

**[CM3070]**  
**Final Project**

**Final Report**

Yau Feng Yuen, Gabriel

# Contents

[1.0] Introduction .....	4
[1.1] Template .....	4
[1.2] Project Overview.....	4
[1.3] Motivation .....	5
[1.4] Related Projects.....	6
[1.4.1] DataFlair: Detecting Fake News with Python and Machine Learning .....	6
[1.4.2] ClaimBuster.....	7
[1.4.3] Snopes .....	8
[1.5] Project Aims.....	8
[2.0] Literature Review .....	9
[2.1] Machine Learning Approaches .....	9
[2.11] Passive-Aggressive Classifier .....	9
[2.12] Random Forest Classifier .....	10
[2.2] BERT .....	11
[2.3] Hybrid Approaches .....	12
[2.31] ClaimBuster.....	13
[2.32] FaKnow .....	14
[2.4] Human-operated Fact Checking Websites.....	15
[2.5] Evaluation .....	16
[3.0] Design .....	18
[3.1] Project Overview.....	18
[3.2] Domain & Users .....	18
[3.21] Domain.....	18
[3.22] Users .....	19
[3.3] Justification for Design Choices .....	22
[3.31] Command Line Interface Application.....	22
[3.32] Bidirectional Representations from Transformers (BERT) .....	22
[3.4] Overall Structure.....	23
[3.5] Technologies & Methods.....	24
[3.6] Work Plan.....	25
[3.7] Testing & Evaluation.....	27
[3.71] Test Plan .....	27
[3.72] Evaluation Metrics .....	27

[4.0] Feature Prototype .....	28
[4.1] Overview .....	<b>Error! Bookmark not defined.</b>
[4.2] Implementation .....	<b>Error! Bookmark not defined.</b>
[4.21] Dataset Preparation .....	<b>Error! Bookmark not defined.</b>
[4.22] Model Training .....	<b>Error! Bookmark not defined.</b>
[4.23] Evaluation .....	<b>Error! Bookmark not defined.</b>
[4.24] Explainability .....	<b>Error! Bookmark not defined.</b>
[4.3] Results of Initial Implementation .....	<b>Error! Bookmark not defined.</b>
[4.4] Further Iterations .....	<b>Error! Bookmark not defined.</b>
[4.41] Mixed Precision Training .....	<b>Error! Bookmark not defined.</b>
[4.42] Gradient Accumulation and Clipping .....	<b>Error! Bookmark not defined.</b>
[4.5] Future Improvements .....	<b>Error! Bookmark not defined.</b>
[5.0] References .....	29

# **[1.0] Introduction**

## **[1.1] Template**

[CM3060] Natural Language Processing – “Fake News Detection”

## **[1.2] Project Overview**

The spread of misinformation in the modern era presents significant and multifaceted challenges. Fake news and its ability to rapidly propagate falsehoods and distort perceptions has eroded trust, destabilised institutions, and undermined democratic processes.

This project aims to address this issue by leveraging Natural Language Processing (NLP) techniques to design a scalable, accurate & user-friendly fake news detection system that will empower a diverse range of primary users – including journalists and educators – to swiftly and reliably access the credibility of news they consume.

# [1.3] Motivation

This project is motivated by the need to address the real-world consequences of fake news that span several domains:

## **1. Erosion of trust in journalism**

The increasing prevalence of fake news undermines public confidence in legitimate news sources. Trustworthy journalism is vital for an informed society, and the erosion of trust in journalism polarises communities, fosters hostility and impedes constructive discourse.

## **2. Economic disruptions**

Fake news can manipulate markets and harm businesses. For instance, baseless rumours about a company's fiscal health can lead to a sudden crash in stock prices, unfairly & disproportionately affecting vulnerable stakeholders, retail investors and small businesses.

## **3. Public health risks**

The COVID-19 pandemic highlighted how misinformation can incite panic – for example, in Singapore, misinformation about supply shortages prompted the local populace to hoard masks and staple foods unnecessarily, straining supply chains and increasing social anxiety unnecessarily.

## **4. Political impact**

Disinformation campaigns are often weaponised to distort public opinion, promote political agendas and undermine elections. These campaigns erode the public's faith in governance.

With the increasing volume of misinformation and the speed of which it is spread, the demand for tools to quickly and effectively verify the credibility of digital news is higher than ever before. This project is tailored to address the distinct needs of two primary user groups:

### **1. Journalists**

Journalists need tools to verify the credibility of sources and combat the spread of false information. An effective solution will improve the credibility and impact of journalistic work in safeguarding the integrity of public discourse.

### **2. Educators**

Educators need tools to empower individuals to navigate constantly-evolving, complex digital landscapes by enhancing media literacy and critical thinking ability.

## [1.4] Related Projects

### [1.4.1] DataFlair: Detecting Fake News with Python and Machine Learning

DataFlair uses a passive-aggressive classifier to train a model on a small dataset of 7796 rows. [1]

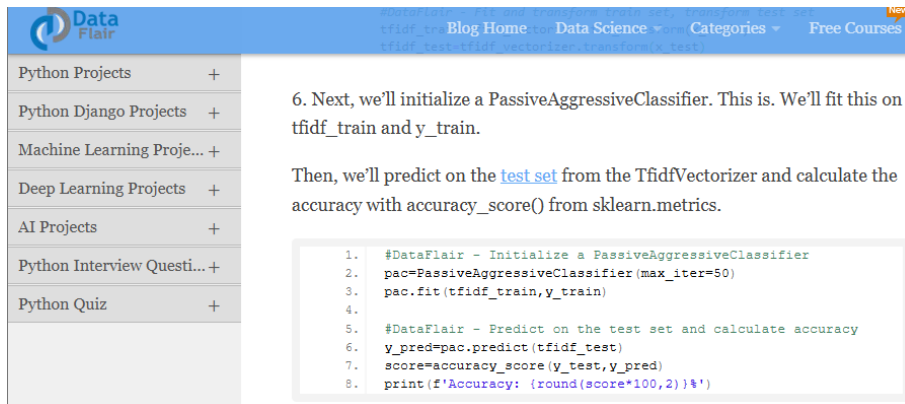


Figure 1: Screenshot of project tutorial on DataFlair [1]

While this project serves as an excellent technical demonstration for beginners, the pipeline demonstrated in this project is too simple for deployment in real-world applications. The model is also trained on an extremely small dataset, which further limits its ability to generalise to real-world applications.

As such, DataFlair's project emphasises the need for a solution that leverages state-of-the-art technologies, and is trained on larger datasets that would not inherently inhibit the model's ability to generalise.

## [1.4.2] ClaimBuster

ClaimBuster applies Natural Language Processing techniques to identify and evaluate factual claims primarily within political discourse.

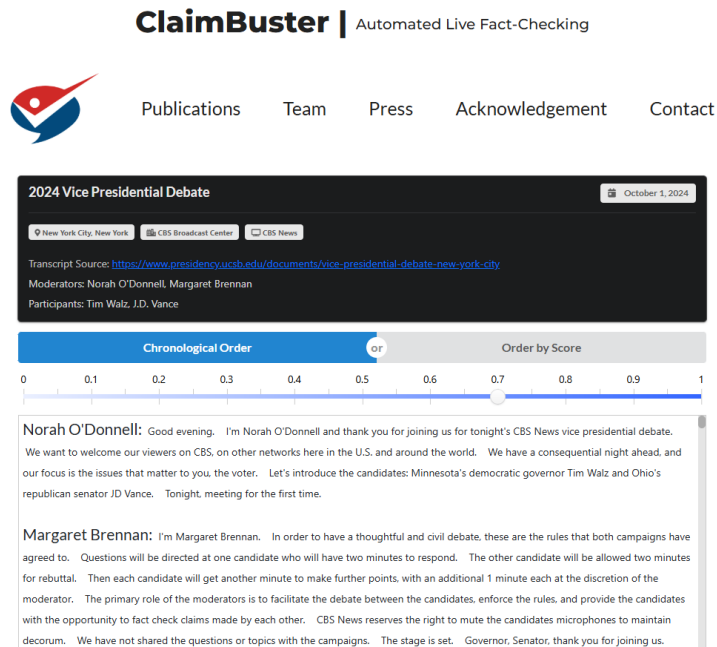


Figure 2: Screenshot of ClaimBuster's user interface [2]

While ClaimBuster is effective in detecting factual claims, its narrow focus on political contexts limits its applicability to broader domains, underscoring need for a solution that is generalisable and scalable across various topics and misinformation domains.

## [1.4.3] Snopes

Snopes is a widely-recognised platform for human fact-checking that relies on human subject-matter experts to manually verify claims, evaluate sources, and analyse misinformation across various topics.

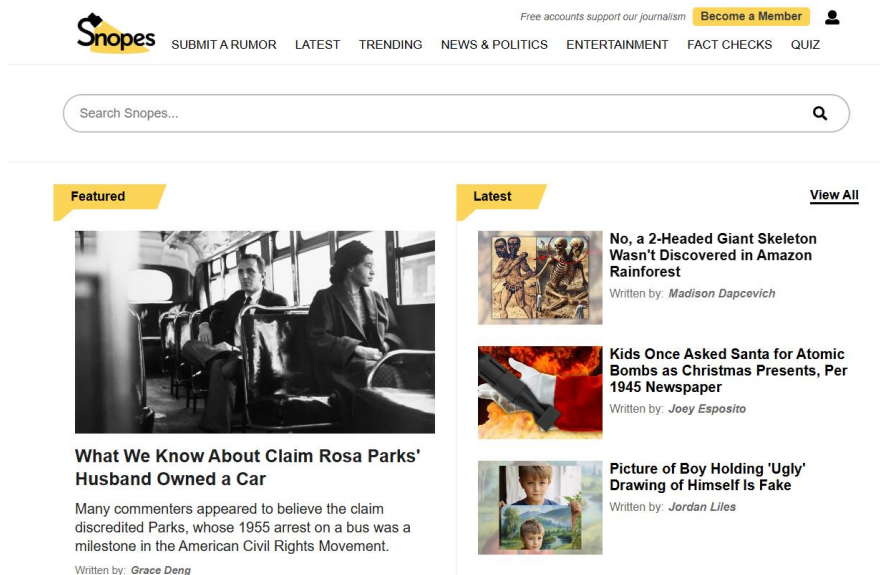


Figure 3: Screenshot of Snopes' homepage [3]

While Snopes is highly credible and thorough, its manual nature is inherently slow and not scalable for real-time verification tasks. Thus, its labor-intensive nature highlights the need for algorithmic alternatives that can complement and augment human expertise.

## [1.5] Project Aims

This project aims to bridge the identified gaps by introducing a comprehensive solution with the following core contributions:

### 1. Scalability through transformer-based models

This project aims to address the limitations of human-intensive methodologies by leveraging cutting-edge machine learning & deep learning techniques that are adaptable to high volumes of information, such as transformer models (e.g. BERT).

### 2. Cross-domain generalisability

This project will be designed to detect misinformation over a wide range of topics and formats, to ensure its relevance and utility in addressing fake news.



## [2.0] Literature Review

Detecting and mitigating the spread of fake news has become a critical area of research within the discipline of natural language processing (NLP). This literature review examines existing approaches and technologies in fake news detection with a focus on their relevance to their implications on the design and implementation of my solution.

### [2.1] Machine Learning Approaches

Several projects have tackled the challenge of fake news detection with machine learning approaches.

#### [2.11] Passive-Aggressive Classifier

The project brief included an example project by DataFlair that uses a small dataset of 7796 rows with news articles labelled as true/fake. TF-IDF vectorisation is first performed to pre-process text data, then a Passive-Aggressive Classifier is applied to train a machine learning model. [1]

This resource serves as an excellent technical demonstration for beginners as it introduces & explains the fundamental steps of data pre-processing, feature extraction, model training, & pipeline building.

However, this resource has too many limitations:

1. **The dataset is not representative of real-world scenarios** – it is too small and lacks variety, thus limiting the model's ability to generalise effectively.
2. **The model is excessively simplistic** – Passive-Aggressive Classifiers are easy to implement but is inadequate to handle nuanced language patterns & context in fake news, and lacks the sophistication of modern NLP models. This is further evidenced by Chang's (2024) comparison of machine learning & deep learning algorithms, where Passive-Aggressive Classifiers were ranked 7<sup>th</sup> in fake news detection performance compared to other algorithms and outperformed by all deep learning algorithms [4].

These limitations restrict the model's ability to handle real-time detection tasks or large-scale datasets and thus is not suitable for deployment in real-world settings.

## [2.12] Random Forest Classifier

Random forest classifiers aggregate multiple decision trees and have been widely employed in fake news detection research.

A potential problem with Random Forest Classifiers is that as their performance declines on imbalanced datasets. Huh (2021) demonstrated that standard Random Forest classifiers achieved a relatively high false-negative rate, a recall of 0.102 and precision of 0.365. This illustrates that 36.5% of predictions in the minority class are correct, and that standard Random Forest classifiers failed to predict approximately 90% of the minority class when presented with an imbalanced dataset. [5]

While it is possible to improve the performance of Random Forest Classifiers when dealing with imbalanced datasets, in Chang's 2024 comparison of ML and DL algorithms, random forest classifiers ranked 9<sup>th</sup>, outperformed by all tested all deep learning algorithms. [4]

## [2.2] BERT

BERT (Bidirectional Encoder Representations from Transformers), first introduced by Devlin et al. (2019), is designed to pre-train deep bidirectional representations from unlabelled text by jointly conditioning on both left and right context in all layers [6], allowing BERT models to effectively model contextual relationships in textual data and discern nuanced linguistic patterns.

BERT's capabilities make it very suitable for fake news detection:

### 1. BERT outperforms other algorithms.

Chang's (2024) comparison of machine learning & deep learning algorithms has shown that BERT is the highest-performing algorithm amongst those tested, with unprecedented scores of 99.95% accuracy, precision, recall & F1 score. [4]

### 2. BERT is extremely flexible.

BERT's transformer architecture allows fine-tuning for task-specific applications, improving domain adaptability.

BERT models, however, have their limitations:

### 1. BERT has high computational demands.

Transformer architectures like BERT have high computational demands. BERT's pre-training corpus included Google's BookCorpus and English Wikipedia, comprising 800 million and 2500 million words respectively for a total of 3300 million words [6].

However, innovations such as MobileBERT, a thinner model that requires less computational power [7], have surfaced to meet the rising demand to leverage BERT's capabilities within smaller computational environments.

### 2. BERT (and other transformer models) are inherently un-interpretable.

Korolev et. al (2023) asserts that state-of-the-art models such as BERT are highly parameterised black boxes [8]. The opacity of transformer-based models has driven interest & demand in explainable AI techniques such as attention visualisation.

Thus, BERT's unparalleled performance establishes it as a cornerstone for this project. However, in implementing BERT, it is important to address scalability concerns through the implementation of lighter variants such as MobileBERT that are less computationally intensive, and to address interpretability challenges by incorporating explainable AI techniques to foster user trust in the model.

## **[2.3] Hybrid Approaches**

Hybrid approaches to fake news detection typically combine content-based and context-based analyses. They often integrate multimodal data – such as text, images, social network metadata – and utilise diverse sources of information to enhance their accuracy and performance.

Hybrid approaches aim to provide a more comprehensive perspective on misinformation propagation through the integration of complimentary data streams.

## [2.31] ClaimBuster

ClaimBuster is a web-based, automated, live fact-checking tool developed by the University of Texas [9].

Hassan et al. (2017) detailed how ClaimBuster works in their published paper “ClaimBuster: The first-ever end-to-end fact-checking system”. [10]

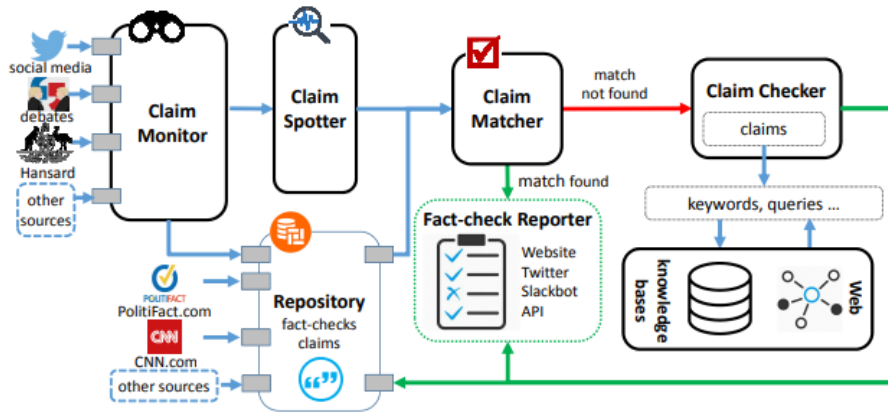


Figure 4: ClaimBuster's system architecture [10]

ClaimBuster implements:

1. A **'claim monitor'** that continuously monitors and retrieves texts from various sources including broadcast media through a decoding device to extract closed captions, social media, and websites (such as the transcripts of Australian parliament proceedings). [10]
2. A **'claim spotter'** that scores sentences' likelihood of containing a factual claim using a classification & scoring model. [10]
3. A **'claim matcher'** takes an important factual claim and searches a fact-check repository and returns fact-checks matching the claim. If a matching fact-check cannot be found, the claim checker queries external knowledge bases and the Web with a question generation tool, before a fact-check reporter finally delivers a report to users through the project website. [10]

ClaimBuster is effective in structured environments such as political debate analysis, where domain-specific knowledge bases facilitate precise fact-checking. Additionally, as pre-verified claims are combined with linguistic patterns, ClaimBuster offers a high degree of accuracy.

However, due to the manual effort required for database curation, scalability of ClaimBuster is constrained, and its applicability in rapidly-evolving contexts may thus be limited. Additionally, as it relies on domain-specific datasets, its ability to adapt to diverse misinformation scenarios may be restricted.

## [2.32] FaKnow

FaKnow is a library designed to standardise the development & evaluation of fake news detection algorithms by integrating various fake news detection algorithms. [11] The library includes a variety of widely-used models, categorised into content-based & social context-based approaches. Additionally, it also offers functionalities for data processing, model training, evaluation, visualisation and logging to enhance reproducibility and reduce redundancy in fake news detection research.

Graves and Cherubini (2016) emphasise that reproducibility is a persistent challenge in misinformation research due to the lack of standardised datasets and frameworks [12]. FaKnow addresses this critical issue by standardising the implementation of various fake news detection algorithms. Additionally, by encompassing both content-based and context-based models, researchers can use FaKnow to explore & evaluate various different approaches within a single platform.

