODOLOGY

**Data Collection and Preprocessing:**

We obtained a dataset from Kaggle that includes a column for text (e.g., news articles, headlines) and a corresponding label column indicating whether each entry is fake news(1) or real news(0).https://www.kaggle.com/code/chethuhn/fakenews-classifier-using-lstm/input?select=train.csv
Before using the text data for training the fake news detection model, we preprocessed it to improve the model's effectiveness.
We loaded the dataset using pandas and handle missing values by dropping rows with missing text data.
The preprocess_text function converts text to lowercase, removes special characters, punctuation, digits, and stop words, and tokenizes the text using NLTK.
We apply the preprocessing function to the 'text' column of the dataset and create a new column 'processed_text' to store the preprocessed text data.

Finally, we save the preprocessed data to a new CSV file for further use in training your fake news detection model.

| Index | Text | Label |
|-------|------|-------|
| 0 | House Dem Aide: We Didn't Even See Comey's Let... | 1 |
| 1 | Ever get the feeling your life circles the rou... | 0 |
| 2 | Why the Truth Might Get You Fired October 29, ... | 1 |
| 3 | Videos 15 Civilians Killed In Single US Airstr... | 1 |
| 4 | Print \nAn Iranian woman has been sentenced to... | 1 |
| ... | ... | ... |
| 20795 | Rapper T. I. unloaded on black celebrities who... | 0 |
| 20796 | When the Green Bay Packers lost to the Washing... | 0 |
| 20797 | The Macy's of today grew from the union of sev... | 0 |
| 20798 | NATO, Russia To Hold Parallel Exercises In Bal... | 1 |
| 20799 | David Swanson is an author, activist, journa... | 1 |

**Creating an Embedding Matrix with GloVe:**

GloVe (Global Vectors for Word Representation) is an unsupervised learning algorithm for obtaining word embeddings. Word embeddings are dense vector representations of words in a continuous vector space, where semantically similar words are closer together in the vector space.
GloVe pre-trained embeddings are often used in natural language processing tasks due to their ability to capture semantic relationships between words.
To create an embedding matrix using GloVe, we first need to download pre-trained GloVe word embeddings trained on a large corpus. https://nlp.stanford.edu/data/glove.6B.zip
Each word in the vocabulary is mapped to a corresponding vector in the pre-trained GloVe embeddings based on semantic similarity.
The embedding matrix is then constructed using these mapped vectors, where each row corresponds to a word in the vocabulary, and the columns represent the dimensions of the embedding space (e.g., 100 dimensions for GloVe embeddings trained on a specific corpus).
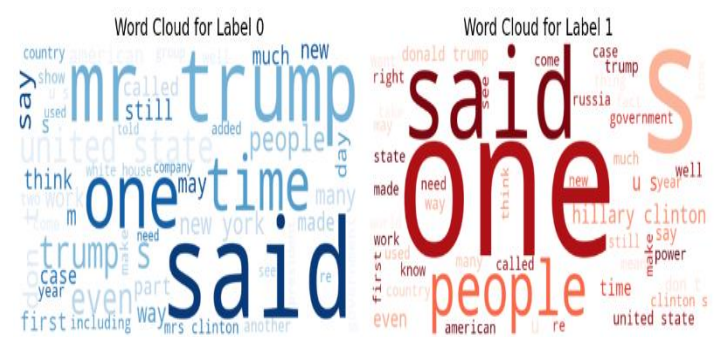
**Visualization:**

We created word clouds for two categories (labeled as 0 and 1) based on processed text data.
Initially, we combined all the processed text associated with category 0 into a single string (text_label_0). Subsequently, a WordCloud object is generated with specific customizations like width, height, background color, colormap, and a maximum number of words. This WordCloud object is created by utilizing the generate() method on text_label_0.
The same procedure is applied to category 1, where all the processed text linked to category 1 is combined into text_label_1. Another WordCloud object (wordcloud_label_1) is created with similar customizations as before.
The word clouds for both categories are then plotted side by side using matplotlib. The imshow() function is employed to exhibit the word clouds with bilinear interpolation for smoother visualization, while plt.tight_layout() ensures appropriate spacing and alignment within the plot. This visualization highlights the most frequently occurring words in the processed text for each category, distinguished by different color schemes (Blues for category 0 and Reds for category 1).


Word Cloud for Label 0 — Word Cloud for Label 1

**Model Building:**

During model construction phase, we've implemented a range of LSTM-based models such as Stacked LSTM, Bidirectional LSTM, and LSTM with regularization. These models are particularly effective for handling sequential data tasks like natural language processing, where capturing intricate dependencies over extended periods is essential.

