# 3. Methodology

The methodology adopted for developing the Fake News Prediction System involves a systematic approach that includes defining hardware and software requirements, system design, dataset preparation, algorithm implementation, exploratory data analysis, and model evaluation. Each stage contributes to ensuring the model is accurate, efficient, and scalable for future extensions.

## Hardware and Software Requirements

The hardware setup used for this project includes a system equipped with an Intel Core i5 processor or higher, a minimum of 8 GB RAM, and at least 20 GB of free storage space. The system runs on a 64-bit operating system. The software environment consists of Python 3.9 or above, executed through platforms such as Google Colab and Jupyter Notebook. The development utilized several Python libraries: Pandas and NumPy for data handling and numerical computation, NLTK for text preprocessing like stopword removal and stemming, Scikit-learn for feature extraction and model training, Matplotlib and Seaborn for visualization, and the Regular Expressions module (re) for text cleaning tasks.

## System Design

The overall system design for the Fake News Prediction System follows a logical flow from data collection to prediction. The pipeline begins with dataset collection from Kaggle, which contains article titles, authors, and corresponding labels indicating whether the news is real or fake. The next step involves data preprocessing, which includes cleaning the text, removing unnecessary symbols and stopwords, and converting the text to a lower case form. The cleaned data is then transformed into numerical features using the TF-IDF vectorization technique. These features are used to train a Logistic Regression model capable of distinguishing between fake and real news. Once trained, the model is evaluated using performance metrics, and finally, it is employed to predict the class of new unseen news articles. A Data Flow Diagram and a conceptual block diagram depict this flow from raw text input to final classification output.

## Dataset Used

The dataset used in this project is the publicly available Fake News dataset from Kaggle. It includes three main attributes: Author, Title, and Label. The Label column represents whether the news article is fake (1) or real (0). During preprocessing, missing values were replaced with blank strings, and the Author and Title columns were merged into a single Content column for better context. The text data was normalized by converting it to lowercase and removing special characters. Stopwords were removed using the NLTK library, and stemming was applied through the Porter Stemmer to reduce words to their root form, which enhances the model's ability to generalize effectively.

## Exploratory Data Analysis and Visualization

Exploratory Data Analysis (EDA) was conducted to understand the data distribution, identify potential biases, and visualize text-based patterns. The analysis examined the balance between fake and real news samples, the most frequent words, and author-level statistics. Visual tools such as bar charts, histograms, and word clouds were created using Matplotlib and Seaborn libraries. These visualizations provided insights into word frequency, headline structure, and the overall dataset composition, which helped guide preprocessing and feature engineering decisions.

## Algorithm

The project uses a Logistic Regression algorithm, a supervised machine learning approach suitable for binary classification problems. The implementation process begins by converting the preprocessed text into TF-IDF feature vectors, followed by splitting the dataset into training and testing sets. The Logistic Regression model is trained on the training data and later tested using the test data. The performance is evaluated using metrics such as accuracy, precision, recall, and the confusion matrix. Logistic Regression was selected because it is simple, efficient, and performs well on text classification tasks with high-dimensional sparse data. It also serves as a strong baseline model before experimenting with deep learning techniques such as BERT or LSTM in future work.

## Model Evaluation and Future Work

After model training, testing was conducted using unseen data to ensure generalization capability. The model achieved high accuracy in distinguishing between fake and real news articles. Evaluation metrics like precision, recall, and accuracy confirmed the reliability of the system. The results highlight that traditional machine learning models, when combined with proper preprocessing techniques like TF-IDF, can effectively detect fake news. In the future, the system can be extended by integrating deep learning models such as BERT or LSTM for better contextual understanding, and by developing a user-friendly web interface for real-time fake news detection.