However, FaKnow has several limitations within the context of my project:

### **1. PyTorch framework dependency**

As FaKnow is built on PyTorch, researchers using other deep learning frameworks (e.g. TensorFlow) may face challenges integrating existing models/workflows to this library. Dependency on a single framework may limit adoption amongst researchers with established preferences or institutional constraints.

### **2. Continuous maintenance demand**

As fake news detection algorithms rapidly advance, continuous updates to the library are required to incorporate the latest models and techniques. This presents maintenance challenges for both users and developers of FaKnow.

Ultimately, while FaKnow presents limitations in its PyTorch dependency and continuous maintenance demand to stay abreast of the latest advancements, its strengths in standardisation and usability make it a valuable tool.

## [2.4] Human-operated Fact Checking Websites

Human-operated fact checking websites such as PolitiFact, FactCheck.org, and Snopes, play a critical role in combating misinformation by employing teams of human experts to evaluate the truthfulness of news content and factual claims.

These websites have earned a wealth of perceived credibility and user trust in two ways:

1. **The rigorous methodologies that human fact-checkers employ ensure credibility due to the high accuracy in identifying and categorising false information.**

For example, PolitiFact displays a “Truth-O-Meter” that provides a granular evaluation of claims which serves to aid public understanding of misinformation.

2. **These websites provide detailed context and explanations for their evaluations.**

This in turn promotes media literacy and critical thinking amongst users, but also transparency which helps to foster user trust.

Liu et al. (2023) assert that users “doubt the ability of machines to adjudicate factual disputes and thus perceive AI-based fact-checkers as less credible than human-based fact-checking services”. [13] Yang et al. argue that this is because automated methods are error-prone [14].

However, human fact-checking also has inherent limitations:

1. **Human fact-checking is not scalable.** Lin et al. (2023) assert that traditional manual fact-checking is time-consuming and labor-extensive, which cannot scale with the unprecedented amount of disinformation and misinformation on social media [15]. This limits their impact in countering real-time viral narratives.
2. **Perceived (or actual) biases in the evaluation of factual claims can undermine user trust,** especially amongst ideologically polarised audiences.

Ultimately, this is relevant to my solution as insights from the strengths of human-operated systems can help to inform the integration of human-like explainability features to foster user trust.

# [2.5] Evaluation

Method/Approach	Strengths	Weaknesses	Example Use Case
Passive-Aggressive Classifier (ML)	<ul style="list-style-type: none"><li>- Simple implementation</li></ul>	<ul style="list-style-type: none"><li>- Limited generalisability</li><li>- Does not capture nuance</li></ul>	<ul style="list-style-type: none"><li>- <u>DataFlair</u> example project</li></ul>
Random Forest Classifier (ML)	<ul style="list-style-type: none"><li>- Robust to overfitting</li></ul>	<ul style="list-style-type: none"><li>- Poor performance on imbalanced datasets</li></ul>	<ul style="list-style-type: none"><li>- Academic experiments</li></ul>
BERT (Deep Learning)	<ul style="list-style-type: none"><li>- State-of-the-art contextual understanding</li></ul>	<ul style="list-style-type: none"><li>- High computational demand</li><li>- Lack of interpretability</li></ul>	<ul style="list-style-type: none"><li>- Cutting-edge NLP tasks</li></ul>
<u>ClaimBuster</u>	<ul style="list-style-type: none"><li>- High accuracy in structured contexts</li></ul>	<ul style="list-style-type: none"><li>- Limited scalability</li><li>- Limited domain applicability</li></ul>	<ul style="list-style-type: none"><li>- Political debate fact-checking</li></ul>
<u>FaKnow</u>	<ul style="list-style-type: none"><li>- Standardised evaluation</li><li>- Reproducible</li></ul>	<ul style="list-style-type: none"><li>- Requires continuous maintenance/updates</li><li>- <u>PyTorch</u> dependency limits flexibility</li></ul>	<ul style="list-style-type: none"><li>- Academic research</li></ul>
Human fact-checking	<ul style="list-style-type: none"><li>- High credibility and user trust</li></ul>	<ul style="list-style-type: none"><li>- Labour-intensive</li><li>- Not scalable</li></ul>	<ul style="list-style-type: none"><li>- Snopes, PolitiFact</li></ul>

Figure 5: Tabular summary of fake news detection methods

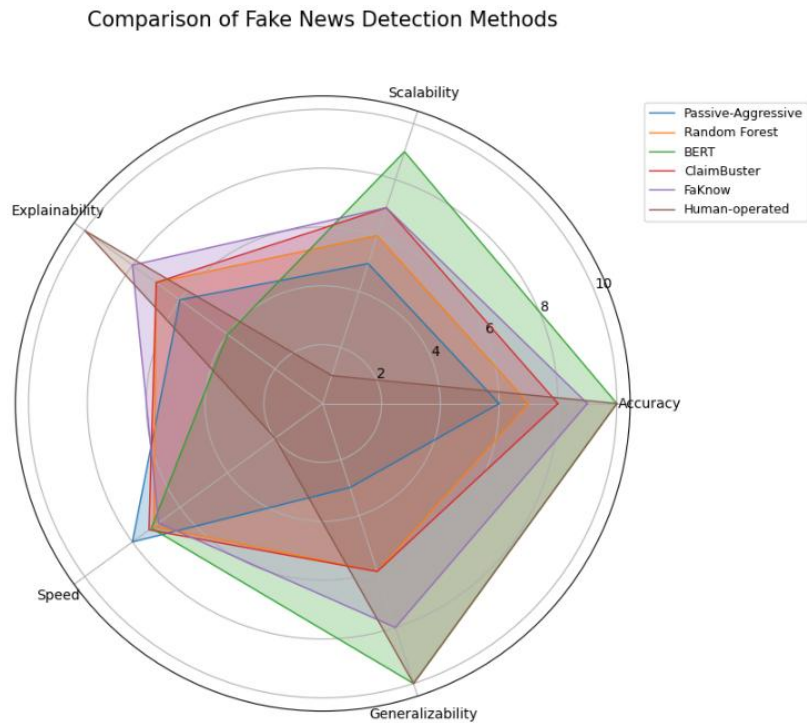


Figure 6: Radar chart comparing key attributes of fake news detection methods



My review of existing methodologies highlights several critical areas for innovation:

### **1. Scalability**

My solution should address the scalability limitations of human fact-checking that drives demand for algorithmic fact-checking in the first place, through incorporating optimising transformer architectures (i.e. BERT).

### **2. Transparency**

To foster user trust & credibility, the development of interpretable AI methods to address the “black box” nature of existing models is crucial. Ideally, these methods would be more human-like, such as through short paragraphs as opposed to attention heatmaps.

### **3. Cross-domain applicability**

My solution should be generalisable to address misinformation across various topics and formats.

## **[3.0] Design**

### **[3.1] Project Overview**

Project Template: [CM3060] Fake News Detection

This project aims to address the issue of misinformation by designing an AI-powered Fake News Detection system that will leverage advanced Natural Language Processing (NLP) techniques to classify user-inputted articles as fake or real. This will be delivered as a Command-Line Interface (CLI) application that will empower users to verify the credibility of online content in real-time.

To achieve that goal, the primary objective is to build a data pipeline to process and link data effectively, so it can be coupled with a query generation strategy that will deliver on accuracy.

### **[3.2] Domain & Users**

#### **[3.21] Domain**

This project falls under the domain of information verification within the field of journalism & digital media.

The challenges this domain addresses include:

- identifying misinformation,
- improving public media literacy, and
- mitigating the societal impact of fake news.

## [3.22] Users

My proposed solution is built with the view that a diverse group spanning various demographics would find value in this solution.

These users would all have varied requirements and use cases, and as such I have elected to focus on groups of primary users:

### **1. Journalists**

My proposed solution will create value for journalists. Their demographic profile spans ages 25 to 60, with a media scope that includes traditional newspapers, digital platforms, freelance journalism & multimedia reporting.

Journalists operating in different regions will have different regional challenges. For example, journalists in the Middle East face misinformation challenges addressing religious & ethnic sensitivities, increasing the need for source validation tools, whereas Western journalists struggle to target digital disinformation – specifically, journalists in democracies such as the United States have challenges combating election-related fake news.

Journalists would use my solution to:

- fact-check sources,
- build credibility in reporting, and,
- investigate information cascades.

Meanwhile, there are journalists that will not use my application. One such sub-demographic of journalists exists in highly state-controlled environments, such as North Korea, where independent news verification is not feasible or allowed. Additionally, journalists involved in creating or propagating fake news are certainly not going to use nor promote a fake news detection solution.

With all the above factors considered, I expect journalists to be primary users and adopters of my solution as they have a direct need for tools that aid in fact-checking and verifying sources in their everyday work. Given the immediacy of their nature of work, they require efficient & accurate solutions.

## **2. Educators**

Educators that are primarily based in academia or non-profit organisations would gain value from my solution, and are likely to heavily use and/or promote it.

Media literacy educators would be foremost adopters amongst this demographic, as they are responsible for equipping students & communities with tools to identify misinformation. They would promote adoption of the solution by recommending (and thus, marketing) it to students and workshop attendees alike.

That said, adoption of my solution would not be exclusive to media literacy educators. Educators spanning various disciplines, teaching various subjects to various demographics of students (ranging from high school to university-level students) have an interest in teaching students to identify fake news when consuming media, and they would also promote adoption of my solution to their students.

Not all educators should be expected to promote or adopt my solution. Educators in some disciplines, such as early childhood education, for instance, would face different challenges and their lesson plans would naturally have very different learning outcomes that would not include discerning news consumption.

As educators act as ‘force multipliers’ by teaching communities of students & professionals to critically evaluate news, they play a vital role in adopting and promoting my solution and are thus expected to be primary users & adopters of my solution.

On top of the primary users, I have identified groups of secondary users who are likely to adopt my application.

## **3. Researchers (secondary users)**

Typically aged 30-60, researchers, often affiliated with universities or think tanks, often operate in data-rich environments with access to large datasets, computational resources and tools for statistical and NLP analysis.

While researchers - especially those investigating the spread & impact of misinformation, the role of algorithmic amplification, and its societal implications - will have a strong interest in analysing misinformation patterns, they are expected to be secondary users as their use of a CLI application may be occasional or project-specific rather than daily (as expected of primary users such as journalists).

Additionally, as their use case relates to academic analysis, tailoring tools to focus on academic rigor and reproducibility requires the implementation of advanced customisation ability that, if included, would make the solution excessively complex for primary users.

#### **4. Content moderators for social media platforms (secondary users)**

Content moderators working for social media platforms are tasked with identifying and removing misinformation from their platforms. They play a critical role in curbing the spread of fake news. Demographically, they tend to be professionals in their 20s to 40s, often with backgrounds in communication or technology.

Content moderators often handle large volumes of flagged content daily under strict organisational guidelines. They are secondary users as they already utilise other tools primarily, usually GUI-based systems integrated into a larger moderation platform. The proposed solution may be added to their workflow as a supplemental tool for cases where deeper analysis is required.

## **[3.3] Justification for Design Choices**

The design choices are meticulously aligned with user needs, and the demands of the domain.

### **[3.31] Command Line Interface Application**

Command-line interface (CLI) applications are lightweight and platform-independent, meaning that the application can operate seamlessly across various operating systems, ensuring broad accessibility for technically-proficient users.

While a CLI application is not ideal for non-technical users, the primary and secondary users of this application are expected to be technically proficient, and either already experienced with CLI or are easily teachable.

However, a possible avenue for future work is the inclusion of a Graphical User Interface (GUI) to cater to more audiences.

### **[3.32] Bidirectional Representations from Transformers (BERT)**

BERT is chosen as the machine learning model for several reasons – firstly, its ability to capture nuanced word meanings based on context will enhance classification accuracy for complex, ambiguous tests.

Additionally, BERT is expected to have high domain adaptability – fine-tuning BERT on a domain-specific dataset is expected to increase performance in fake news detection.

Benchmarks consistently demonstrate BERT's superiority over traditional NLP models, indicating it is likely the optimal model to implement in my solution. However, its high computational demands and lack of inherent interpretability needs to be addressed through using lightweight variants (e.g. MobileBERT) and explainable AI tools respectively.

## [3.4] Overall Structure

The project's architecture comprises the following layers:

### 1. Data Layer

The data layer incorporates datasets from multiple sources such as labelled misinformation repositories and domain-specific datasets. Structured storage mechanisms (e.g. SQLite) will be used to manage raw data and processed features for easy access & reproducibility.

The pre-processing pipeline is also implemented here – with tokenisation, stopword removal, stemming & lemmatisation, and feature engineering (TF-IDF vectorisation & embeddings) to transform raw data to model-ready input.

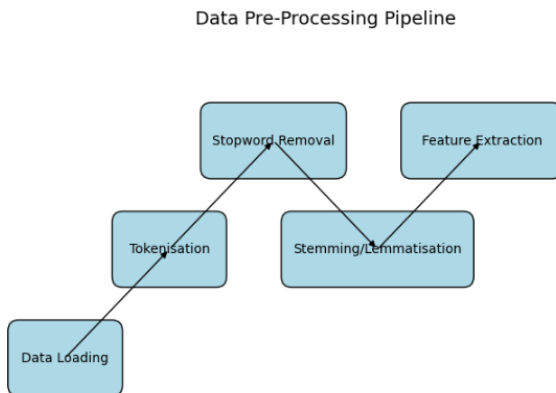


Figure 7: Data pre-processing pipeline in data layer

### 2. Modelling layer

The modelling layer comprises traditional machine learning models (e.g. logistic regression, SVM) for comparative benchmarking, as well as deep learning models (i.e. BERT) fine-tuned for fake news detection. These models are combined to improve overall accuracy and robustness.

Additionally, this layer also incorporates additional NLP techniques such as attention mechanisms to enhance context understanding.

### 3. Application layer

The application layer comprises the Command-Line Interface (CLI) to accept user inputs, process articles, and deliver results in real-time.

Additionally, it will also integrate modules like SHAP (Shapley Additive exPlanations) and/or LIME (Local Interpretable Model-Agnostic Explanations) for transparency. As these modules provide users with insights into model decisions, the added transparency enhances user trust and interpretability.

4. Evaluation & monitoring layer

This layer comprises the evaluation framework that will track performance metrics (i.e. precision, recall, F1 score, ROC-AUC) during training and deployment phases.

This layer will also include error analysis tools to identify common misclassifications and potential biases to guide iterative improvements, and tools to log system performance metrics to ensure reliability and scalability.

5. Deployment layer

The solution will be packaged into a deployable Python package, with dependencies managed through Conda.

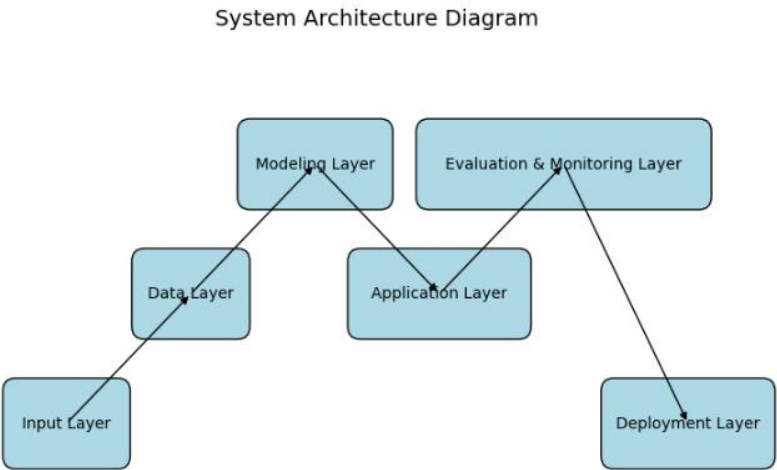


Figure 8: System architecture diagram

[3.5] Technologies & Methods

Programming Language	Python
Key Libraries	<b>Scikit-learn</b> – tools for traditional machine learning <b>Transformers (Hugging Face)</b> – fine-tuning BERT for contextual NLP tasks <b>Pandas</b> – for data cleaning & manipulation <b>NumPy</b> – for numerical computation <b>Matplotlib</b> – for data visualisation <b>Click</b> – to simplify CLI creation & improve user interaction
Evaluation Techniques	Cross-validation ROC curves SHAP/LIME for explainability

Figure 9: Table of technologies & methods used



# [3.6] Work Plan

## Milestones

*The position column, charts should be ordered sequentially.*

*Enter the date for a milestone in this column.*

*Enter a milestone description in this column. These descriptions will appear in the chart.*

No.	Position	Date	Milestone
1	1	16/12/2024	Preliminary Report Submission
2	2	7/1/2025	Jan: Monthly Check In (approx)
3	3	27/1/2025	Draft Report Submission
4	4	10/2/2025	Feb: Monthly Check In (approx)
4	4	3/3/2025	Mar: Monthly Check In (approx)
4	4	10/3/2025	Written Examination
5	5	24/3/2025	Final Report Submission

Figure 10: Project milestones

## Tasks

*Enter the start date for each task below. For best results sort this column in ascending order.*

*Enter the end date for each task or activity below, in this column.*

*Enter tasks and/or activities in this column.*

No.	Start Date	End Date	Task
1	16/12/2024	16/12/2024	Preliminary Report Submission
2	17/12/2024	13/1/2025	Dataset preparation & pre-processing pipeline
3	2/1/2025	13/1/2025	Refine project scope & finalise methodologies based on supervisor and grader feedback
4	8/1/2025	26/1/2025	Implement all models
5	15/1/2025	16/1/2025	Evaluate models using small test data subset
6	16/1/2025	23/1/2025	Refine models
7	1/1/2025	16/1/2025	Draft 1: Preliminary Report
8	24/1/2025	9/2/2025	Continue refining models
9	28/1/2025	9/2/2025	Document interim results & challenges between interim report completion and supervisor check-in
10	11/2/2025	2/3/2025	Integrate feedback
11	11/2/2025	25/2/2025	Finalise pipeline
12	25/2/2025	9/3/2025	Simple front-end
13	9/3/2025	23/3/2025	Final Report
14	11/3/2025	23/3/2025	Final tests & fixes