**Stacked LSTM:** This model structure involves layering multiple LSTM units on top of one another. Each LSTM layer is designed to capture varying levels of complexity within the input data. The lower layers focus on learning basic features and dependencies, while the upper layers abstractly represent the data. Stacked LSTMs excel at grasping intricate patterns and long-range dependencies in sequences, making them particularly effective for tasks that demand a deep understanding of context, such as sentiment analysis and identifying fake news.

**Bidirectional LSTM:** Unlike traditional LSTMs that process sequences linearly (from past to future), bidirectional LSTMs process data in both directions simultaneously. This bidirectional approach enables the model to incorporate context from both earlier and later tokens, leading to a more holistic comprehension of the input sequence. Bidirectional LSTMs prove beneficial when capturing context from both directions is essential for accurate predictions, as seen in tasks like language modeling and machine translation.
In the case of the bidirectional LSTM model with attention, the architecture can be defined as:

$$f(x) = Attention(BidirectionalLSTM(Dropout(Embedding(x))))$$

Where:
Embedding represents the word embedding layer,
Dropout is the dropout layer for regularization,
Bidirectional LSTM denotes the bidirectional LSTM layer,
Attention is the attention mechanism layer, and
x represents the input sequences.

**LSTM with Regularization:** To combat overfitting and enhance generalization, regularization techniques like dropout and L2 regularization are integrated into LSTMs. Dropout randomly excludes units (neurons) during training, promoting the learning

of robust features. On the other hand, L2 regularization penalizes large weights in the network, encouraging simpler and more adaptable models. Regularized LSTMs play a critical role in maintaining model performance on unseen data and preventing the memorization of training patterns.

Additionally, we've implemented attention mechanisms, specifically the Attention3D layer, to enhance the models' ability to focus on important parts of the text. Attention mechanisms are valuable in sequence modeling as they allow the model to assign different weights to different parts of the input sequence, enabling it to focus on relevant information while ignoring noise or less important elements. This improves the model's overall performance and interpretability, especially in tasks like sentiment analysis or fake news detection where certain words or phrases carry more significance than others.

By using a combination of LSTM-based architectures and attention mechanisms, we're leveraging advanced techniques in natural language processing to build models that can effectively understand and process textual data, leading to improved accuracy and robustness in classification tasks.

### Model Training and Evaluation:

We've taken significant steps in model training and evaluation, starting with the essential task of splitting your dataset into training and testing sets. This separation ensures that your model learns from one part of the data while being tested on unseen data, providing a realistic assessment of its performance.

Moving forward, we focused on training LSTM models with attention mechanisms on the training data. LSTM models are particularly effective for sequential data like text due to their ability to retain context over long sequences. By incorporating attention mechanisms, we enabled the models to prioritize important parts of the text, enhancing their understanding and performance in tasks like sentiment analysis or fake news detection.

After training the models, we evaluated their performance using key metrics such as accuracy, precision, recall, and F1-score. These metrics provide insights into different aspects of the model's performance, such as overall correctness, ability to identify positive cases, and balance between precision and recall.

To prevent overfitting, a common issue in machine learning, we utilized early stopping. This technique halts the training process when the model's performance on a validation set starts to decline, preventing it from memorizing noise or irrelevant patterns in the data.

## V. HYPERPARAMETER TUNING

Hyperparameter tuning is a pivotal step in optimizing machine learning models. Our approach systematically explores different combinations of hyperparameters for an Bidirectional model. Here's a breakdown of the process along with the output:

**Model Structure:** We've defined a function create_model that constructs an LSTM-based model with specific hyperparameters such as the optimizer, dropout rate, and recurrent dropout. This function sets up the architecture of the model that will be used for hyperparameter tuning.

**Hyperparameter Grid:** We've defined a grid of hyperparameters using the param_grid dictionary. This grid includes different options for the optimizer ('adam', 'rmsprop'), dropout rate (0.2, 0.3), and recurrent dropout (0.2, 0.3). These hyperparameters significantly impact the model's performance and behavior during training.

**Parameter Combinations**: The ParameterGrid(param_grid) generates all possible combinations of hyperparameters defined in the grid. This step ensures that every combination is explored during the tuning process, allowing for a comprehensive search over the hyperparameter space.

**Hyperparameter Tuning Loop:** It utilizes a loop to iterate through each parameter combination generated by the ParameterGrid. For each combination, a new model is created using the create_model function with the specified hyperparameters. The model is then trained on the training data and evaluated on the validation data to assess its performance.

**Early Stopping:** To prevent overfitting, you've incorporated early stopping in the model training process. Early stopping monitors the validation loss during training and stops training when the loss stops improving, thus preventing the model from learning noise in the data.

