# ACROPOLIS INSTITUTE OF TECHNOLOGY AND RESEARCH

**Department of Computer Science Engineering (Data Science)**

**Synopsis**

on

***"Fake News Detection Using Natural Language Processing"***

## 1. INTRODUCTION

### 1.1. Overview

The proliferation of "fake news"—intentionally false information presented as authentic—on social media threatens societal stability and necessitates sophisticated, automated detection systems. This project presents the design of a robust, multi-faceted system for this purpose, leveraging state-of-the-art Natural Language Processing (NLP) and Graph Neural Networks (GNNs).

Our core philosophy is to move beyond simplistic content analysis by embracing a hybrid methodology. This approach scrutinizes both the linguistic characteristics of a news article and the contextual patterns of its propagation across networks. The project will be implemented in phases, beginning with content-based classification to establish a performance baseline using a traditional TF-IDF with Support Vector Machine model. We will then advance to a deep learning model (BERT) for its superior contextual understanding. A second phase will develop a Graph Neural Network to analyze social context, source credibility, and spread dynamics.

This dual approach ensures a comprehensive analysis resilient to sophisticated disinformation. The final deliverable will be an interactive web dashboard that provides users with not only an article's classification but also interpretable insights into the model's decision-making process using SHAP (SHapley Additive exPlanations), fostering transparency and user trust.

### 1.2. Purpose and Scope

This project's primary purpose is to architect an intelligent, scalable, and explainable tool to combat the pervasive issue of digital misinformation. This system serves as a critical defense mechanism for a diverse range of stakeholders, from social media platforms to individual consumers. Given that manipulated content can sway public opinion and erode trust in institutions, such a tool is paramount.

The goals are multifaceted. For content platforms and news aggregators, it provides an automated first-line-of-defense to pre-filter and flag malicious content at scale, reducing the manual workload for moderation teams. For the public, it is an accessible tool to independently assess the credibility of online news, fostering critical readership and enhancing digital literacy. The system also acts as a powerful assistive technology for journalists and fact-checkers, helping them prioritize articles for investigation. For academic researchers, it creates a robust framework to analyze the complex interplay between linguistic patterns and network dynamics in misinformation, contributing to a deeper understanding of how fake news is constructed and disseminated.

The scope is specifically defined for feasibility. The analysis is confined to textual, English-language news articles and will not extend to detecting multimodal misinformation, such as doctored images or deepfake videos. The final output will provide a probabilistic classification—the likelihood of an article being real or fake—and an explanation for that classification, rather than a comprehensive journalistic fact-check.

## 2. LITERATURE SURVEY

### 2.1. Existing Problem

Detecting fake news is a complex and evolving challenge due to several significant hurdles that limit the effectiveness of current solutions.

- **Scale and Velocity:** The primary challenge is the sheer volume and speed of information online. Manual fact-checking is impractical, as millions of posts are created every minute. By the time a false story is debunked, it has often already achieved viral reach, rendering corrections ineffective.

- **Increasing Sophistication:** Disinformation has grown more advanced. Modern campaigns mimic the style of legitimate news and leverage advanced AI like Large Language Models (LLMs) to generate convincing fake articles at scale. These AI-generated texts can easily evade older, rule-based detection systems.
- **Limitations of Content-Only Analysis:** Automated systems relying solely on text can be fooled by well-written but factually incorrect articles and struggle with satire or propaganda. Crucially, they are blind to the context of *who* is sharing the information and *how* it is spreading—key indicators of coordinated disinformation campaigns.
- **Cognitive Bias:** The problem is also human. Social media algorithms create "echo chambers" that reinforce pre-existing beliefs, making users more susceptible to sharing misinformation that confirms their biases. An effective system must be robust enough to counteract these powerful cognitive effects.

## 2.2. Proposed Solution

To address the multifaceted nature of the fake news problem, we propose a comprehensive, hybrid system that integrates content-based analysis with graph-based propagation analysis. This phased approach allows for a progressively deeper understanding of the information, moving from "what is being said" to "how it is being spread."