Figure 11: Project tasks

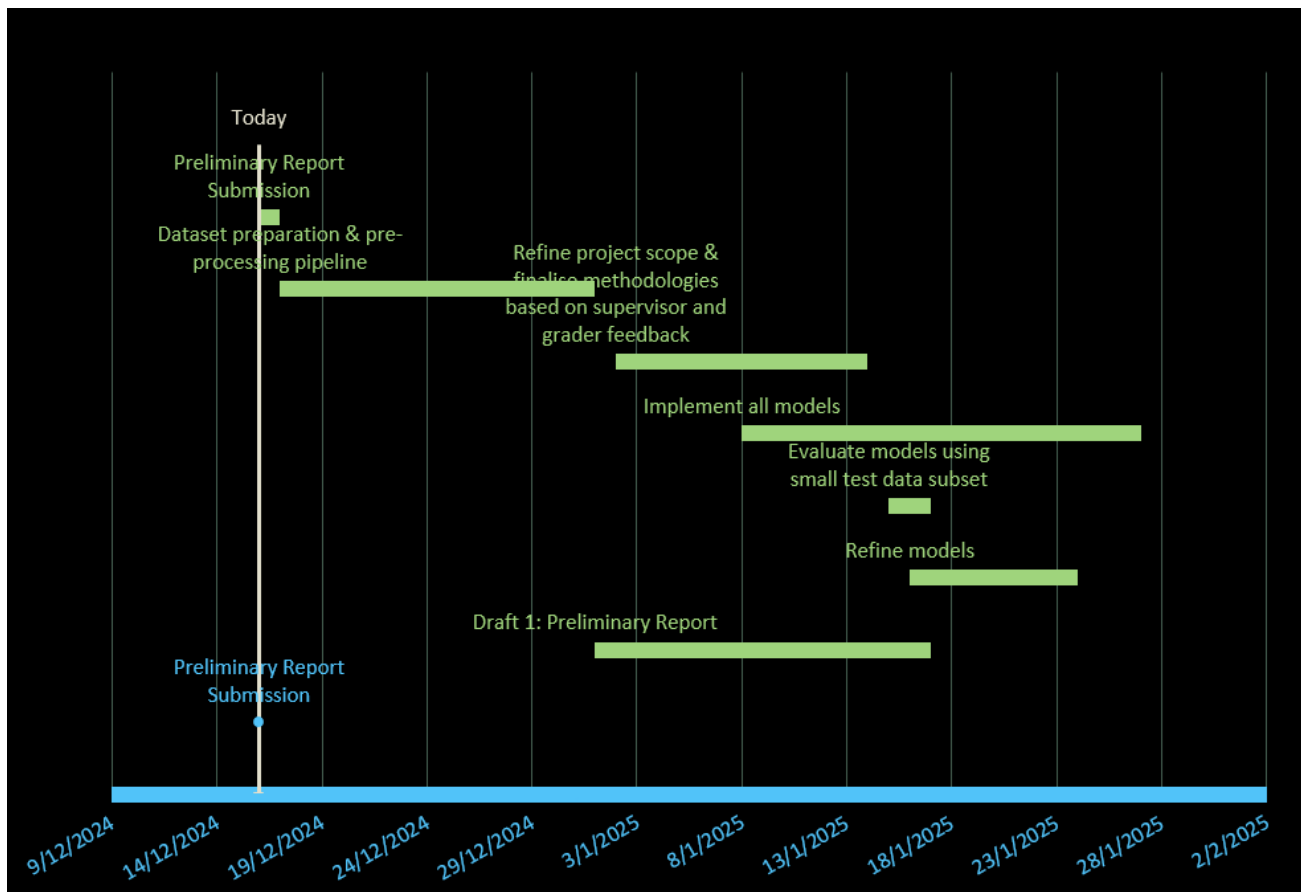


Figure 12: Gantt chart for project timeline

# [3.7] Testing & Evaluation

## [3.71] Test Plan

### 1. **Functional testing**

- Verify end-to-end functionality of the system, from input to output.
- Test BERT model on separate test dataset to evaluate metrics such as accuracy, precision, recall, F1 score.
- Additionally, use confusion matrices to analyse false positives/negatives.

### 2. **User testing**

- Collect feedback from target users on usability & effectiveness
- Conduct surveys with target users assess ease of use, clarity of results & overall satisfaction.
- Distribute a SUS (System Usability Scale) questionnaire to quantitatively evaluate usability
- Include task-based performance metrics (e.g. time on task, error rates) to measure efficiency.

### 3. **Security testing**

- Perform input validation tests to test against SQL injection and other vulnerabilities
- Ensure the application handles malformed or adversarial inputs gracefully

## [3.72] Evaluation Metrics

### Model Evaluation Metrics

1. Accuracy
2. Precision
3. Recall (Sensitivity)
4. F1 score
5. ROC-AUC Curve

### User Evaluation Metrics

1. System Usability Score (SUS) score – to assess overall usability
2. Interpretability Score – users rate the clarity of flagged explanations on a 1-5 scale
3. Time on task – to measure how efficiently users complete fake news identification tasks
4. Error rate – to track user misunderstandings of system outputs
5. Net Promoter Score (NPS) – to assess overall user satisfaction, and likelihood of recommending the system

## **[4.0] Implementation**

## **[5.0] Evaluation**

## **[6.0] Conclusion**

## **[7.0] Appendix**

### **[7.1] Work Done (versioning log)**

#### **[7.10] Feature Prototype (v0)**

##### **[7.101] Feature Prototype Overview**

The feature prototype demonstrates the implementation of fine-tuning a MobileBERT [16] model for the purpose of fake news detection.

MobileBERT is a variant of BERT chosen for its computational efficiency while retaining its performance in natural language processing tasks. The objective of this prototype is to train MobileBERT to classify political statements as either ‘fake’ or ‘real’ using a simplified, binarized version of the LIAR dataset’s labels.

This prototype’s implementation serves to assess the viability of employing MobileBERT in a fake news detection application with real-world use cases.

## [7.102] Dataset Preparation

The LIAR dataset comprises political statements categorised into six classes – “pants-fire”, “false”, “barely-true”, “half-true”, “mostly-true” and “true” [17] – that encapsulate a spectrum of truthfulness. In order to simplify the classification task for this prototype, these classes were merged into simplified binary labels:

1. Fake – “pants-fire”, “false”, “barely-true”
2. Real – “half-true”, “mostly-true”, “true”

The dataset was then prepared over multiple steps:

### 1. Data Loading

```
# Load dataset to pandas DataFrame

# import pandas \o/
import pandas as pd

# Load train, test, validation datasets
# for the purposes of this demo, we'll be using LIAR dataset :D
train_ds = "liar_dataset/train.tsv"
test_ds = "liar_dataset/test.tsv"
valid_ds = "liar_dataset/valid.tsv"

# now, i'll use pandas to read TSV files :D
# columns are as according to the README in liar_dataset directory :D

columns = [
    "id", "label", "statement", "subject", "speaker", "speaker_job_title",
    "state_info", "party_affiliation", "barely_true_counts", "false_counts",
    "half_true_counts", "mostly_true_counts", "pants_on_fire_counts", "context"
]

train_df = pd.read_csv(train_ds, sep='\t', names=columns)
test_df = pd.read_csv(test_ds, sep='\t', names=columns)
valid_df = pd.read_csv(valid_ds, sep='\t', names=columns)
```

Training, validation, and test splits were imported using pandas. This way, data would be cleanly separated, and evaluation of the model would be unbiased.



## 2. Label binarisation

```
# binarising labels!

# since the labels have multiple classes,
# for the sake of this feature prototype,
# i'll just simplify them to binary true/fake labels :)

# map labels to binary classes! :D
# 'pants-fire', 'false', 'barely-true' -> fake (0)
# others -> real (1)

def binarise(df):
    # validate expected labels exist before applying transformation!
    expected_labels = ["pants-fire", "false", "barely-true", "half-true", "mostly-true", "true"]
    unexpected_labels = set(df['label']) - set(expected_labels)
    if unexpected_labels:
        raise ValueError(f"Unexpected labels found: {unexpected_labels}")
    df['label'] = df['label'].apply(lambda x: 0 if x in ['pants-fire', 'false', 'barely-true'] else 1)
    return df

train_df = binarise(train_df)
test_df = binarise(test_df)
valid_df = binarise(valid_df)
```

To transform multi-class labels to binary labels ‘fake’ (0) and ‘real’ (1), a mapping function was implemented to map “pants-fire”, “false” and “barely-true” to ‘fake’ (0), and “half-true”, “mostly-true” and “true” to ‘real’ (1).

```
# print statement to check df structure!
print(train_df.head())
print(test_df.head())
print(valid_df.head())

# checking that all labels in dataset are valid
print("Unique labels in training data:", train_df['label'].unique())
assert set(train_df['label'].unique()) == {0, 1}, "Labels must be binary (0 or 1)"
```

Some basic validation work was also done to verify that all labels conformed to the expected binary format.

### 3. Tokenisation

Hugging Face's AutoTokenizer was employed to tokenise the text statements. This would ensure uniformity in length through truncation and padding to a maximum sequence length of 128 tokens.

This process preserves critical semantic features without compromising computational efficiency.

```
# tokenise statements
# i'll tokenise statements using Hugging Face's tokenizer!

# import autotokeniser
from transformers import AutoTokenizer

# Load tokenizer
tokenizer = AutoTokenizer.from_pretrained("google/mobilebert-uncased")

# tokenise data
def tokenise(df, tokenizer, max_length=128):
    return tokenizer(
        df['statement'].tolist(),
        truncation=True,
        padding=True,
        max_length=max_length,
        return_tensors="pt"
    )

train_encodings = tokenise(train_df, tokenizer)
test_encodings = tokenise(test_df, tokenizer)
valid_encodings = tokenise(valid_df, tokenizer)
```

## [7.103] Model Training

MobileBERT is then fine-tuned to perform binary classification in the task-specific setting on detecting fake news.

The fine-tuning pipeline included the following key components:

### 1. Dataset creation

```
# prepare our dataset for pytorch!
import torch

class LIARDataset(torch.utils.data.Dataset):
    def __init__(self, encodings, labels):
        self.encodings = encodings
        self.labels = labels

    def __len__(self):
        return len(self.labels)

    def __getitem__(self, idx):
        item = {key: val[idx] for key, val in self.encodings.items()}
        item['labels'] = torch.tensor(self.labels[idx])
        return item

train_dataset = LIARDataset(train_encodings, train_df['label'].tolist())
test_dataset = LIARDataset(test_encodings, test_df['label'].tolist())
valid_dataset = LIARDataset(valid_encodings, valid_df['label'].tolist())
```

The tokenised text and corresponding labels were encapsulated into PyTorch-compatible datasets to enable it to be seamlessly integrated with PyTorch's DataLoader, thus streamlining the training and evaluation workflows.

### 2. Model configuration

```
# now that our dataframes are tokenised,
# let's load pre-trained BERT.

from transformers import AutoModelForSequenceClassification

# Load our model :D
model = AutoModelForSequenceClassification.from_pretrained("google/mobilebert-uncased", num_labels=2)
```

A custom classification head with two output labels was initialised to adapt MobileBERT's architecture for binary classification. This modification would leverage MobileBERT's pre-trained language understanding capabilities while ensuring alignment with the binary classification task.

### 3. Training setup

```
# this codeblock for our dataloader! :D

# num_workers to use multiple cpu cores, pin_memory as we are training on GPU
train_loader = DataLoader(train_dataset, batch_size=16, shuffle=True)
test_loader = DataLoader(test_dataset, batch_size=16)
valid_loader = DataLoader(valid_dataset, batch_size=16)

# train_loader = DataLoader(train_dataset, batch_size=16, shuffle=True, num_workers=4, pin_memory=True)
# test_loader = DataLoader(test_dataset, batch_size=16, num_workers=4, pin_memory=True)
# valid_loader = DataLoader(valid_dataset, batch_size=16, num_workers=4, pin_memory=True)

# this codeblock for optimiser!
optimizer = AdamW(model.parameters(), lr=5e-5)

# this codeblock for scheduler!
num_training_steps = len(train_loader) * 10 # 10 epochs
lr_scheduler = get_scheduler("linear", optimizer=optimizer, num_warmup_steps=0, num_training_steps=num_training_steps)

# this codeblock for device config!
device = torch.device("cuda") if torch.cuda.is_available() else torch.device("cpu")
model.to(device)
```

The AdamW optimiser is employed with a learning rate of  $5e-5$  to ensure stable gradient updates.

A linear learning rate scheduler with no warm-up steps adapts the learning rate dynamically throughout training.

The training process uses a batch size of 16. Shuffling was applied to the training set to prevent learning biases and to improve generalisation.

## 4. Training loop

```
# this is our training loop.
# i will also implement early stopping here
from tqdm import tqdm

patience = 3 # stop training after 1 epoch without improvement!
best_loss = float('inf') # track the best loss achieved so far
patience_counter = 0

model.train()
for epoch in range(10): # max of 10 epochs, but we might stop earlier!
    epoch_loss = 0
    loop = tqdm(train_loader, leave=True)
    for batch in loop:
        batch = {k: v.to(device) for k, v in batch.items()}
        outputs = model(**batch)
        loss = outputs.loss
        loss.backward()

        optimizer.step()
        lr_scheduler.step()
        optimizer.zero_grad()

        epoch_loss += loss.item()
        loop.set_description(f"Epoch {epoch}")
        loop.set_postfix(loss=loss.item())

    # check for early stopping!
    if epoch_loss < best_loss:
        best_loss = epoch_loss
        patience_counter = 0
    else:
        patience_counter += 1
        if patience_counter >= patience:
            print("triggering early stopping!")
            break
```

The training process was capped to a maximum of 10 epochs. Early stopping was implemented based on validation loss to prevent overfitting.

Epoch 0: 100%	640/640	[01:43<00:00,	6.21it/s,	loss=0.681]
Epoch 1: 100%	640/640	[01:42<00:00,	6.23it/s,	loss=0.737]
Epoch 2: 100%	640/640	[01:43<00:00,	6.18it/s,	loss=0.715]
Epoch 3: 100%	640/640	[01:42<00:00,	6.27it/s,	loss=0.545]
Epoch 4: 100%	640/640	[01:44<00:00,	6.15it/s,	loss=0.749]
Epoch 5: 100%	640/640	[01:41<00:00,	6.33it/s,	loss=0.688]
Epoch 6: 100%	640/640	[01:42<00:00,	6.24it/s,	loss=0.555]
Epoch 7: 100%	640/640	[01:42<00:00,	6.27it/s,	loss=0.685]
Epoch 8: 100%	640/640	[01:42<00:00,	6.23it/s,	loss=0.681]
Epoch 9: 100%	640/640	[01:42<00:00,	6.24it/s,	loss=0.618]

Loss values were monitored per epoch, and gradients were computed and optimised after each batch. Early stopping criteria is based on validation loss, and training would stop if no improvement was observed over three consecutive epochs.

## [7.104] Evaluation

```
# evaluating the model!

# evaluate the model on test data :D
from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score, roc_auc_score

model.eval()
predictions, true_labels = [], []

with torch.no_grad():
    for batch in test_loader:
        batch = {k: v.to(device) for k, v in batch.items()}
        outputs = model(**batch)
        logits = outputs.logits
        predictions.extend(torch.argmax(logits, dim=-1).tolist())
        true_labels.extend(batch['labels'].tolist())

# metrics
accuracy = accuracy_score(true_labels, predictions)
precision = precision_score(true_labels, predictions, zero_division=0)

recall = recall_score(true_labels, predictions)
f1 = f1_score(true_labels, predictions)
roc_auc = roc_auc_score(true_labels, predictions)
```

The test set was used to assess the model's performance using standard classification metrics – accuracy, precision, recall, F1 score and ROC-AUC score.

```
# visualise results

from sklearn.metrics import ConfusionMatrixDisplay, confusion_matrix

cm = confusion_matrix(true_labels, predictions)
labels = train_df['label'].unique()
disp = ConfusionMatrixDisplay(confusion_matrix=cm, display_labels=labels)
disp.plot(cmap="Blues")
```

Additionally, a confusion matrix was generated to represent the distribution of true positives, true negatives, false positives and false negatives. The visualisation is meant to facilitate error analysis and highlight potential areas for model improvement.

## [7.105] Explainability

To enhance interpretability and build user trust, LIME (Local Interpretable Model-agnostic Explanations) was integrated into the prototype to identify and highlight the most influential input features contributing to the model’s predictions.

This approach aims to introduce transparency to the “black box” transformer models tend to be, which would allow stakeholders to understand the rationale behind its classifications.

```
# define a wrapper for model predictions
class ModelWrapper:
    def __init__(self, model, tokenizer, max_length=128, device='cpu'):
        self.model = model
        self.tokenizer = tokenizer
        self.max_length = max_length
        self.device = device

    def predict_proba(self, texts):
        # tokenize the input texts
        encodings = self.tokenizer(
            texts,
            truncation=True,
            padding=True,
            max_length=self.max_length,
            return_tensors="pt"
        )
        input_ids = encodings['input_ids'].to(self.device)
        attention_mask = encodings['attention_mask'].to(self.device)

        # get model predictions
        with torch.no_grad():
            outputs = self.model(input_ids, attention_mask=attention_mask)
            logits = outputs.logits

        # convert Logits to probabilities
        probs = torch.softmax(logits, dim=-1).cpu().numpy()
        return probs
```

A model wrapper class was created for LIME to interact with the MobileBERT model. It tokenises input text, and model outputs are transformed into interpretable probability scores for “fake” and “real” classes.

```
# initialize the LIME explainer
explainer = LimeTextExplainer(class_names=["Fake", "Real"])

# wrap model and tokenizer
wrapper = ModelWrapper(model, tokeniser, device=device)
```

LimeTextExplainer is initialised with class names “fake” and “real”.

```
# Select a sample from your test data to explain
sample_index = 0 # Example: explain the first test sample
sample_text = test_df['statement'].iloc[sample_index]
sample_label = test_df['label'].iloc[sample_index]

# generate explanation
explanation = explainer.explain_instance(
    sample_text,
    wrapper.predict_proba,
    num_features=10, # number of features to include in explanation
    top_labels=1 # focus on the top predicted label
)
```

A test sample is selected from the dataset to generate an explanation. LIME identifies the top 10 features that most influenced the model's prediction for the sample.

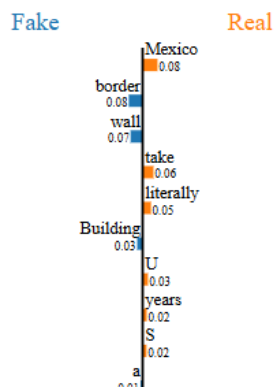
```
# display the explanation
print(f"Original Text: {sample_text}")
print(f"True Label: {'Real' if sample_label == 1 else 'Fake'}")

# show weights of top features for the predicted label
explanation.show_in_notebook(text=True)
```

LIME provides a weighted list of influential features for each explained instance.

Original Text: Building a wall on the U.S.-Mexico border will take literally years.  
True Label: Real

Prediction probabilities



Text with highlighted words

Building a wall on the U.S.-Mexico border will take literally years.