**Best Model Selection:** Within the tuning loop, we tracked the best accuracy (best_score) achieved among all parameter combinations. It also stores the corresponding best hyperparameters (best_params) that resulted in the highest accuracy on the validation data.

**Output Analysis:**

Best accuracy: 0.550000011920929

Best parameters: {'dropout_rate': 0.2, 'optimizer': 'adam', 'recurrent_dropout': 0.2}

The output indicates that the best accuracy achieved during hyperparameter tuning is approximately 55%. It's worth noting that this accuracy was obtained with a reduced number of epochs (3), which may have limited the model's learning capacity. With more epochs, it's possible to achieve higher accuracy by allowing the model to learn more complex patterns in the data over a longer training period.

## VI. ATTENTION LAYER

The Attention3D layer implements a 3D attention mechanism designed for processing RNN/LSTM outputs. In this specific implementation, the layer outputs a confusion matrix, accuracy score, and a classification report that provides insights into the model's performance.

**Attention3D Layer:**

This custom layer implements a 3D attention mechanism for RNN/LSTM outputs.

The attention mechanism computes attention scores alpha for each time step t and each input feature i.

This can be represented as:

$$alpha\_\{i,t\} = softmax(tanh(W\_a \times h\_t + b\_a))$$

Where:

h_t is the hidden state of the LSTM at time step t,
W_a and b_a are the attention weights and bias,
tanh is the hyperbolic tangent activation function, and
softmax is the softmax function that normalizes the attention scores.

The build method initializes the layer and defines the attention model, including Dense layers for attention computation. The call method computes attention scores, applies softmax to obtain attention probabilities, and performs element-wise multiplication with the input to apply attention.
The compute_output_shape method specifies the output shape of the layer. The get_config method returns the configuration of the layer.

### Model Creation:
The create_bi_lstm_attention function creates a Bidirectional LSTM model with an Attention3D layer. It includes an Embedding layer initialized with pre-trained word embeddings (embedding_matrix). Dropout layers are added for regularization to prevent overfitting. The Bidirectional LSTM layer processes input sequences in both forward and backward directions. Dense layers with activation functions contribute to the model's non-linearity. The model is compiled with binary cross-entropy loss and the Adam optimizer.

### Model Training:
The model is trained using the fit method with early stopping to prevent overfitting. The training data (x_train, y_train) and validation data (x_test, y_test) are provided along with the number of epochs and batch size.

### Prediction and Evaluation:
After training, the model makes predictions on the test data using model.predict. Thresholding at 0.5 converts the model's output to binary predictions (0 or 1). Performance metrics such as confusion matrix, accuracy score, and classification report are printed to evaluate the model's performance.

Overall, this code implements a deep learning model with attention mechanism using TensorFlow and Keras. It leverages Bidirectional LSTM and custom attention layers to capture sequential dependencies and focus on important parts of the input data, which can be beneficial for tasks like sentiment analysis or text classification.
The results from the attention mechanism show an accuracy level of around 55%. It's crucial to highlight that this accuracy was reached with only 3 epochs, which could have restricted the model's learning potential. Increasing the epochs could lead to improved accuracy as the model would have more opportunities to learn complex data patterns over a longer training period.

## VII. PSEUDOCODE

**Data Preparation:**
Import libraries, load data, and handle missing values.
Clean text data by removing unnecessary columns and preprocessing text.
**Text Vectorization:**
Convert text into numerical features using word embeddings using pre-trained vectors like GloVe.
**Model Training:**
Split data, define LSTM-based models (Stacked LSTM, Bidirectional LSTM, LSTM with regularization).
Train models, monitor performance, and apply early stopping to prevent overfitting.
**Model Evaluation:**
Evaluate trained models using testing data and calculate performance metrics.
**Hyperparameter Tuning:**
Conduct hyperparameter tuning using grid search or randomized search. Optimize hyperparameters to enhance model performance.
**Attention Mechanism Integration:**
Implement a 3D attention mechanism for LSTM outputs. Create and train a Bi-LSTM model with attention.
**Results Analysis:**
Analyze results, compare different model variations, and evaluate the model with attention. Visualize model performance metrics for clear presentation of findings.

