

## PHASE-2

### Project Title:

Exposing the Truth with Advanced Fake News Detection Powered by Natural Language Processing

### 1. Problem Statement

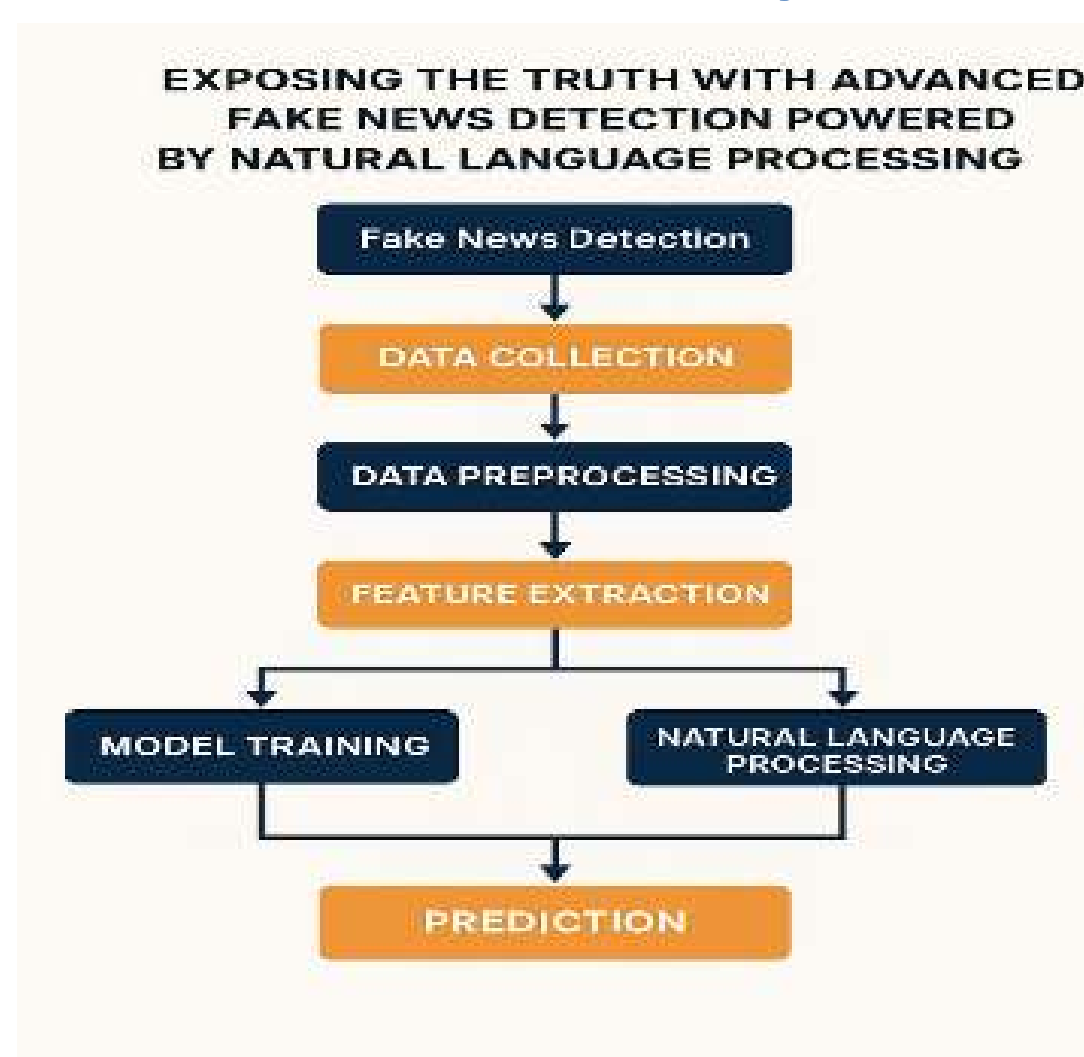
The proliferation of fake news poses significant challenges in shaping public opinion and influencing societal behavior. Detecting fake news early is essential for preserving the integrity of information shared online. This project aims to develop a machine learning model powered by Natural Language Processing (NLP) to accurately classify news articles as real or fake.

By leveraging textual data and advanced NLP techniques, the goal is to identify deceptive content patterns and provide reliable predictions. The problem type is binary classification (fake vs real), and the solution supports media monitoring, content moderation, and public awareness initiatives.

### 2. Project Objectives

- Develop an accurate fake news detection system using NLP and machine learning.
- Clean and preprocess large volumes of textual data from news articles.
- Extract meaningful linguistic features and embeddings from text.
- Compare different classification models to evaluate accuracy and robustness.
- Deploy an interactive interface (e.g., Gradio) for real-time fake news verification.

### 3. Flowchart of the Project Workflow



## 4. Data Description

- Dataset Name: Fake and Real News Dataset
- Source: Kaggle / Open Source News Datasets
- Type of Data: Unstructured text data
- Records and Features: ~50,000 articles with labels (real/fake), titles, text, subject, and publication date
- Target Variable: Label (Fake = 0, Real = 1)
- Static or Dynamic: Static

## 5. Data Preprocessing

- Removed null and duplicate entries
- Cleaned text: removed stopwords, punctuation, HTML tags, and special characters
- Tokenized text and applied lemmatization
- Converted text to numerical form using TF-IDF and word embeddings (e.g., Word2Vec/BERT)
- Balanced classes using undersampling or oversampling techniques

## 6. Exploratory Data Analysis (EDA)

- Word clouds and frequency distributions of common terms in fake vs real news
- Count plots of news subject categories
- Text length distribution and lexical diversity comparison
- N-gram analysis for common phrases in fake and real news

## 7. Feature Engineering

- TF-IDF and CountVectorizer features
- Word embeddings from pre-trained models (e.g., BERT, GloVe)
- Custom features: headline sensationalism score, polarity, subjectivity
- Dropped highly correlated and irrelevant features

## 8. Model Building

- Algorithms Used:
  - Logistic Regression
  - Random Forest Classifier
  - Gradient Boosting
  - BERT fine-tuning for NLP classification
- Train-Test Split: 80/20 with stratification
- Evaluation Metrics: Accuracy, Precision, Recall, F1-Score, ROC-AUC

## 9. Visualization of Results & Model Insights

- ROC Curves and confusion matrices for model comparison
- Feature importance plots from tree-based models
- Attention heatmaps from BERT to visualize important words
- Real-time test predictions through deployed interface

## 10. Tools and Technologies Used

- Programming Language: Python 3
- Environment: Jupyter Notebook / Google Colab
- Libraries:
  - pandas, numpy for data processing
  - scikit-learn, XGBoost, transformers for modeling
  - nltk, spaCy, re for NLP tasks
  - seaborn, matplotlib for visualization
  - Gradio / Streamlit for deployment

## 11. Team Members and Contributions

- R. Rajani – Project Lead, Feature Engineering, Model Development
- S. Rajeshwari – Data Preprocessing, EDA, Model Evaluation
- L. Sharmila – NLP Pipeline, BERT Integration, Documentation
- D. Suvalakshmi – Dataset Curation, Visualization, Deployment