## [7.106] Results of Initial Implementation

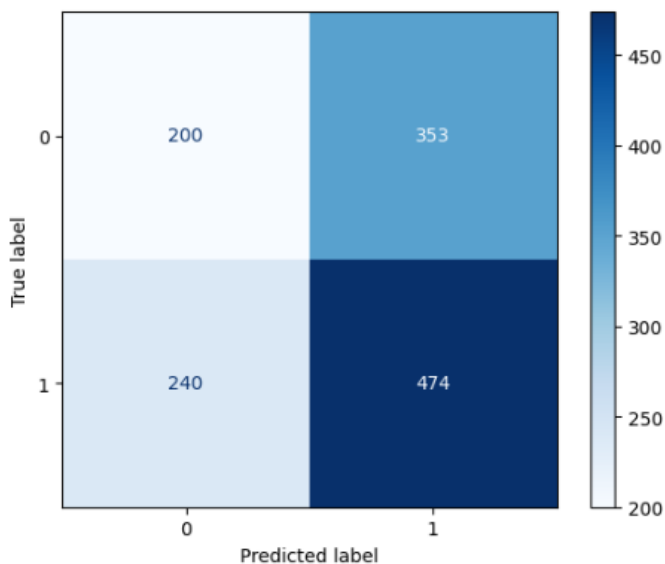
The prototype demonstrated the feasibility of fine-tuning MobileBERT for the task of fake news detection.

The key performance metrics achieved during evaluation are as follows:

```
Accuracy: 0.5320
Precision: 0.5732
Recall: 0.6639
F1 Score: 0.6152
ROC-AUC: 0.5128
```

Accuracy: 53.20%  
Precision: 57.32%  
Recall: 66.39%  
F1 Score: 61.52%  
ROC-AUC: 51.28%

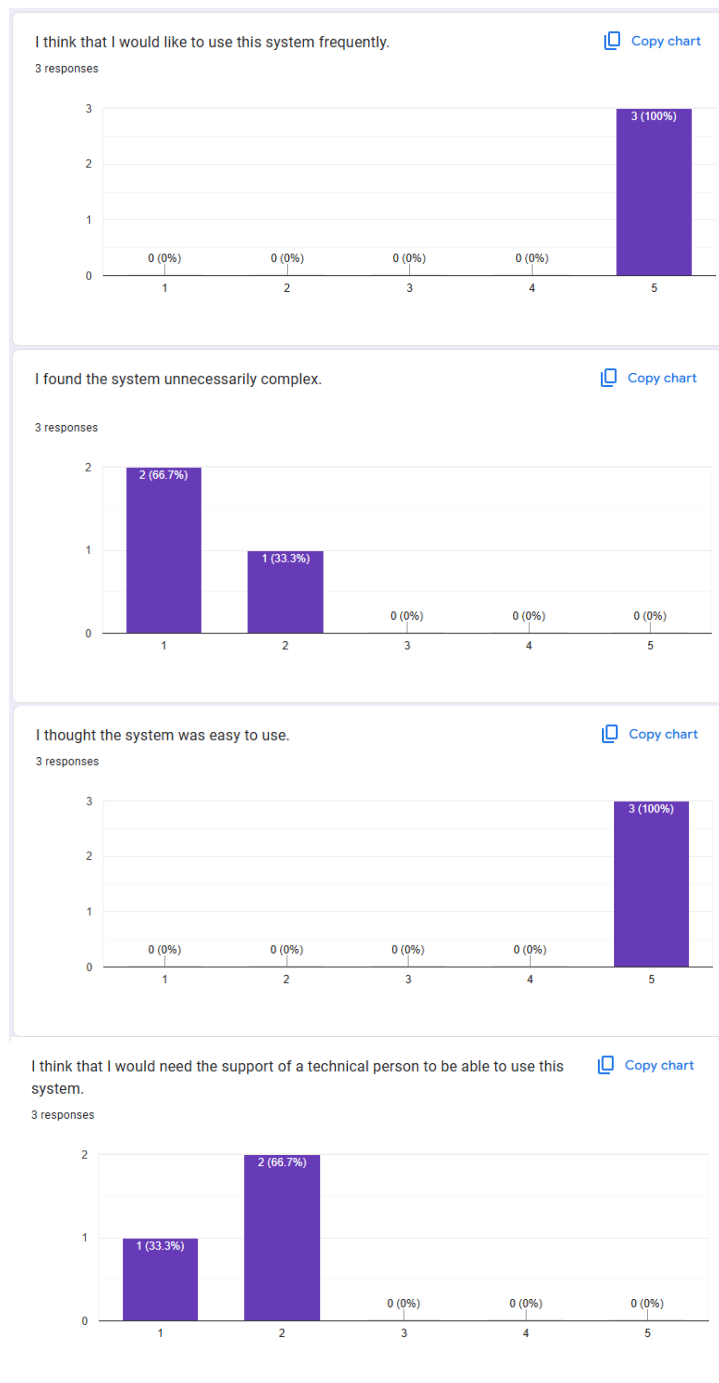
These values indicate that while the model is more accurate than randomly guessing, there is large room for improvement, especially in optimising prediction and recall.



The confusion matrix reveals that the model exhibits a tendency towards higher false positives, thus significantly affecting precision. The high recall score indicates the model is capable of capturing real statements, but improving precision will be a main goal moving forward in future iterations.

Additionally, a small group survey was conducted with a group of three student jourliasts.

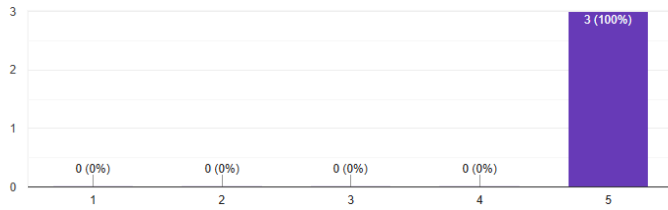
Participants filled in a System Usability Scale (SUS) questionnaire and were asked to rate how much they understood why statements were being flagged.



I found the various functions in this system were well integrated.

 [Copy chart](#)

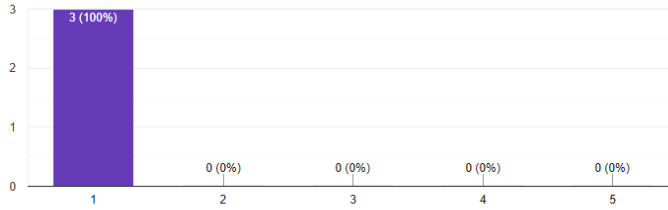
3 responses



I thought there was too much inconsistency in this system.

 [Copy chart](#)

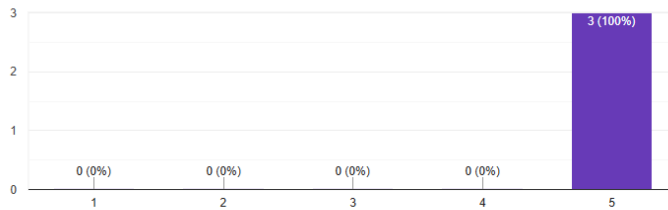
3 responses



I would imagine that most people would learn to use this system very quickly.

 [Copy chart](#)

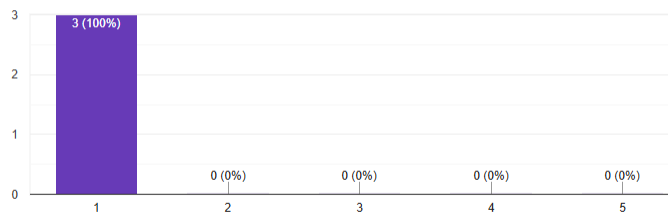
3 responses



I found the system very cumbersome to use.

 [Copy chart](#)

3 responses





From the results of the questionnaire, I can deduce that:

**1. The participants wanted more clarity as to why statements are flagged.**

This can be done by introducing a more human explainability element, like a short statement instead of just the LIME output, to be more conducive to non-technical users.

**2. The participants felt a CLI application was too intimidating/technical.**

Implementing a simple GUI may be a way to bridge the gap and instil confidence in non-technical users.

## [7.107] Further Iterations

As the original implementation did not perform exceptionally well according to model evaluation metrics, in subsequent iterations of the prototypes, more techniques were attempted to improve the performance of the model.

### [7.1071] Mixed Precision Training

The original training loop was worked on to implement mixed precision training using PyTorch's `torch.cuda.amp` module to enable efficient training by performing certain operations in lower 16-bit floating point precision, while maintaining stability in others with a 32-bit floating point precision.

The potential benefit of implementing mixed precision training is to reduce computation time by utilising tensor cores on GPUs, and lower memory usage so larger batch sizes or models would be able to fit in the same memory.

```
# gradient scaler for mixed precision
scaler = GradScaler()

patience = 3 # stop training after 3 epochs without improvement!
best_loss = float('inf') # track the best loss achieved so far
patience_counter = 0

model.train()
for epoch in range(10): # max of 10 epochs, but we might stop earlier!
    epoch_loss = 0
    loop = tqdm(train_loader, leave=True)

    for batch in loop:
        # move batch to device
        batch = {k: v.to(device) for k, v in batch.items()}

        # forward pass with mixed precision
        with autocast():
            outputs = model(**batch)
            loss = outputs.loss # calculate loss

        # backward pass with gradient scaling
        scaler.scale(loss).backward()

        # gradient clipping. i do this to avoid exploding gradients.
        clip_grad_norm_(model.parameters(), max_norm=1.0)

        # update model parameters
        scaler.step(optimizer) # apply scaled gradients to optimizer
        scaler.update() # update the scaler for next iteration
        optimizer.zero_grad() # reset gradients after updating
        lr_scheduler.step() # step the scheduler

    # accumulate epoch loss
    epoch_loss += loss.item()

    # progress bar logic
    loop.set_description(f"Epoch {epoch}")
    loop.set_postfix(loss=loss.item())

# check for early stopping
if epoch_loss < best_loss:
    best_loss = epoch_loss
    patience_counter = 0
else:
    patience_counter += 1
    if patience_counter >= patience:
        print("triggering early stopping!")
        break
```

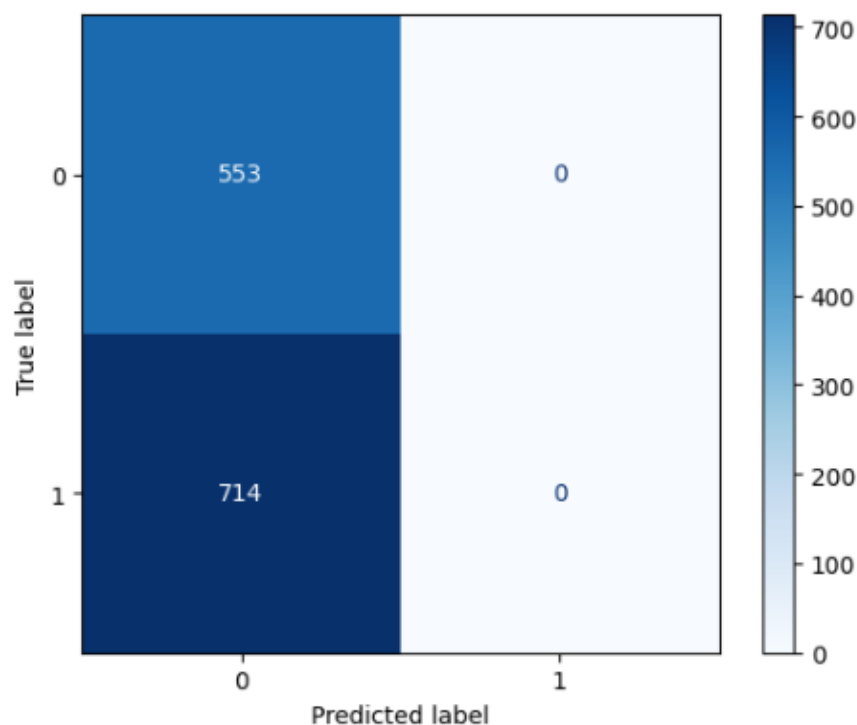
To implement mixed precision training,

1. A gradient scaler ‘GradScaler()’ is initialised.
2. ‘autocast’ is used to wrap the forward pass in the training loop.
3. The scaler is used to scale the loss before calling ‘backward’, and to handle the optimizer step and update.

```
Epoch 0: 100%|██████████████████████████████████████████████████████████████████████████████| 640/640 [01:40<00:00, 6.38it/s, loss=nan]
Epoch 1: 100%|██████████████████████████████████████████████████████████████████████████████| 640/640 [01:38<00:00, 6.50it/s, loss=nan]
Epoch 2: 100%|██████████████████████████████████████████████████████████████████████████████| 640/640 [01:34<00:00, 6.78it/s, loss=nan]
triggering early stopping!
```

Accuracy: 0.4365  
Precision: 0.0000  
Recall: 0.0000  
F1 Score: 0.0000

```
<sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x1e5ffe6b090>
```



The results of the model trained with mixed precision training indicated a serious issue in the training process, as evidenced by the NaN loss during training across all epochs. This suggests that the model encountered numerical instability during the training loop, preventing the model from learning any meaningful patterns.

Looking at the evaluation metrics, the precision, recall and F1 score are all zero, indicating that the model trained with mixed precision training did not predict any class correctly. The confusion matrix indicated that all predictions were consistently 'o', regardless of the true label.

Attempts were made to reduce the learning rate, monitor gradients and adjust gradient clipping. However, unfortunately, the same issues persisted and time constraints restricted any further troubleshooting and debugging to make mixed precision training work at this time.

## [7.1072] Gradient Accumulation and Clipping

Gradient accumulation enables training with larger effective batch sizes without increasing memory requirements by accumulating gradients over multiple batches before performing an update, instead of updating the model parameters after each batch.

The potential improvement of implementing this is improved model stability, convergence, and numerical stability due to the prevention of exploding gradients. This also improves resource efficiency, as GPU memory usage is reduced as compared to directly increasing the batch size.

```
model.train()
for epoch in range(10): # max of 10 epochs, but we might stop earlier!
    epoch_loss = 0
    loop = tqdm(train_loader, leave=True)
    optimizer.zero_grad() # reset gradients at the start of the epoch

    for i, batch in enumerate(loop):
        batch = {k: v.to(device) for k, v in batch.items()}

        # forward pass
        outputs = model(**batch)
        loss = outputs.loss / accumulation_steps # normalize loss for accumulation
        loss.backward() # backpropagation

        # perform optimization step after accumulation_steps
        if (i + 1) % accumulation_steps == 0 or (i + 1) == len(loop):
            # gradient clipping
            torch.nn.utils.clip_grad_norm_(model.parameters(), max_norm)

            # update model parameters
            optimizer.step()
            lr_scheduler.step()
            optimizer.zero_grad() # after update, reset gradients

        # accumulate epoch loss
        epoch_loss += loss.item() * accumulation_steps # scale loss back to original value

        # progress bar update logic
        loop.set_description(f"Epoch {epoch}")
        loop.set_postfix(loss=loss.item() * accumulation_steps)

    # early stopping logic
    if epoch_loss < best_loss:
        best_loss = epoch_loss
        patience_counter = 0
    else:
        patience_counter += 1
        if patience_counter >= patience:
            print("triggering early stopping!")
            break
```

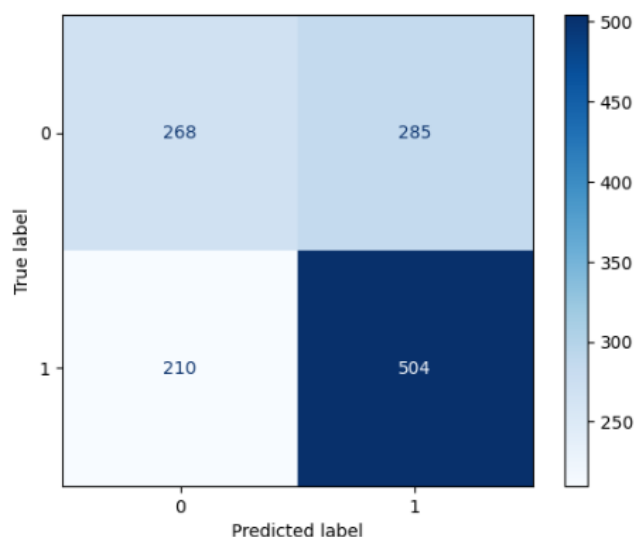
The original training loop was adapted to implement gradient accumulation and clipping, through the following key changes:

1. The loss is divided by accumulation steps to ensure that gradients are appropriately scaled during accumulation,
2. Gradients are clipped before the optimiser step to prevent exploding gradients. 'max\_norm' sets the maximum gradient norm.
3. The model parameters are updated only after 'accumulation\_steps' batches.
4. 'optimizer.zero\_grad()' is called after the optimiser step to reset gradients.

Epoch 0: 100%	640/640 [01:19<00:00, 8.04it/s, loss=0.775]
Epoch 1: 100%	640/640 [01:20<00:00, 7.94it/s, loss=0.757]
Epoch 2: 100%	640/640 [01:19<00:00, 8.05it/s, loss=0.548]
Epoch 3: 100%	640/640 [01:19<00:00, 8.07it/s, loss=0.767]
Epoch 4: 100%	640/640 [01:20<00:00, 7.94it/s, loss=0.614]
Epoch 5: 100%	640/640 [01:22<00:00, 7.76it/s, loss=0.447]
Epoch 6: 100%	640/640 [01:21<00:00, 7.88it/s, loss=0.436]
Epoch 7: 100%	640/640 [01:20<00:00, 7.96it/s, loss=0.275]
Epoch 8: 100%	640/640 [01:17<00:00, 8.30it/s, loss=0.141]
Epoch 9: 100%	640/640 [01:16<00:00, 8.34it/s, loss=0.0166]

During training, the loss steadily decreases across epochs, starting from 0.775 in epoch 0 and finally converging to 0.0166, indicating that the model is learning the training data effectively.

Accuracy: 0.6093  
Precision: 0.6388  
Recall: 0.7059  
F1 Score: 0.6707  
ROC-AUC: 0.5953



The accuracy value of 60.93% has improved compared to the initial implementation (53.2%), as has the precision score (63.88%) compared to the initial implementation (57.32%). The recall and F1 score have also increased, though the current values and confusion matrix indicates the model fails to correctly classify more true instances, and suggests that while gradient accumulation and clipping have improved training stability, the balance between precision and recall is not fully optimised yet.

While the significant decrease in training loss has translated to better evaluation metrics, overfitting may be occurring as the loss has converged too aggressively to near-zero, suggesting that the model may have potentially memorised the training data while not generalising well to unseen data.

Attempts to tune the learning rate, increase the patience for early stopping, and using weighted loss functions were made to improve on the results, however the evaluation metrics did not see tangible improvements, with time constraints prohibiting further troubleshooting.