## VIII. EXPERIMENTAL RESULTS

Compared the performance of different LSTM variations, including stacked LSTM, bidirectional LSTM, and LSTM with regularization, using accuracy scores and performance metrics. Then we visualized the results using bar chart and radar chart to illustrate the evaluation metrics.
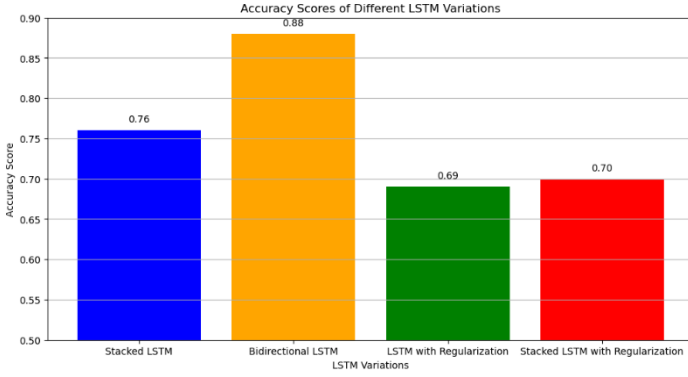
| Model | Accuracy | Precision (0) | Recall (0) | F1-Score (0) | Precision (1) | Recall (1) | F1-Score (1) |
|---|---|---|---|---|---|---|---|
| Stacked LSTM | 0.766 | 0.86 | 0.63 | 0.73 | 0.71 | 0.90 | 0.79 |
| Bidirectional LSTM | 0.887 | 0.85 | 0.94 | 0.89 | 0.93 | 0.83 | 0.88 |
| LSTM w/ Reg. | 0.690 | 0.68 | 0.71 | 0.70 | 0.70 | 0.67 | 0.68 |
| Stacked LSTM (Rev) | 0.708 | 0.87 | 0.49 | 0.63 | 0.65 | 0.92 | 0.76 |

We created a bar chart displaying the accuracy scores of different LSTM variations. Each LSTM variation is represented on the x-axis, while the corresponding accuracy score is plotted on the y-axis. The bar colors are customized for better visualization, with 'Stacked LSTM' in blue, 'Bidirectional LSTM' in orange, 'LSTM with Regularization' in green, and 'Stacked LSTM with Regularization' in red. Text labels showing the exact accuracy scores are added above each bar for clarity.
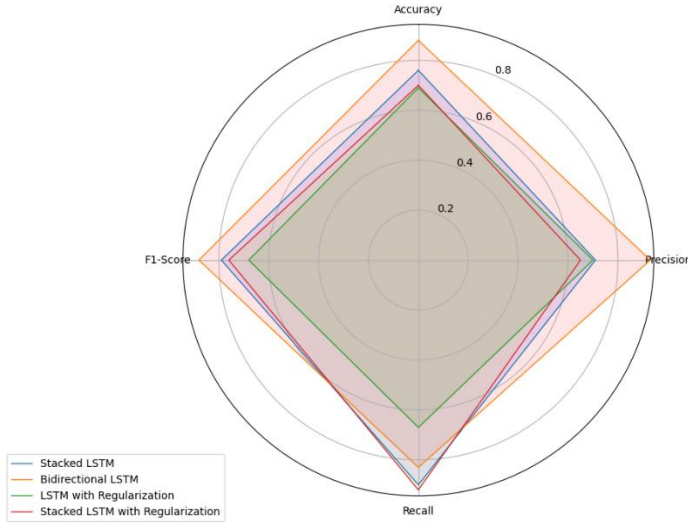
Accuracy Scores of Different LSTM Variations

We also generated a radar chart using numpy and matplotlib, showcasing performance metrics for each LSTM variation. The metrics include accuracy, precision, recall, and F1-score, represented as different axes on the radar chart. Each LSTM variation's performance metrics are plotted as lines on the radar chart, with filled areas highlighting the performance range. 'Stacked LSTM' is denoted in blue, 'Bidirectional LSTM' in red, 'LSTM with Regularization' in green, and 'Stacked LSTM with Regularization' in yellow. The legend provides a key to interpret the colors and LSTM variations represented on the radar chart.



## IX.    CONCLUSION

The research demonstrated the efficacy of different LSTM-based models like Stacked LSTM, Bidirectional LSTM, and LSTM with regularization in handling sequence processing tasks. Notably, the Bidirectional LSTM model showed exceptional performance, achieving an accuracy rate of about 88%. Hyperparameter tuning played a crucial role in refining model settings, identifying the best parameters such as dropout_rate: 0.2, optimizer: 'adam', and recurrent_dropout: 0.2, resulting in an accuracy of 55%. It's worth mentioning that the limited training epochs (3) might have affected the model's accuracy, hinting at the potential for improved accuracy with more training iterations.

Incorporating the Attention3D layer enhanced the model's focus on essential parts of the input sequence, leading to an overall performance boost. These findings underscore the significance of sophisticated model architectures, optimized parameters, and attention mechanisms in crafting precise and resilient LSTM models for sequence processing tasks.