**Phase 1: Content-Based Analysis**

This initial phase focuses on extracting linguistic and semantic features directly from the text. We will develop two models in parallel to benchmark performance:

- **Baseline Model (TF-IDF/SVM):** A classical and robust machine learning pipeline. It uses Term Frequency-Inverse Document Frequency (TF-IDF) to convert text into numerical vectors, which are then classified by a Support Vector Machine (SVM).
- **Advanced Model (BERT):** To achieve state-of-the-art performance, we will fine-tune a pre-trained Bidirectional Encoder Representations from Transformers (BERT) model. BERT's ability to learn a deep, contextual

understanding of language is expected to significantly outperform the baseline by capturing nuance and subtle semantic cues.

**Phase 2: Graph-Based Propagation Analysis**

This advanced phase addresses the critical limitations of content-only analysis by modeling the social context of news dissemination.

- **Graph Neural Network (GNN) Model:** We will construct a graph where articles, users, and sources are nodes, and their interactions are edges. A GNN will be trained on this structure to learn patterns indicative of misinformation, such as rapid spread by bot-like accounts.
- **Hybrid Feature Integration:** The contextual embeddings generated by our BERT model will be used as initial features for the article nodes, creating a highly sophisticated hybrid model informed by both linguistic content and network dynamics.

**Phase 3: Interpretable Deployment**

The final phase focuses on making the system accessible and transparent.

- **Interactive Dashboard:** The system will be deployed as a user-friendly web application using Streamlit, providing real-time classification.
- **Explainable AI (XAI):** We will integrate SHAP (SHapley Additive exPlanations) to generate visualizations that highlight the words or phrases most strongly contributing to a classification, thereby building user trust and providing valuable insights.

## 3. THEORETICAL ANALYSIS

### 3.1. Block Diagram

The architecture of the proposed system is designed as a modular, end-to-end pipeline, illustrated in the conceptual workflow below. The process begins with data ingestion from diverse sources, which is then fed into two parallel analytical streams: content-based analysis and graph-based analysis.

The content-based stream preprocesses the raw text (tokenization, lemmatization, stop-word removal) and then feeds it into both the TF-IDF/SVM and the BERT models for training and evaluation. The graph-based stream

constructs a network graph from user-article interaction data, using BERT embeddings as features for the article nodes. This graph is used to train a GNN model. The predictions from the superior model are then passed to the backend of a Streamlit application. The frontend allows users to input text, which the backend processes to deliver a classification and a SHAP-based visualization explaining the prediction. Key performance metrics such as Accuracy, Precision, Recall, and F1-Score will be used throughout to evaluate and compare the models.
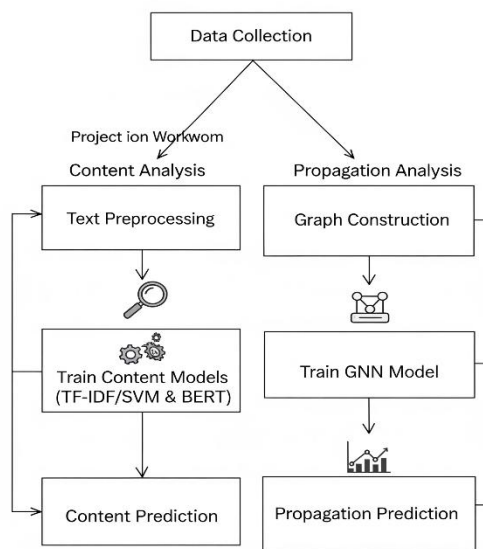


*Fig. 1. Project Workflow for Hybrid Fake News Detection System*

3.2. Hardware/Software Designing

The successful development and execution of this project are contingent upon a specific set of hardware and software resources.

**Software Requirements:**

- **Programming Language:** Python 3.8+ will be the exclusive language for all development, scripting, and modeling tasks.

- **Core Data Science Libraries:**

    o **Pandas:** For high-performance data manipulation, cleaning, and structuring of the news datasets into DataFrames.

- **NLTK / spaCy:** For essential NLP preprocessing tasks, including tokenization, part-of-speech tagging, named entity recognition, and lemmatization.

- **Scikit-learn:** To implement the entire classical machine learning pipeline, including the TF-IDF vectorizer, the SVM classifier, and a suite of tools for model evaluation and metrics calculation.