## [7.108] Future Improvements from this point

While time constraints limited the amount of iterative work done on the prototype, several possible areas for improvement were identified for future work.

### **1. Addressing class imbalance**

The dataset had an uneven distribution of labels. This potentially skewed the model's performance, and can be addressed by techniques such as oversampling or applying class weights during training to mitigate bias and improve the model's robustness.

### **2. Precision improvements**

The high rate of false positives indicates a need to refine decision boundaries, possibly by adding more nuanced examples of fake statements to the dataset, or applying advanced loss functions to penalise false positives.

### **3. Expanding evaluation metrics**

Incorporating additional evaluation metrics and producing more comprehensive user evaluation metrics could provide more comprehensive insights and a better understanding of the model's reliability. This would be critical when the model is deployed for real-world applications, where confidence in predictions is crucial.

### **4. More human-like explainability features**

To instil confidence and foster user trust, explainability features need to be understood by users. Implementing short and easily-understandable textual explanations as to why content is flagged would be more useful for non-technical users who may feel alienated by LIME's charts.

## [7.11] Version 1

### [7.111] Overview & Aims

The goal of Version 1 is to simply improve the functional performance of the Initial Prototype.

This is done through changing the training framework. While the Initial Prototype used a custom PyTorch training loop where each training step was handled manually, this version uses Hugging Face's Trainer API that abstracts away specific implementation details.

Additionally, some changes were made to the Evaluation pipeline. The metrics computation function is now directly implemented into the training loop, so that it is possible for me to assess the model's performance at each training epoch.

As the changes in this version are only done to back-end logic and there are no changes to user presentation, I will not be conducting user evaluation/testing for this version of the product as I do not expect any meaningful insights to emerge.

## [7.112] Changes from Previous Version (Feature Prototype)

The main change from the previous version (Feature Prototype) is as follows:

### 1. New training logic

```
# Initialize Trainer
trainer = Trainer(
    model=model,
    args=training_args,
    train_dataset=train_dataset,
    eval_dataset=valid_dataset,
    compute_metrics=compute_metrics,
)

# Train the model
trainer.train()

# Save the final model and tokenizer
model_path = "mobilebert_fake_news_final"
trainer.save_model(model_path)
tokenizer.save_pretrained(model_path)
print(f"Model and tokenizer saved to {model_path}")
```

A new training loop has been implemented that leverages Hugging Face's Trainer API.

This reduces implementation complexity as the Trainer API handles many details that would normally require manual implementation (e.g. manual hyperparameter tuning), making the codebase more concise, and easy to maintain.

Additionally, as the Trainer API is more widely tried and tested than a proprietary/custom implementation, it should be less prone to implementation errors, making it a more robust solution.

Additionally, evaluation metrics are displayed at each epoch so I could assess the model's performance at each training epoch.

[6400/6400 29:36, Epoch 10/10]

Epoch	Training Loss	Validation Loss	Accuracy	Precision	Recall	F1	Roc Auc
1	39373.743800	0.652544	0.629283	0.615942	0.763473	0.681818	0.623620
2	0.630000	0.682776	0.613707	0.590717	0.838323	0.693069	0.604227
3	0.676600	0.711157	0.626947	0.632168	0.676647	0.653651	0.624849
4	0.660500	0.952788	0.591900	0.571856	0.857784	0.686228	0.580678
5	0.331300	1.169161	0.626947	0.620999	0.726048	0.669427	0.622764
6	0.368400	1.930946	0.599688	0.589953	0.755988	0.662730	0.593091
7	0.108300	3.072587	0.602025	0.595849	0.730539	0.656355	0.596601
8	0.058800	4.806461	0.610592	0.598824	0.761976	0.670619	0.604202
9	0.050600	4.869033	0.609034	0.606138	0.709581	0.653793	0.604790
10	0.046500	5.081384	0.610592	0.607417	0.711078	0.655172	0.606351

[7.113] Training Arguments & Evaluation Results

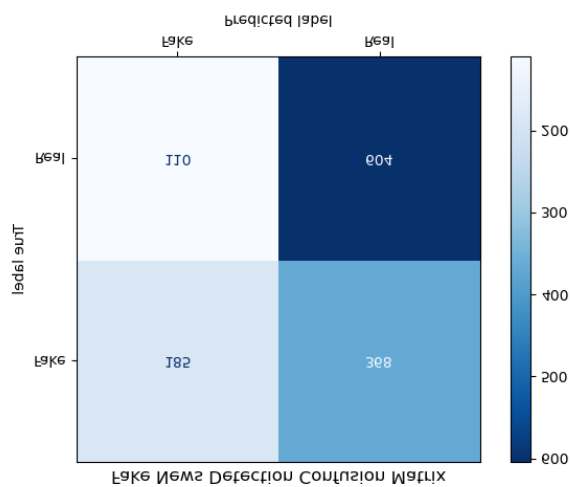
To try and improve performance, training arguments were changed and the evaluation results were assessed to try and determine the most optimal training arguments for this model.

1. First notable set of arguments (Argument Set 1)

```
# Configure training arguments
training_args = TrainingArguments(
    output_dir="mobilebert_fake_news",
    eval_strategy="epoch",
    save_strategy="epoch",
    logging_strategy="epoch",
    per_device_train_batch_size=16,
    per_device_eval_batch_size=16,
    num_train_epochs=10,
    weight_decay=0.01,
    load_best_model_at_end=True,
    metric_for_best_model="f1",
    push_to_hub=False,
    report_to="none" # disable wandb or other reporting to simplify
)
```

Epoch	Training Loss	Validation Loss	Accuracy	Precision	Recall	F1	Roc Auc
1	39373.743800	0.652544	0.629283	0.615942	0.763473	0.681818	0.623620
2	0.630000	0.682776	0.613707	0.590717	0.838323	0.693069	0.604227
3	0.676600	0.711157	0.626947	0.632168	0.676647	0.653651	0.624849
4	0.660500	0.952788	0.591900	0.571856	0.857784	0.686228	0.580678
5	0.331300	1.169161	0.626947	0.620999	0.726048	0.669427	0.622764
6	0.368400	1.930946	0.599688	0.589953	0.755988	0.662730	0.593091
7	0.108300	3.072587	0.602025	0.595849	0.730539	0.656355	0.596601
8	0.058800	4.806461	0.610592	0.598824	0.761976	0.670619	0.604202
9	0.050600	4.869033	0.609034	0.606138	0.709581	0.653793	0.604790
10	0.046500	5.081384	0.610592	0.607417	0.711078	0.655172	0.606351

Metric	loss	accuracy	precision	recall	f1	roc_auc	runtime	samples_per_second	steps_per_second	epoch
Value	0.6625	0.6227	0.6214	0.8459	0.7165	0.5902	6.0581	209.143	13.206	10.0

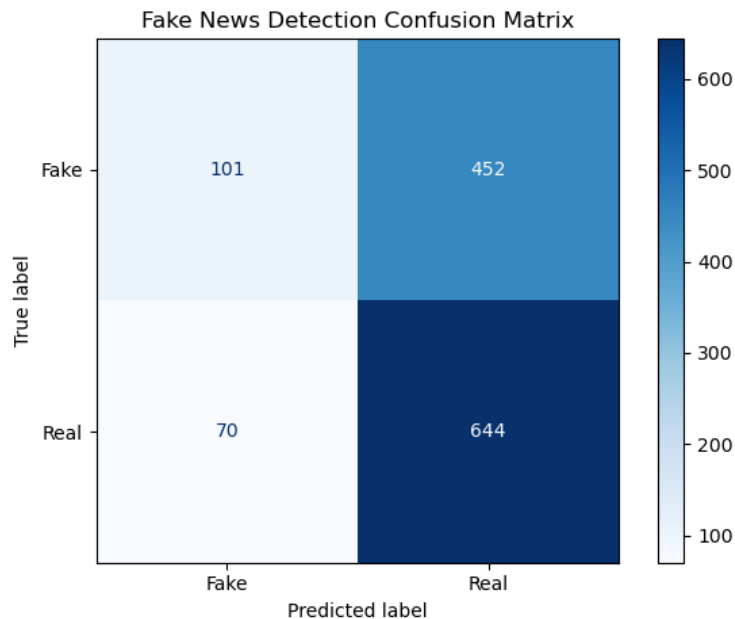


## 2. Second notable set of arguments (Argument Set 2)

```
# Configure training arguments
training_args = TrainingArguments(
    output_dir="mobilebert_fake_news",
    eval_strategy="epoch",
    save_strategy="epoch",
    logging_strategy="epoch",
    per_device_train_batch_size=16,
    per_device_eval_batch_size=16,
    num_train_epochs=4, # change to 4 from 10
    weight_decay=0.1, # change to 0.1 from 0.01
    load_best_model_at_end=True,
    metric_for_best_model="f1",
    push_to_hub=False,
    report_to="none", # disable wandb or other reporting to simplify
)
```

Epoch	Training Loss	Validation Loss	Accuracy	Precision	Recall	F1	Roc Auc
1	7183.722700	0.688509	0.580997	0.559415	0.916168	0.694665	0.566850
2	0.634100	0.672759	0.619159	0.596340	0.829341	0.693801	0.610288
3	0.568800	0.687449	0.633956	0.621622	0.757485	0.682861	0.628743
4	0.477500	0.745430	0.623832	0.615770	0.736527	0.670757	0.619075

Metric	loss	accuracy	precision	recall	f1	roc_auc	runtime	samples_per_second	steps_per_second	epoch
Value	0.6738	0.588	0.5876	0.902	0.7116	0.5423	5.8943	214.954	13.572	4.0



## [7.114] Evaluating Differences Between Both Sets of Arguments

Reviewing the evaluation metrics from the models trained on both sets of arguments above reveals some fundamental differences:

### 1. Trade-off between precision & recall

Argument set 2 sacrifices precision (-5.4%) for an improvement in recall (+6.6%).

This roughly means that Argument Set 1 trained a model that is more balanced and careful about making “real news” predictions, whereas Argument Set 2 is more aggressive to classify content as “real news”.

The confusion matrix for the model trained with Argument Set 2 confirms this tendency – 452 fake news items were incorrectly classified as real (high false positive rate), whereas only 70 real news items were incorrectly classified, suggesting that set 2 has a strong bias towards “real” classification.

### 2. Training efficiency

The model trained with Argument Set 2 achieved comparable F1 score (0.7% lower) while only being trained across 4 epochs (as opposed to 10 in Argument Set 1). This indicates that despite fewer training epochs, higher weight decay effectively prevented overfitting.

I suspect this is likely to have important implications for computational efficiency, as this demonstrates that strong regularisation can achieve comparable results with less training.

While the higher accuracy and ROC-AUC of the model trained by Argument Set 1 suggests that it has an overall better ability to discriminate between real and fake news, the F1 scores are quite close, indicating that both approaches have merit.

## [7.115] Conclusions and Future Work from this point

After making many changes to the work, I believe that while further hyperparameter tuning may improve the model, it would make more sense for my next step to be implementing a more robust data pipeline that leverages more datasets.

No matter what I do, the model overfits far too quickly, leading me to believe there is either insufficient data volume in the LIAR dataset alone, or “data leakage” – superficial patterns in the dataset that allow the model to “cheat” during training (e.g. certain sources or writing styles or vocabulary choices may strongly correlate with labels in ways that do not generalise to real-world scenarios).

As such, the most immediate future work that I will implement is a more robust data pipeline so that there is more data in the hopes that the model will learn generalisable features as opposed to memorising training examples.

## [7.12] Version 2

### [7.121] Overview & Aims

In this stage of work, I aim to implement a more robust data pipeline. This is done by introducing the ISOT Fake News Dataset [new2], a compilation of several thousand fake news and truthful articles, obtained from different legitimate news sites and sites flagged as unreliable by Politifact.com. [new2]

Additionally, motivated by evaluation results after immediately introducing the ISOP Fake News Dataset into the pipeline (along with the LIAR dataset), dataset balancing features were added, and hyperparameter tuning was introduced.

User evaluation is not performed for this version, as no changes were made to user-facing elements of the application, hence I do not expect any meaningful insights to emerge.

### [7.122] Integrating the ISOT dataset in the data layer

In order to expand the data layer to include more datasets, the implementation approach needs to be more sophisticated as a simple concatenation approach would fail to address that both the ISOT and LIAR datasets have fundamentally different structures: the LIAR dataset contains short claims, whereas the ISOT dataset contains full news articles with title/text separations.

I identified that a new pipeline at this stage would have several key components:

#### **1. Dataset-specific preprocessing**

I would need to create pre-processing modules for each dataset. These modules would handle dataset-specific formats and extract valuable features from each dataset, before transforming dataset-specific fields into a common schema.

#### **2. Unified feature representation**

I would need to design a common schema that would be filled with information that both datasets have (e.g. core text content, real/fake binary label, etc.)

#### **3. Evaluation that reports overall performance across all datasets, and dataset-specific performance**

This would be important to check that the model is not overly-generalising to one dataset – for example, if the model performed well with ISOT dataset but did not perform well on the LIAR dataset, I would need to be able to identify this situation.

## [7.123] Initial Successful Execution: Implementation Details

For this section, while I have performed lots of trial and error, I will only document the two most significant executions in the report as there was little value/insights gleaned from other execution runs.

### 1. Dataset-specific pre-processing modules

```
# module for Loading and preprocessing LIAR dataset
# arguments:
# train_path, valid_path, test_path - paths to LIAR dataset files
# return - tuple of (train_df, valid_df, test_df) with preprocessed data
def load_liar_dataset(train_path, valid_path, test_path):
    # step 1: define column names for LIAR dataset
    column_names = ['id', 'label', 'statement', 'subject', 'speaker', 'job_title', 'state_info',
                    'party_affiliation', 'barely_true_counts', 'false_counts', 'half_true_counts',
                    'mostly_true_counts', 'pants_on_fire_counts', 'context']

    # step 2: Load datasets
    train_df = pd.read_csv(train_path, sep='\t', names=column_names)
    valid_df = pd.read_csv(valid_path, sep='\t', names=column_names)
    test_df = pd.read_csv(test_path, sep='\t', names=column_names)

    # step 3: convert Labels to binary (0 for fake, 1 for real)
    # map 'pants-fire', 'false', 'barely-true' to fake (0)
    # map 'half-true', 'mostly-true', 'true' to real (1)
    label_map = {
        'pants-fire': 0, 'false': 0, 'barely-true': 0,
        'half-true': 1, 'mostly-true': 1, 'true': 1
    }

    # apply mapping to each dataset
    train_df['label'] = train_df['label'].map(label_map)
    valid_df['label'] = valid_df['label'].map(label_map)
    test_df['label'] = test_df['label'].map(label_map)

    # step 4: add source identifier
    train_df['source'] = 'liar'
    valid_df['source'] = 'liar'
    test_df['source'] = 'liar'

    return train_df, valid_df, test_df

# module for Loading and preprocessing ISOT dataset
# arguments:
# fake_path - path to Fake.csv
# real_path - path to Real.csv
# test_split - fraction of data to use for testing (default: 0.2)
# valid_split - fraction of data to use for validation (default: 0.1)
# return - tuple of (train_df, valid_df, test_df) with preprocessed data
def load_isot_dataset(fake_path, real_path, test_split=0.2, valid_split=0.1):
    # step 1: Load fake and real news
    fake_df = pd.read_csv(fake_path)
    real_df = pd.read_csv(real_path)

    # step 2: add binary Labels (0 for fake, 1 for real)
    fake_df['label'] = 0
    real_df['label'] = 1

    # step 3: add source identifier
    fake_df['source'] = 'isot'
    real_df['source'] = 'isot'

    # step 4: combine text and title for ISOT dataset
    fake_df['statement'] = fake_df['title'] + " " + fake_df['text']
    real_df['statement'] = real_df['title'] + " " + real_df['text']

    # step 5: combine datasets
    combined_df = pd.concat([fake_df, real_df], ignore_index=True)

    # step 6: shuffle data
    combined_df = combined_df.sample(frac=1, random_state=42).reset_index(drop=True)

    # step 7: split into train, validation, and test sets
    test_size = int(len(combined_df) * test_split)
    valid_size = int(len(combined_df) * valid_split)

    test_df = combined_df[:test_size]
    valid_df = combined_df[test_size:test_size+valid_size]
    train_df = combined_df[test_size+valid_size:]

    return train_df, valid_df, test_df
```



## 2. Unified Dataset Class to handle both LIAR and ISOT Datasets

```
# unified dataset class that handles both LIAR and ISOT datasets
class FakeNewsDataset(Dataset):
    # initialize dataset with dataframe and tokenizer
    # arguments:
    # df - dataframe with unified schema
    # tokenizer - tokenizer to use for text encoding
    # max_length - maximum sequence length (default: 128)
    def __init__(self, df, tokenizer, max_length=128):
        self.df = df
        self.texts = df['text'].tolist()
        self.labels = df['label'].tolist()
        self.sources = df['source'].tolist()
        self.tokenizer = tokenizer
        self.max_length = max_length

    # return number of samples in dataset
    def __len__(self):
        return len(self.texts)

    # get item at specified index
    # argument:
    # idx - index of item to get
    # return - dictionary with encoded inputs and label
    def __getitem__(self, idx):
        text = self.texts[idx]
        label = self.labels[idx]
        source = self.sources[idx]

        # encode text using tokenizer
        encoding = self.tokenizer(
            text,
            max_length=self.max_length,
            padding='max_length',
            truncation=True,
            return_tensors='pt'
        )

        # squeeze tensor dimensions
        return {
            'input_ids': encoding['input_ids'].squeeze(),
            'attention_mask': encoding['attention_mask'].squeeze(),
            'labels': torch.tensor(label, dtype=torch.long),
            'source': source # include source identifier
        }
```

The rest of the pipeline remained functionally similar, with some code re-factoring to tidy up the code and to ensure it would be able to support the integration of both the LIAR and ISOT Fake News datasets.

[7.124] Initial Successful Execution: Evaluation

1. Training

Epoch	Training Loss	Validation Loss	Accuracy	Precision	Recall	F1	Roc Auc
1	4034.274900	0.150010	0.915815	0.918870	0.906395	0.912590	0.915538
2	0.201800	0.152343	0.914602	0.875326	0.960700	0.916028	0.915959
3	0.137300	0.182318	0.915988	0.899517	0.930690	0.914838	0.916421
4	0.109200	0.266266	0.911831	0.884745	0.940693	0.911861	0.912680
5	0.066000	0.408894	0.910272	0.894227	0.924259	0.908995	0.910683

2. Performance on Test Set

Overall Performance Metrics:

Metric	loss	accuracy	precision	recall	f1	roc_auc	runtime	samples_per_second	steps_per_second	epoch
Value	0.0845	0.9547	0.9322	0.9776	0.9543	0.9554	58.9909		173.688	10.866 5.0

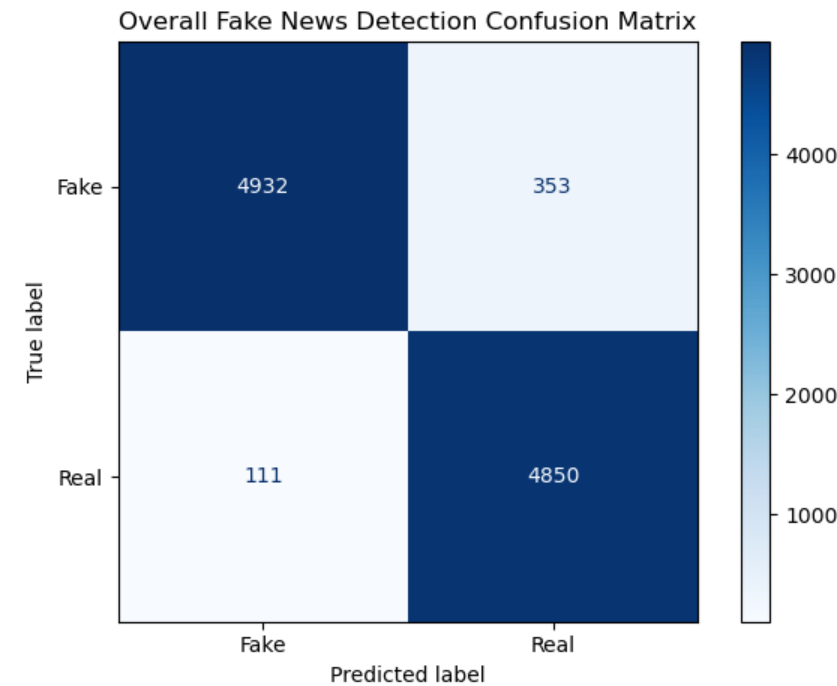
ISOT Dataset Performance Metrics:

Metric	accuracy	precision	recall	f1	roc_auc
Value	0.9998	0.9998	0.9998	0.9998	0.9998

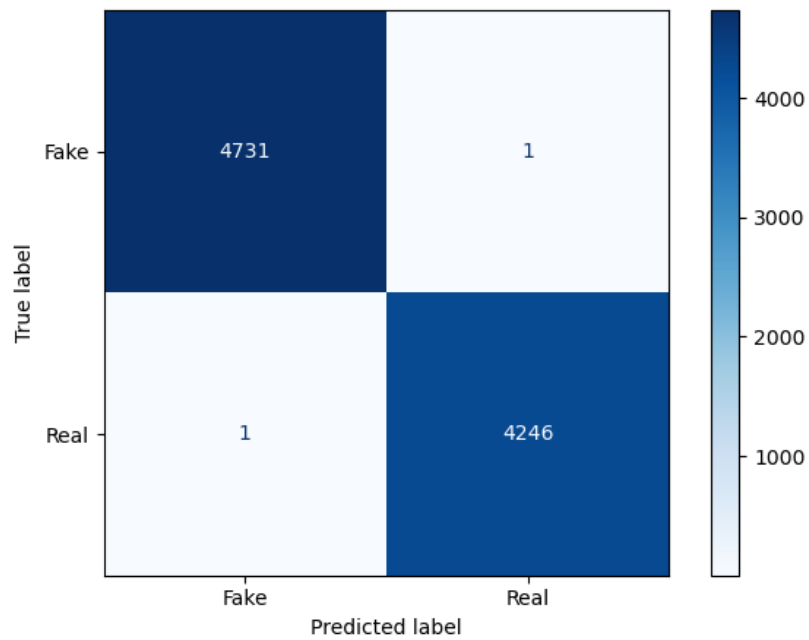
LIAR Dataset Performance Metrics:

Metric	accuracy	precision	recall	f1	roc_auc
Value	0.6354	0.6318	0.8459	0.7234	0.6047

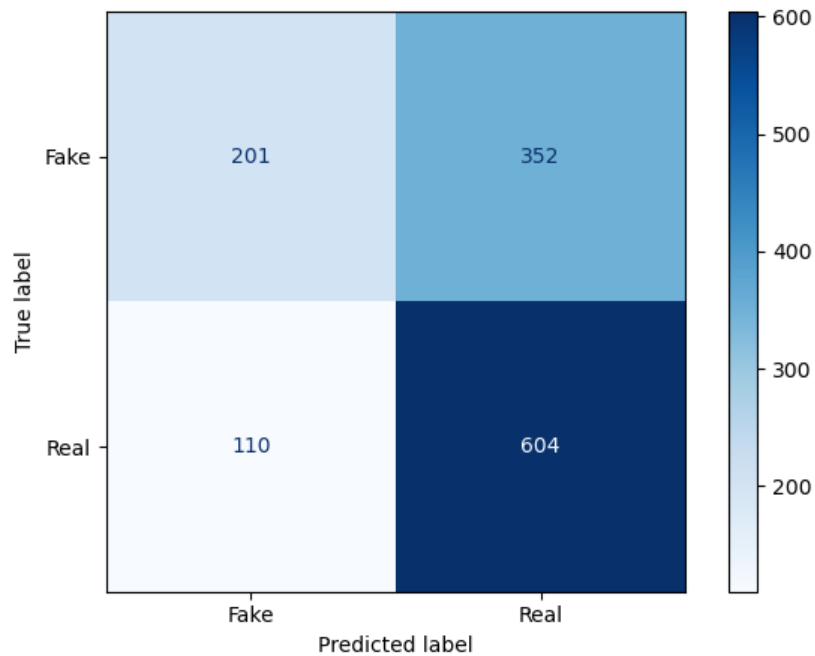
3. Confusion Matrices



ISOT Dataset Fake News Detection Confusion Matrix



LIAR Dataset Fake News Detection Confusion Matrix



There are various observations that can be gleaned from the evaluation results:

**1. The combined dataset is significantly imbalanced.**

The ISOT dataset dominates with 31,430 training samples, whereas LIAR dataset contributes only 10240 samples. This imbalance is reflected in the model's performance.

**2. There is a severe performance discrepancy between datasets.**

The model exhibits nearly-perfect performance on the ISOT dataset, with 99.9% across all metrics, and only 2 misclassifications across 8979 test samples.

Meanwhile, it performs worse on the LIAR dataset, with 63.54% accuracy. The confusion matrix shows that 352 fake news items from LIAR were misclassified as real, and 110 real statements were misclassified as fake.

From the above observations, I can draw the following insights:

**1. There is a potential data leakage in ISOT dataset.**

The model's performance on ISOT is suspiciously excellent, suggesting possible issues (e.g. obvious patterns or features in ISOT that make classification trivial, duplicates between train/test splits, unintended data leakage, etc.)

**2. The LIAR dataset is inherently more challenging.**

The LIAR dataset contains political statements. Compared to the ISOT dataset, this means shorter text segments with less context, as well as more nuanced language that requires nuanced understanding of subtleties. Additionally, political claims may require external knowledge to verify, whereas ISOT may contain more obvious linguistic patterns that distinguish fake news.

**3. There is model bias from the training set imbalance.**

As the model is primarily trained on ISOT data, it has optimised for patterns present in the ISOT dataset and given less weight to the characteristics of fake news in LIAR.

As such, the overall metrics are heavily misleading as they are skewed by the dominant (and likely easier) ISOT dataset.

## [7.125] Second Successful Execution: Implementation Changes

To address potential model bias from the training set imbalance, I have implemented two complimentary approaches to balance the influence of each dataset:

### 1. Weighted sampler for DataLoader

The `create_balanced_sampler` function ensures that samples from the smaller LIAR dataset would be selected with higher probability during training, thus balancing the influence of both datasets.

```
[5]: # weighted sampler to balance dataset influence during training
# arguments:
# train_df - training dataframe with 'source' column
# return - PyTorch sampler that balances dataset representation
def create_balanced_sampler(train_df):
    # step 1: calculate weights for each sample based on source
    dataset_counts = train_df['source'].value_counts()
    total_samples = len(train_df)

    # calculate weights (inverse of dataset frequency)
    weights = []
    for source in train_df['source']:
        source_count = dataset_counts[source]
        # weight = total_samples / (num_datasets * source_count)
        weight = total_samples / (2 * source_count) # 2 datasets: LIAR and ISOT
        weights.append(weight)

    # step 2: create weighted sampler
    weights = torch.FloatTensor(weights)
    sampler = torch.utils.data.WeightedRandomSampler(
        weights=weights,
        num_samples=len(weights),
        replacement=True
    )

    return sampler
```

## 2. Class weights for Loss Function

Meanwhile, the `calculate_class_weights` function searches through different combinations of learning rates, weight decay values, batch sizes and number of epochs. The function records and displays all results sorted by F1 score.

```
# calculate class weights for loss function to handle class imbalance
# arguments:
# train_df - training dataframe with 'label' column
# return - class weights tensor for balanced loss calculation
def calculate_class_weights(train_df):
    # step 1: count classes
    class_counts = train_df['label'].value_counts().to_dict()
    total_samples = len(train_df)

    # step 2: calculate weights (inverse of class frequency)
    num_classes = len(class_counts)
    weights = []

    for class_idx in range(num_classes):
        if class_idx in class_counts:
            weight = total_samples / (num_classes * class_counts[class_idx])
        else:
            weight = 1.0
        weights.append(weight)

    # step 3: convert to tensor
    class_weights = torch.FloatTensor(weights)
    return class_weights
```

I implemented weighted sampling and class weights with the aim to help ensure the LIAR dataset (with only 10,240 samples) would not have excessive influence over the larger ISOT dataset in training. In doing so, I hoped that both datasets would contribute meaningfully to the model's learning process while being less biased towards patterns in the ISOT dataset.

Additionally, I have also implemented hyperparameter tuning with grid search.

### 3. Hyperparameter tuning with grid search

```
# hyperparameter tuning with grid search
# arguments:
# model_name - pretrained model name
# train_dataset - dataset for training
# valid_dataset - dataset for validation
# train_sampler - sampler for balanced training
# return - best hyperparameters found
def tune_hyperparameters(model_name, train_dataset, valid_dataset, train_sampler):
    # step 1: define hyperparameter grid
    param_grid = {
        "learning_rate": [1e-5, 2e-5, 3e-5, 5e-5],
        "weight_decay": [0.01, 0.05, 0.1, 0.2],
        "batch_size": [8, 16, 32],
        "epochs": [3, 4, 5]
    }

    # step 2: initialize tracking variables
    best_f1 = 0
    best_params = None
    results = []

    # step 3: manual grid search (using Trainer for each combination)
    print("Starting hyperparameter tuning...")
    total_combos = (len(param_grid["learning_rate"]) *
                    len(param_grid["weight_decay"]) *
                    len(param_grid["batch_size"]) *
                    len(param_grid["epochs"]))

    # To reduce search time, we'll select a smaller subset of combinations
    # This approach tries each parameter while keeping others at default values
    # It's a simplified grid search that's more efficient than testing all combinations
    default_lr = 2e-5
    default_wd = 0.1
    default_bs = 16
    default_epochs = 4

    # build combinations list
    combinations = []

    # vary learning rate
    for lr in param_grid["learning_rate"]:
        combinations.append({
            "learning_rate": lr,
            "weight_decay": default_wd,
            "batch_size": default_bs,
            "epochs": default_epochs
        })

    # vary weight decay
    for wd in param_grid["weight_decay"]:
        if wd != default_wd: # Skip duplicate of default combo
            combinations.append({
                "learning_rate": default_lr,
                "weight_decay": wd,
                "batch_size": default_bs,
                "epochs": default_epochs
            })

    # vary batch size
    for bs in param_grid["batch_size"]:
        if bs != default_bs: # Skip duplicate of default combo
            combinations.append({
                "learning_rate": default_lr,
                "weight_decay": default_wd,
                "batch_size": bs,
                "epochs": default_epochs
            })

    # vary epochs
    for ep in param_grid["epochs"]:
        if ep != default_epochs: # Skip duplicate of default combo
            combinations.append({
                "learning_rate": default_lr,
                "weight_decay": default_wd,
                "batch_size": default_bs,
                "epochs": ep
            })

    print(f"Testing {len(combinations)} hyperparameter combinations instead of {total_combos}")
```

The `tune_hyperparameters` function will search through different combinations of learning rate, weight decay values, batch sizes and number of epochs, then records and displays all results sorted by F1 score. This made it easier for me to understand the effects of different hyperparameters.

#### 4. Early stopping

```
# set up early stopping
early_stopping_callback = EarlyStoppingCallback(
    early_stopping_patience=2,
    early_stopping_threshold=0.001
)

# initialize Trainer for final training
trainer = Trainer(
    model=model,
    args=training_args,
    train_dataset=train_dataset,
    eval_dataset=valid_dataset,
    compute_metrics=compute_metrics,
    callbacks=[early_stopping_callback]
)
```

Early stopping was implemented to automatically stop training the model if improvements stop, with the goal of preventing overfitting and saving time.

If the evaluation metric (F1 score) did not improve for 2 evaluation steps, training would stop. A threshold of 0.001 was used in this version.

This is then added to the Trainer to tell the trainer to monitor evaluation performance and use the callback to stop training early if required.



## [7.126] Second Successful Execution: Process

The hyperparameter tuning process attempted eleven (11) combinations of hyperparameters:

The initial set of hyperparameters yielded these results:

```
Trying combination 1/11:  
Learning rate: 1e-05  
Weight decay: 0.1  
Batch size: 16  
Epochs: 4
```

[10420/10420 20:36, Epoch 4/4]

Epoch	Training Loss	Validation Loss	Accuracy	Precision	Recall	F1	Roc Auc
1	2466.217300	276.600342	0.857093	0.998485	0.706324	0.827370	0.852657
2	8.536100	3.574726	0.526763	0.506446	0.940336	0.658329	0.538931
3	0.412600	0.391899	0.906288	0.855703	0.970347	0.909426	0.908173
4	0.334300	0.267923	0.908540	0.877368	0.943194	0.909091	0.909559

[caption all imgs]

Hyperparameter combination 2 was the next improvement over the initial set of hyperparameters:

```
Trying combination 2/11:  
Learning rate: 2e-05  
Weight decay: 0.1  
Batch size: 16  
Epochs: 4
```

[10420/10420 20:11, Epoch 4/4]

Epoch	Training Loss	Validation Loss	Accuracy	Precision	Recall	F1	Roc Auc
1	0.312000	0.169905	0.906808	0.864560	0.957842	0.908814	0.908309
2	0.259300	0.184365	0.913563	0.894376	0.931761	0.912686	0.914099
3	0.189900	0.246792	0.911658	0.887572	0.936406	0.911335	0.912386
4	0.267600	0.272393	0.912351	0.896576	0.926045	0.911072	0.912754

[361/361 00:14]

```
F1 Score: 0.9127  
New best F1 score: 0.9127
```

Ultimately, hyperparameter combination 6 yielded the best F1 score:

```
Trying combination 6/11:  
Learning rate: 2e-05  
Weight decay: 0.05  
Batch size: 16  
Epochs: 4
```

[10420/10420 49:36, Epoch 4/4]

Epoch	Training Loss	Validation Loss	Accuracy	Precision	Recall	F1	Roc Auc
1	0.307700	0.170086	0.909233	0.866109	0.961415	0.911277	0.910768
2	0.255000	0.188364	0.917374	0.893559	0.941765	0.917029	0.918092
3	0.197700	0.239377	0.913217	0.896480	0.928189	0.912059	0.913657
4	0.155300	0.263501	0.914083	0.900522	0.924973	0.912584	0.914403

[361/361 00:32]

```
F1 Score: 0.9170  
New best F1 score: 0.9170
```

```
Hyperparameter tuning complete!
Best F1 score: 0.9170
Best parameters: {'learning_rate': 2e-05, 'weight_decay': 0.05, 'batch_size': 16, 'epochs': 4}

All Results (sorted by F1 score):
F1: 0.9170, Params: {'learning_rate': 2e-05, 'weight_decay': 0.05, 'batch_size': 16, 'epochs': 4}
F1: 0.9152, Params: {'learning_rate': 2e-05, 'weight_decay': 0.1, 'batch_size': 16, 'epochs': 3}
F1: 0.9145, Params: {'learning_rate': 2e-05, 'weight_decay': 0.1, 'batch_size': 32, 'epochs': 4}
F1: 0.9143, Params: {'learning_rate': 2e-05, 'weight_decay': 0.1, 'batch_size': 16, 'epochs': 5}
F1: 0.9136, Params: {'learning_rate': 2e-05, 'weight_decay': 0.01, 'batch_size': 16, 'epochs': 4}
F1: 0.9135, Params: {'learning_rate': 2e-05, 'weight_decay': 0.2, 'batch_size': 16, 'epochs': 4}
F1: 0.9128, Params: {'learning_rate': 2e-05, 'weight_decay': 0.1, 'batch_size': 8, 'epochs': 4}
F1: 0.9127, Params: {'learning_rate': 2e-05, 'weight_decay': 0.1, 'batch_size': 16, 'epochs': 4}
F1: 0.9106, Params: {'learning_rate': 3e-05, 'weight_decay': 0.1, 'batch_size': 16, 'epochs': 4}
F1: 0.9094, Params: {'learning_rate': 1e-05, 'weight_decay': 0.1, 'batch_size': 16, 'epochs': 4}
F1: 0.9094, Params: {'learning_rate': 5e-05, 'weight_decay': 0.1, 'batch_size': 16, 'epochs': 4}
```

The model was then trained with the identified optimal hyperparameters.