- **Deep Learning Frameworks:**

  - **PyTorch / TensorFlow:** To serve as the backend for building, training, and running the deep learning models. We will select one based on library support and team familiarity.

  - **Hugging Face Transformers:** A critical library that provides a high-level API for accessing and fine-tuning thousands of pre-trained transformer models, including BERT.

- **Graph Modeling Libraries:**

  - **NetworkX:** For the construction, manipulation, and study of the structure and dynamics of the propagation network graph.

  - **PyTorch Geometric (PyG) / Deep Graph Library (DGL):** Specialized libraries for implementing and training Graph Neural Networks efficiently on top of PyTorch or TensorFlow.

- **Deployment and Visualization:**

  - **Streamlit:** A Python framework for rapidly building and deploying interactive data science web applications and dashboards with minimal front-end coding.

  - **SHAP:** The core library for implementing model interpretability and generating visualizations that explain model predictions.

  - **Matplotlib / Seaborn:** For creating static charts, graphs, and plots for data analysis and inclusion in the final project report.

**Hardware Requirements:**

- **CPU:** A modern multi-core processor (e.g., Intel Core i7 / AMD Ryzen 7 or better) is required for efficient data preprocessing and handling of non-GPU intensive tasks.

- **RAM:** A minimum of **16 GB of RAM**

- **GPU:** A dedicated NVIDIA GPU with CUDA support and **8 GB of VRAM**

- **Cloud Alternative:** As an alternative to local hardware, cloud-based GPU instances such as Google Colab Pro or AWS EC2 will be considered to ensure timely model training.

- **Storage:** A minimum of 50 GB of fast SSD storage is required to house the datasets, development environments, Python libraries, and the saved trained model files.

## 4. APPLICATIONS

The developed hybrid system for fake news detection has a broad spectrum of practical applications across various sectors:

- **Social Media Platforms:** It can be integrated into content moderation pipelines to automatically flag, down-rank, or label potential misinformation, curbing its viral spread.

- **News Aggregators and Search Engines:** The system can be used to vet sources and articles, promoting higher-quality information and helping users distinguish credible journalism from propaganda.

- **Financial Markets:** It can scan financial news and social media for rumours or fabricated reports designed to manipulate prices, protecting investors from fraud.

- **Public Health Sector:** The tool can identify and flag dangerous medical misinformation during health crises, helping public health officials to counter it with accurate information.

- **Brand Reputation Management:** Corporations can use the system to monitor online conversations and detect the spread of false information about their products or services, allowing for rapid responses.

- **Digital Literacy and Education:** The dashboard can be used as an educational tool to teach students how to critically evaluate online sources and become more discerning consumers of information.

# REFERENCES

[1] Devlin, J., Chang, M.W., Lee, K., and Toutanova, K., "BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding," in *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics*, Minneapolis, MN, 2019.

[2] Jurafsky, D. and J.H. Martin, *Speech and Language Processing*, 3rd ed. Upper Saddle River, NJ: Prentice Hall, 2023.

[3] Kipf, T. N., and Welling, M., "Semi-Supervised Classification with Graph Convolutional Networks," in *International Conference on Learning Representations (ICLR)*, Toulon, France, 2017.

[4] Lundberg, S. and S. Lee. (2017, Dec. 25). *A Unified Approach to Interpreting Model Predictions* [Online]. Available: https://github.com/slundberg/shap.

[5] Shu, K., Sliva, A., Wang, S., Tang, J., and Liu, H., "Fake News Detection on Social Media: A Data Mining Perspective," *ACM SIGKDD Explorations Newsletter*, vol. 19, issue 1, pp. 22-36, Sep 2017.

[6] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I., "Attention Is All You Need," in *Advances in Neural Information Processing Systems 30*, 2017, pp. 5998-6008.

**Internal Guide:**

Prof. Dr. Manish Vyas

**External Guide:**

Ms. Surbhi Jain

**Group Members:**

Aditi Jain (0827CD231002)

Anshika Jain (0827CD231008)

Anurag Hamora (0827CD231011)

Dev Mohan Saxena (0827CD231023)