Epoch	Training Loss	Validation Loss	Accuracy	Precision	Recall	F1	Roc Auc
1	0.306400	0.177676	0.907847	0.856111	0.973562	0.911067	0.909780
2	0.245900	0.185932	0.912697	0.891772	0.933190	0.912011	0.913300
3	0.180800	0.244151	0.910272	0.898080	0.919257	0.908545	0.910536

[7.127] Second Successful Execution: Evaluation

1. Performance on Test Set

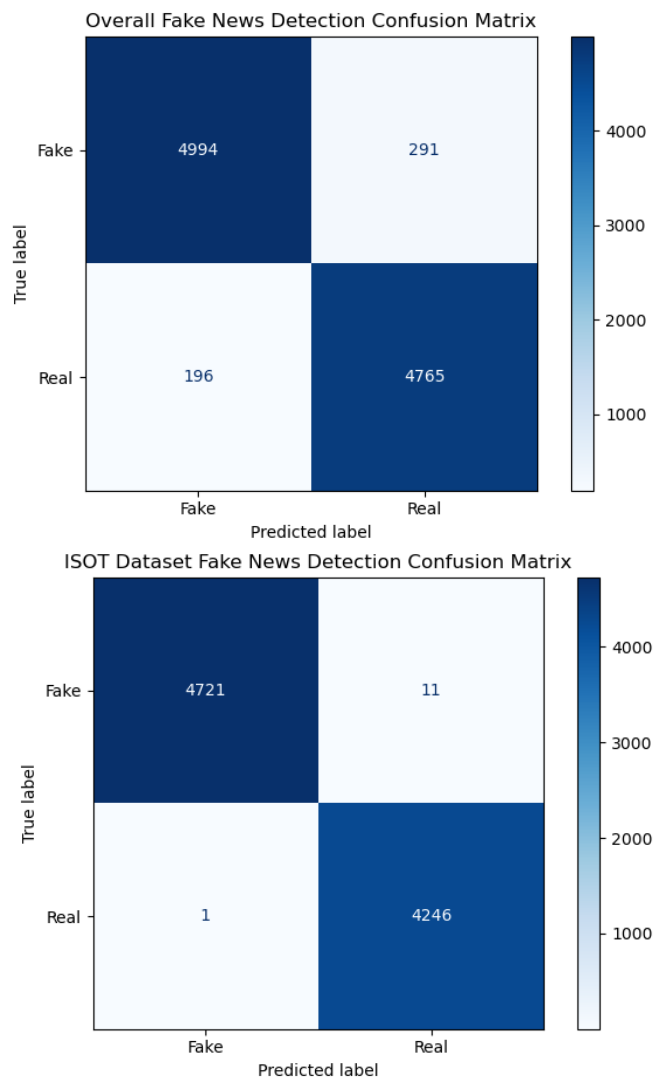
```
Overall Performance Metrics:
Metric  loss  accuracy  precision  recall    f1  roc_auc  runtime \
Value   0.1009   0.9525    0.9424   0.9605  0.9514  0.9527  59.874

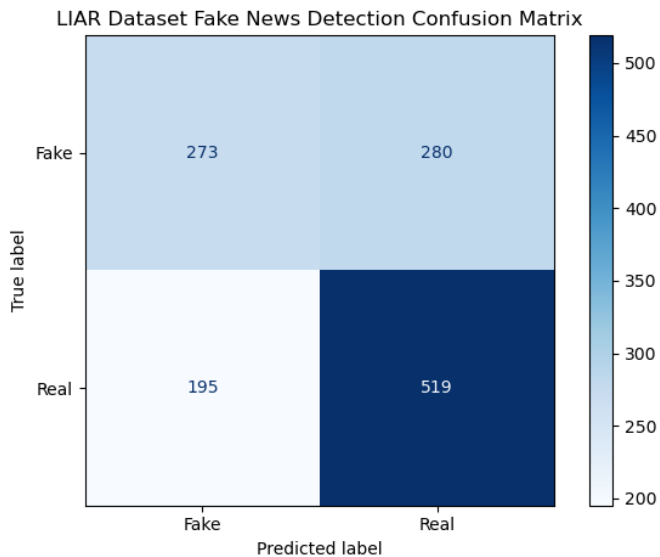
Metric  samples_per_second  steps_per_second  epoch
Value                171.126         10.706    3.0

ISOT Dataset Performance Metrics:
Metric  accuracy  precision  recall    f1  roc_auc
Value   0.9987    0.9974   0.9998  0.9986  0.9987

LIAR Dataset Performance Metrics:
Metric  accuracy  precision  recall    f1  roc_auc
Value   0.6251    0.6496   0.7269  0.6861  0.6103
```

2. Confusion matrices





### 3. Optimal hyperparameters (from hyperparameter tuning)

```
Hyperparameter tuning complete!
Best F1 score: 0.9170
Best parameters: {'learning_rate': 2e-05, 'weight_decay': 0.05, 'batch_size': 16, 'epochs': 4}

All Results (sorted by F1 score):
F1: 0.9170, Params: {'learning_rate': 2e-05, 'weight_decay': 0.05, 'batch_size': 16, 'epochs': 4}
F1: 0.9152, Params: {'learning_rate': 2e-05, 'weight_decay': 0.1, 'batch_size': 16, 'epochs': 3}
F1: 0.9145, Params: {'learning_rate': 2e-05, 'weight_decay': 0.1, 'batch_size': 32, 'epochs': 4}
F1: 0.9143, Params: {'learning_rate': 2e-05, 'weight_decay': 0.1, 'batch_size': 16, 'epochs': 5}
F1: 0.9136, Params: {'learning_rate': 2e-05, 'weight_decay': 0.01, 'batch_size': 16, 'epochs': 4}
F1: 0.9135, Params: {'learning_rate': 2e-05, 'weight_decay': 0.2, 'batch_size': 16, 'epochs': 4}
F1: 0.9128, Params: {'learning_rate': 2e-05, 'weight_decay': 0.1, 'batch_size': 8, 'epochs': 4}
F1: 0.9127, Params: {'learning_rate': 2e-05, 'weight_decay': 0.1, 'batch_size': 16, 'epochs': 4}
F1: 0.9106, Params: {'learning_rate': 3e-05, 'weight_decay': 0.1, 'batch_size': 16, 'epochs': 4}
F1: 0.9094, Params: {'learning_rate': 1e-05, 'weight_decay': 0.1, 'batch_size': 16, 'epochs': 4}
F1: 0.9094, Params: {'learning_rate': 5e-05, 'weight_decay': 0.1, 'batch_size': 16, 'epochs': 4}
```

The optimal configuration was determined to be combination 6 (2e-05 learning rate, 0.05 weight decay, 16 batch size, 4 epochs). Though interestingly, the top 8 configurations used a learning rate of 2e-05 and achieved F1 scores between 0.912 and 0.917, suggesting this learning rate is optimal.

To assess the evaluation results of this execution run and identify how the changes from the initial execution has changed the effectiveness of the model, I will draw a table to compare the evaluation results from both execution runs:

#### Overall Performance

Metric	Execution 1	Execution 2	Change
Accuracy	0.9547	0.9525	-0.0022
Precision	0.9322	0.9424	+0.0102
Recall	0.9776	0.9605	-0.0171
F1 Score	0.9543	0.9514	-0.0029
ROC-AUC	0.9554	0.9527	-0.0027

The overall metrics show a slight trade-off: while Execution 2 has improved precision, this comes at a cost of decreased recall and overall accuracy. This indicates the model is more careful in classifying and makes fewer positive predictions.

## ISOT Dataset

Metric	Execution 1	Execution 2	Change
Accuracy	0.9998	0.9987	-0.0011
Precision	0.9998	0.9974	-0.0024
Recall	0.9998	0.9998	0
F1 Score	0.9998	0.9986	-0.0012
ROC-AUC	0.9998	0.9987	-0.0011

ISOT performance remains exceptionally high across both executions, though Execution 2 exhibits very slightly lower performance.

## LIAR Dataset

Metric	Execution 1	Execution 2	Change
Accuracy	0.6354	0.6251	-0.0103
Precision	0.6318	0.6496	+0.0178
Recall	0.8459	0.7269	-0.01190
F1 Score	0.7234	0.6861	-0.0373
ROC-AUC	0.6047	0.6103	+0.0056

The LIAR dataset shows the most drastic changes. There is a substantial drop in recall, along with a slight increase in precision. Interestingly, despite decreases in other metrics, there is a small improvement in ROC-AUC.

From the above comparison between both executions, there are a few key learning points:

- 1. The implementation of a weighted sampler, class weights & hyperparameter tuning created a more balanced classifier for the LIAR dataset.**

Execution 1 displayed a large gap between recall (0.8459) and precision (0.6318), indicating a classifier biased towards labelling content as “real”. This is evident in the confusion matrices. Execution 2 has more balanced precision and recall, suggesting that the dataset balancing techniques addressed the bias, although not to the extent that I hoped.

- 2. The model from Execution 2 might have better at discriminating fake news despite lower accuracy.**

The increased ROC-AUC on Execution 2 was surprising, considering that it exhibited lower accuracy and lower F1 scores. The original model may have been overly optimistic about classifying political statements as real.

- 3. Dataset balancing seems to primarily affect how models handle difficult classification tasks.**

The changes to ISOT performance across both executions were minimal, indicating that dataset balancing primarily affects how models handle difficult classifications while having less impact on more straightforward classification tasks.

## [7.128] Conclusions and Future Work from this point

While the overall metrics for Execution 2 has slightly decreased compared to Execution 1 of Version 2, the demonstrated improved balance between precision and recall on challenging content (LIAR dataset) indicates to me that Version 2 would be more trustworthy in real-world applications, as Execution 1 was overly-optimistic in classifying statements as ‘real’.

Interestingly, the implementation of the weighted sampler, class weights & hyperparameter tuning did not dramatically improve the model’s performance on the LIAR dataset. This suggests to me that balancing datasets alone will not be sufficient to address the model’s present inability to classify political statements from the LIAR dataset.

Currently planned future work from this point will need to achieve the following goals:

1. Expanding the data pipeline to include more data.
2. Improving performance on the LIAR dataset, if possible/feasible.
3. Creating a simple CLI front-end for the end-user to interact with the model.

On top of that, I would personally prefer to make changes to the hyperparameter tuning implementation – specifically, I would like to replace my custom grid search with an established library (e.g. Optuna).

## [7.13] Version 2.1 & 2.2

### [7.131] Overview & Aims

The goal of Version 2.1 is to enhance the hyperparameter tuning logic implemented in the pipeline. Whereas the previous version used a simplified grid search that manually iterated through a limited set of combinations of hyperparameters, this new version integrates Bayesian optimisation using the Optuna library.

The goal of Version 2.2 is to further enhance with the pipeline by implementing a lightweight text pre-processing function appropriate for feeding data into transformer models, and by implementing gradient accumulation.

In Version 2.2, the implementation of mixed precision training was also explored, but was ultimately not kept in the codebase due to its drawbacks.

User evaluation is not performed at this stage as there are no changes to user-facing elements, hence I do not expect meaningful insights to emerge.

### [7.132] Version 2.1: Implementation

I have changed the hyperparameter tuning logic to implement Bayesian optimisation through Optuna.

While the implementation in Version 2 tested each parameter individually, Version 2.1 implements Optuna-based Bayesian optimisation that encompasses several elements:

#### 1. Hyperparameter search method

While Version 2 used a fixed grid of hyperparameter combinations.

Version 2.1 uses Optuna's TPESampler to explore hyperparameter combinations over a defined range.

```
def objective(trial):  
    # step 2.1: sample hyperparameter values  
    lr = trial.suggest_float("learning_rate", 1e-5, 5e-5, log=True)  
    weight_decay = trial.suggest_float("weight_decay", 0.01, 0.2, log=True)  
    batch_size = trial.suggest_categorical("batch_size", [8, 16, 32])  
    epochs = trial.suggest_int("epochs", 2, 3)
```

The number of trials were set to 2. While higher values can be used, limited time and computing resources required me to accelerate experimentation. This would not cause issues because the model already scored well overall in evaluation metrics.

```
# step 6: perform hyperparameter tuning with Optuna  
best_hparams = tune_hyperparameters(  
    model_name,  
    train_dataset,  
    valid_dataset,  
    balanced_sampler,  
    n_trials=2 # adjust this (trials)  
)
```





## 2. Training robustness improvements

Checkpoints after each trial are saved after training to allow for recovery in case of interruption.

```
# save checkpoint for possible resumption
torch.save({
    'epoch': epochs,
    'model_state_dict': model.state_dict(),
    'f1_score': f1_score
}, checkpoint_path)
```

```
# step 2.7: attempt to resume from checkpoint if available
if os.path.exists(checkpoint_path):
    print(f"Resuming from checkpoint for trial {trial.number}")
    checkpoint_state = torch.load(checkpoint_path)
    model.load_state_dict(checkpoint_state['model_state_dict'])
    start_epoch = checkpoint_state['epoch']
```

Training failures are caught and handled gracefully through logging and saving trial state.

```
except Exception as e:
    # handle training failures by saving state
    if start_epoch > 0:
        torch.save({
            'epoch': start_epoch,
            'model_state_dict': model.state_dict(),
            'f1_score': 0 # default for failed runs
        }, checkpoint_path)
    print(f"Training failed: {e}")
    return 0
```

## 3. Storage and reproducibility improvements

The trial results are saved in a persistent SQLite database, improving the reproducibility of trials.

```
# step 3: create and configure the Optuna study
study = optuna.create_study(
    direction="maximize", # maximize F1 score
    sampler=TPESampler(seed=42), # use Tree-structured Parzen Estimator
    pruner=MedianPruner(n_startup_trials=5, n_warmup_steps=5), # prune unpromising trials
    study_name="fake_news_detection",
    storage=f"sqlite:///base_dir/optuna_study.db", # persistent storage
    load_if_exists=True # resume existing study if available
)
```

## [7.133] Version 2.1: Evaluation

The new hyperparameter tuning logic has introduced several benefits:

### 1. Top trials & hyperparameter importance

After hyperparameter tuning, the top trials (ranked by F1 score) are printed, allowing quick assessment of parameter set performance. Optuna's built-in feature importance module also ranks hyperparameters by their impact on F1 score.

```
Optuna hyperparameter tuning complete!
Best F1 score: 0.9141
Best parameters: {'learning_rate': 1.097990803659665e-05, 'weight_decay': 0.13394334706750485, 'batch_size': 16, 'epochs': 3}

Hyperparameter importance:
  batch_size: 0.4054
  learning_rate: 0.2973
  weight_decay: 0.2973

Top 5 trials:
Rank 1: F1=0.9141, Params={'learning_rate': 1.097990803659665e-05, 'weight_decay': 0.13394334706750485, 'batch_size': 16, 'epochs': 3}
Rank 2: F1=0.9141, Params={'learning_rate': 1.827226177606625e-05, 'weight_decay': 0.17254716573280354, 'batch_size': 8, 'epochs': 3}
Rank 3: F1=0.9140, Params={'learning_rate': 1.827226177606625e-05, 'weight_decay': 0.17254716573280354, 'batch_size': 8, 'epochs': 2}
```

### 2. Benefit of checkpointing

If training crashes, the checkpointing system helps to avoid wasted work and allows safe resumption. This turned out to be a boon as I have limited computing resources/power available.

Additionally, the model was trained on the set of hyperparameters that was identified as optimal:

#### 1. Performance on test set

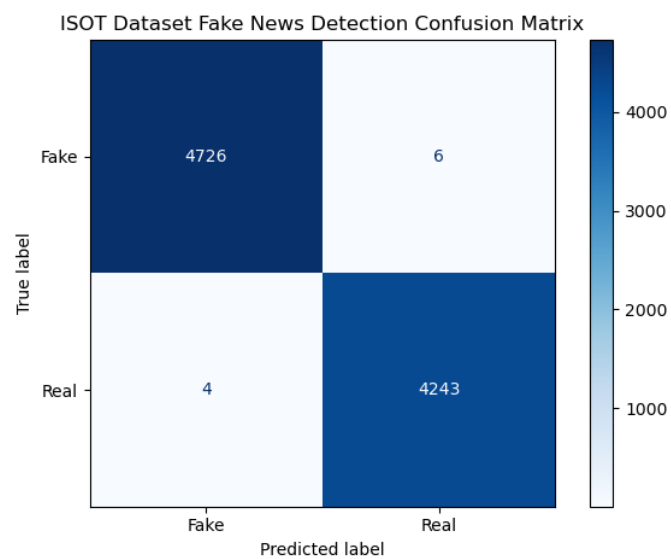
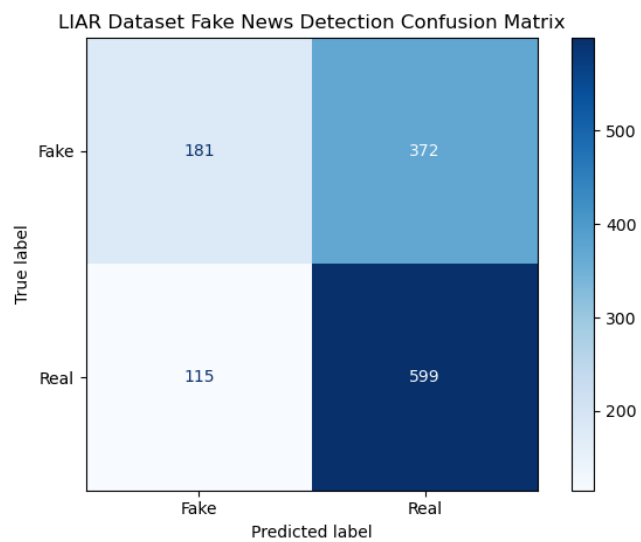
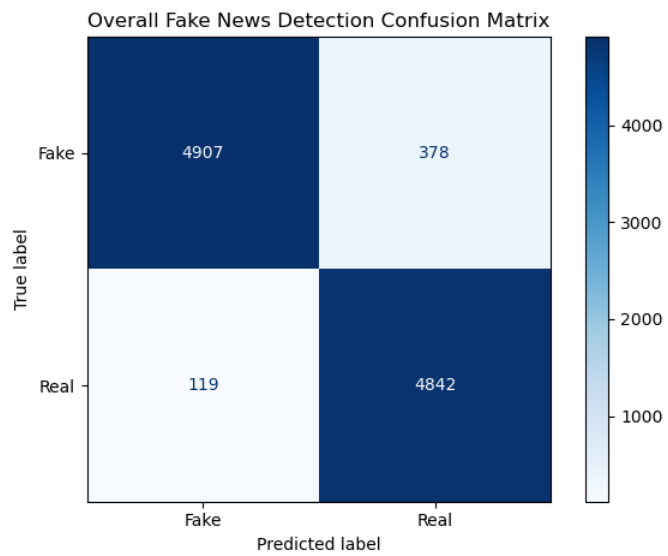
```
Overall Performance Metrics:
Metric  loss  accuracy  precision  recall    f1  roc_auc  runtime \
Value   0.1783   0.9515    0.9276   0.976   0.9512   0.9522   59.5297

Metric  samples_per_second  steps_per_second  epoch
Value           172.116           10.768    3.0

LIAR Dataset Performance Metrics:
Metric  accuracy  precision  recall    f1  roc_auc
Value   0.6156    0.6169   0.8389   0.711  0.5831

ISOT Dataset Performance Metrics:
Metric  accuracy  precision  recall    f1  roc_auc
Value   0.9989    0.9986   0.9991   0.9988  0.9989
```

## 2. Confusion matrices



To assess how the new hyperparameter tuning logic has impacted the model's performance, I can compare the performance of the model trained with the new hyperparameter tuning logic against the performance of the model trained in Version 2 (Execution 2).

## Overall Performance

Metric	Version 2 (exec. 2)	Version 2.1	Change
Accuracy	0.9525	0.9515	-0.0010
Precision	0.9424	0.9276	-0.0148
Recall	0.9605	0.9760	+0.0155
F1 Score	0.9514	0.9512	-0.0002
ROC-AUC	0.9527	0.9522	-0.0005

Interestingly, Version 2.1 shows an increase in recall (+0.0155) at the slight cost of precision (-0.0148). This suggests that the updated approach slightly increases false positives but identifies more instances of fake news.

Despite the lower precision, the F1 score is only slightly decreased (-0.0002), which indicates the overall effectiveness of the model is largely similar. This is an overall positive development, as recall was boosted significantly while maintaining similar overall performance.

## ISOT Dataset

Metric	Version 2 (exec. 2)	Version 2.1	Change
Accuracy	0.9987	0.9989	+0.0002
Precision	0.9974	0.9986	+0.0012
Recall	0.9998	0.9991	-0.0007
F1 Score	0.9986	0.9988	+0.0002
ROC-AUC	0.9987	0.9989	+0.0002

The model's performance on the ISOT dataset remains exceptionally high, indicating that my model architecture is already sufficiently robust to reach the ceiling of possible performance on this dataset.

## LIAR Dataset

Metric	Version 2 (exec. 2)	Version 2.1	Change
Accuracy	0.6251	0.6156	-0.0095
Precision	0.6496	0.6169	-0.0327
Recall	0.7269	0.8389	+0.1120
F1 Score	0.6861	0.7110	+0.0249
ROC-AUC	0.6103	0.5831	-0.0272

The model in Version 2.1 has a substantial increase in recall (+0.112), achieved at the cost of precision (-0.0327), indicating that Version 2.1 is far more aggressive at flagging potential fake news in the challenging dataset, identifying 11% more fake news articles though with increased false positives.

Despite the precision decrease, the improved F1 score (+0.0249) suggests that the recall improvement sufficiently compensates for the precision drop in terms of overall effectiveness on the LIAR dataset.

However, the decrease in ROC-AUC (-0.0272) indicates that the more aggressive model in Version 2.1 comes with some cost to the model's confidence.

The above observations reveal an implication in a fake news detection system – is it better to catch more fake stories even if it means falsely flagging some legitimate news? While Version 2.1's model identifies more fake news (higher recall), it does so at the cost of more false alarms (lower precision).

## [7.134] Version 2.2: Implementation

Some more improvements were added to Version 2.2 of the project to enhance the pipeline. A text pre-processing function was added, and gradient accumulation was implemented to the model training process.

Additionally, mixed precision training was attempted but was ultimately not implemented due to issues that I was unable to resolve.

### 1. Text pre-processing

A simple, very minimalistic text pre-processing function was added to handle basic text issues (e.g. encoding, whitespace) while preserving semantic content that transformer-based models could handle.

Unlike traditional NLP pre-processing, this approach is very deliberately minimalistic to retain stylistic elements and contextual information that transformer models can leverage for better classification performance.

```
# function to perform minimal text preprocessing for transformer models
# arguments:
# text - input text to preprocess (string)
# max_length - maximum number of words to keep (default: 512)
#
# returns preprocessed text string
def minimal_text_preprocessing(text, max_length=512):
    # step 1: handle edge cases
    if not isinstance(text, str):
        return ""

    # step 2: clean encoding issues
    # remove non-ascii characters that could cause problems
    text = text.encode('ascii', 'ignore').decode('ascii')

    # step 3: normalize text format
    # replace multiple whitespace characters with single space
    text = re.sub(r'\s+', ' ', text).strip()

    # step 4: truncate oversized inputs
    # transformer models have context length limits
    words = text.split()
    if len(words) > max_length:
        text = ' '.join(words[:max_length])

    return text
```

## 2. Gradient accumulation

Rather than updating the model parameters after each batch, gradients are now added up across several batches, before a single update is performed. This effectively simulates training with larger batch sizes while circumventing the memory requirements.

My implementation leverages the built-in support for gradient accumulation in Hugging Face's `TrainingArguments` class.

It is added as a hyperparameter in the search space when tuning hyperparameters.

```
def objective(trial):
    # step 2.1: sample hyperparameter values
    lr = trial.suggest_float("learning_rate", 1e-5, 5e-5, log=True)
    weight_decay = trial.suggest_float("weight_decay", 0.01, 0.2, log=True)
    batch_size = trial.suggest_categorical("batch_size", [8, 16, 32])
    epochs = trial.suggest_int("epochs", 2, 3)
    gradient_accumulation_steps = trial.suggest_categorical("gradient_accumulation_steps", [1, 2, 4]) # new (v2.5 step2)
```

It is then applied to training configuration in hyperparameter tuning and in final model training:

```
# configure training with optimal hyperparameters
training_args = TrainingArguments(
    output_dir="./balanced_model",
    eval_strategy="epoch",
    save_strategy="epoch",
    learning_rate=best_hparams["learning_rate"],
    per_device_train_batch_size=best_hparams["batch_size"],
    per_device_eval_batch_size=best_hparams["batch_size"],
    num_train_epochs=best_hparams["epochs"],
    weight_decay=best_hparams["weight_decay"],
    load_best_model_at_end=True,
    metric_for_best_model="f1",
    push_to_hub=False,
    report_to="none",
    gradient_accumulation_steps=best_hparams.get("gradient_accumulation_steps", 1) #new (v2.5 step2): gradient accumulation
)
```

As earlier trials (versions) did not implement gradient accumulation, I use `get()` with a default value of 1 in the final training to handle cases where that parameter did not exist in earlier trials in order to avoid `KeyError` exceptions.

### 3. Mixed Precision Training (ultimately not implemented)

There were attempts to implement mixed precision training. Ultimately, this was not successfully achieved.

In the TrainingArguments setup, passing 'fp16=True' enabled mixed precision training, however this resulted in too many NaN values and the model failing to perform.

Many rounds of changes to fp16\_opt\_level configuration, gradient scaling parameters and loss scaling adjustments did not resolve this issue.

I have several hypotheses to why mixed precision training did not work:

- **MobileBERT is already an optimised model**

The model is already designed for efficiency, and additional attempts to optimise may have interfered.

Additionally, some operations in this model architecture may simply have not been compatible with FP16.

- **Batch sizes were too small for mixed precision training**

Mixed precision training usually offers more pronounced benefits when larger batch sizes are implemented, and small batches thus potentially faced stability issues.



# [7.135] Version 2.2: Evaluation

After implementing all the changes of Version 2.2, running the hyperparameter tuning pipeline found a new set of optimal hyperparameters:

```
Optuna hyperparameter tuning complete!
Best F1 score: 0.9168
Best parameters: {'learning_rate': 1.827226177606625e-05, 'weight_decay': 0.17254716573280354, 'batch_size': 8, 'epochs': 2, 'gradient_accumulation_steps': 2}

Hyperparameter importance:
  learning_rate: 0.5058
  weight_decay: 0.3818
  batch_size: 0.1125

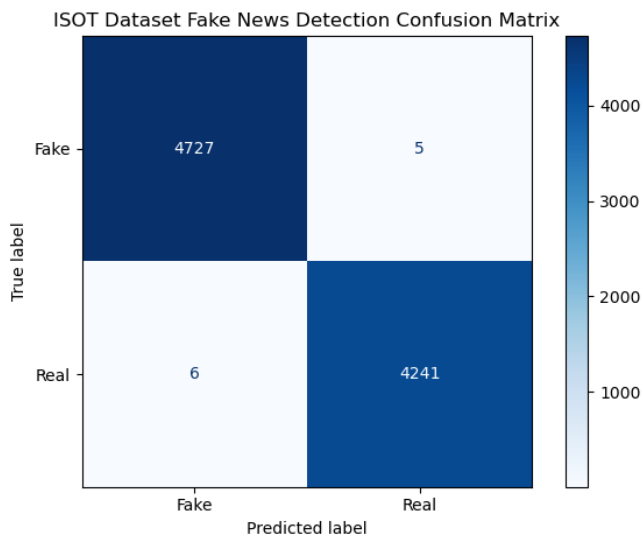
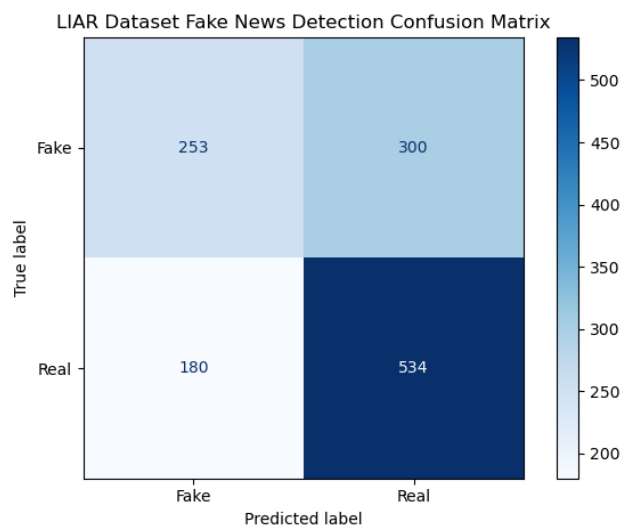
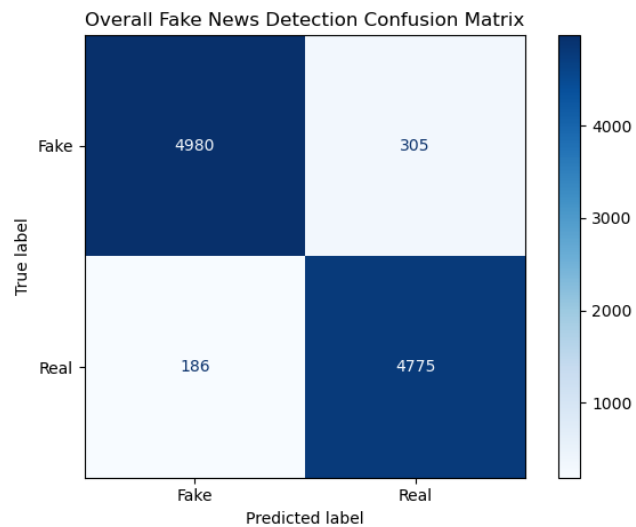
Top 5 trials:
Rank 1: F1=0.9168, Params={'learning_rate': 1.827226177606625e-05, 'weight_decay': 0.17254716573280354, 'batch_size': 8, 'epochs': 2, 'gradient_accumulation_steps': 2}
Rank 2: F1=0.9141, Params={'learning_rate': 1.097990803659665e-05, 'weight_decay': 0.13394334706750485, 'batch_size': 16, 'epochs': 3}
Rank 3: F1=0.9141, Params={'learning_rate': 1.827226177606625e-05, 'weight_decay': 0.17254716573280354, 'batch_size': 8, 'epochs': 3}
Rank 4: F1=0.9140, Params={'learning_rate': 1.827226177606625e-05, 'weight_decay': 0.17254716573280354, 'batch_size': 8, 'epochs': 2}
Rank 5: F1=0.9131, Params={'learning_rate': 1.827226177606625e-05, 'weight_decay': 0.17254716573280354, 'batch_size': 8, 'epochs': 2, 'gradient_accumulation_steps': 2}
```

The newly-discovered set of optimal parameters were used to train the model:

## 1. Performance on test set

Overall Performance Metrics:							
Metric	loss	accuracy	precision	recall	f1	roc_auc	runtime \
Value	0.0926	0.9521	0.94	0.9625	0.9511	0.9524	100.7962
Metric	samples_per_second		steps_per_second		epoch		
Value	101.651		12.709		1.999424		
LIAR Dataset Performance Metrics:							
Metric	accuracy	precision	recall	f1	roc_auc		
Value	0.6212	0.6403	0.7479	0.6899	0.6027		
ISOT Dataset Performance Metrics:							
Metric	accuracy	precision	recall	f1	roc_auc		
Value	0.9988	0.9988	0.9986	0.9987	0.9988		

## 2. Confusion matrices



From the confusion matrices and evaluation results above:

### Overall Performance

Metric	V2.2	V2.1	Change
Accuracy	0.9521	0.9515	+0.0006
Precision	0.9400	0.9276	+0.0124
Recall	0.9625	0.9760	-0.0135
F1 Score	0.9511	0.9512	-0.0001
ROC-AUC	0.9524	0.9522	+0.0002

### ISOT Dataset

Metric	V2.2	V2.1	Change
Accuracy	0.9988	0.9989	-0.0001
Precision	0.9988	0.9986	+0.0002
Recall	0.9986	0.9991	-0.0005
F1 Score	0.9987	0.9988	-0.0001
ROC-AUC	0.9988	0.9989	-0.0001

The model's performance on the ISOT dataset remains exceptionally high, indicating that my model architecture is already sufficiently robust to reach the ceiling of possible performance on this dataset.

### LIAR Dataset

Metric	V2.2	V2.1	Change
Accuracy	0.6212	0.6156	+0.0056
Precision	0.6403	0.6169	+0.0234
Recall	0.7479	0.8389	-0.0910
F1 Score	0.6899	0.7110	-0.0211
ROC-AUC	0.6027	0.5831	+0.0196

On the LIAR dataset, precision has increased (+2.34%), indicating that the model is more conservative in labelling fake news and resulting in fewer false positives. ROC-AUC score has also improved (+1.96%), and accuracy has improved slightly (+0.56%).

However, the recall dropped significantly (-9.10%), meaning the model catches fewer fake news samples. The F1 score on LIAR dropped (-2.11%) to reflect this trade-off, and overall recall also dropped slightly (-1.35%).

The above results indicate that:

1. Gradient accumulation may have possibly over-regularised the model slightly.
2. Noise removal may have helped precision by removing encoding errors and whitespace, cleaning up bad inputs, making the model more confident in predicting true negatives and reducing false positives.
3. Removing non-ASCII characters and truncating long texts may have accidentally stripped contextual clues (e.g. emphasis, sarcasm, slang)
4. Truncation can clip off the end of articles, where fake news cues (e.g. sources) may appear, leading to more false negatives and reducing recall.

These trade-offs highlight how minor decisions relating to model architecture and data pre-processing can significantly influence the balance between model conservativeness and sensitivity.

## [7.14] Version 3

### [7.141] Overview & Aims

In Version 3, I re-implement the Application Layer so that an end-user may interact with the model.

As changes are made to the user-facing elements of the application, I will also go through the user evaluation process for this version.

## [7.2] References

- [1] “Advanced Python Project - Detecting Fake News with Python,” DataFlair, Sep. 16, 2019. <https://data-flair.training/blogs/advanced-python-project-detecting-fake-news/>
- [2] “ClaimBuster,” Uta.edu, 2024. <https://idir.uta.edu/claimbuster/event/70/> (accessed Dec. 15, 2024).
- [3] Snopes, “Snopes.com,” Snopes.com, 2018. <https://www.snopes.com/>
- [4] L.-J. Chang, “Comparison of Machine Learning and Deep Learning Algorithms in Detecting Fake News,” Proceedings of the 28th World Multi-Conference on Systemics, Cybernetics and Informatics, pp. 203–209, Sep. 2024, doi: <https://doi.org/10.54808/wmsci2024.01.203>.
- [5] K. Huh, “Surviving In a Random Forest with Imbalanced Datasets,” SFU Professional Computer Science, Feb. 13, 2021. <https://medium.com/sfu-csmp/surviving-in-a-random-forest-with-imbalanced-datasets-b98b963d52eb>
- [6] J. Devlin, M.-W. Chang, K. Lee, K. Google, and A. Language, “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding,” 2019. Available: <https://aclanthology.org/N19-1423.pdf>
- [7] Z. Sun, H. Yu, X. Song, R. Liu, Y. Yang, and D. Zhou, “MobileBERT: a Compact Task-Agnostic BERT for Resource-Limited Devices,” arXiv:2004.02984 [cs], Apr. 2020, Available: <https://arxiv.org/abs/2004.02984>
- [8] V. Korolev and P. Protsenko, “Accurate, interpretable predictions of materials properties within transformer language models,” Oct. 13, 2023. Accessed: Dec. 15, 2024. [Online]. Available: [https://www.cell.com/patterns/pdfExtended/S2666-3899\(23\)00158-7](https://www.cell.com/patterns/pdfExtended/S2666-3899(23)00158-7)
- [9] “ClaimBuster,” RAND. <https://www.rand.org/research/projects/truth-decay/fighting-disinformation/search/items/claimbuster.html>
- [10] N. Hassan et al., “ClaimBuster,” Proceedings of the VLDB Endowment, vol. 10, no. 12, pp. 1945–1948, Aug. 2017, doi: <https://doi.org/10.14778/3137765.3137815>.
- [11] Y. Zhu, Y. Li, J. Wang, M. Gao, and J. Wei, “FaKnow: A Unified Library for Fake News Detection,” arXiv (Cornell University), 2024. <https://arxiv.org/abs/2401.16441> (accessed Dec. 15, 2024).
- [12] L. Graves and F. Cherubini, “The Rise of Fact-Checking Sites in Europe,” Reuters Institute, 2016. Available: <https://reutersinstitute.politics.ox.ac.uk/sites/default/files/research/files/The%2520Rise%2520of%2520Fact-Checking%2520Sites%2520in%2520Europe.pdf>
- [13] X.-J. Liu, Q. Li, L. Wang, and M. J. Metzger, “Checking the Fact-Checkers: The Role of Source Type, Perceived Credibility, and Individual Differences in Fact-Checking Effectiveness,” Communication Research, Oct. 2023, doi: <https://doi.org/10.1177/00936502231206419>.

- [14] J. Yang, D. A. Vega-Oliveros, T. Seibt, and A. Freitas, “Scalable Fact-checking with Human-in-the-Loop,” arXiv (Cornell University), Dec. 2021, doi: <https://doi.org/10.1109/wifs53200.2021.9648388>.
- [15] Zi Hen Lin, Z. Wang, M. Zhao, Y. Song, and L. Lan, “An AI-based System to Assist Human Fact-Checkers for Labeling Cantonese Fake News on Social Media,” 2022 IEEE International Conference on Big Data (Big Data), Dec. 2022, doi: <https://doi.org/10.1109/bigdata55660.2022.10020949>.
- [16] “google/mobilebert-uncased · Hugging Face,” Hugging Face. <https://huggingface.co/google/mobilebert-uncased>
- [17] tfs4, “GitHub - tfs4/liar\_dataset: dataset liar,” GitHub, 2019. [https://github.com/tfs4/liar\\_dataset](https://github.com/tfs4/liar_dataset) (accessed Dec. 15, 2024).