

Transformer vs. Classical ML for Fake News Detection: A SHAP-Explainable Comparative Study on WELFake

Moqbel Tawhib Salah Abdo

Student ID: 202353460008

Department of Artificial Intelligence

School of Artificial Intelligence and Engineering

Nanjing University of Information Science and

Technology

China 2025

Hosen Md Roki

Student ID: 202353460011

Department of Artificial Intelligence

School of Artificial Intelligence and Engineering

Nanjing University of Information Science and

Technology

China 2025

ABSTRACT

Automated fake news detection requires both high predictive accuracy and interpretable outputs for responsible deployment. This study benchmarks seven models — Logistic Regression, Linear SVM, Random Forest, XGBoost, BERT, RoBERTa, and DistilBERT — under identical experimental conditions on the WELFake corpus (72,134 articles). RoBERTa achieves the best performance (accuracy=96.84%, macro F1=0.968, AUC-ROC=0.9943, MCC=0.937), confirmed via five-fold stratified cross-validation (96.52%±0.27%) and McNemar's pairwise significance tests ($\chi^2 > 6.63$, $p < 0.01$). A SHAP explainability framework identifies sensationalist framing, anonymous sourcing, and absence of quantitative specificity as the dominant linguistic markers of fabricated content. An ablation study demonstrates an 8.65 pp gain from combining title and body text. Together, these results establish that predictive accuracy and interpretability are complementary objectives in responsible fake news AI.

Keywords

Fake News Detection BERT RoBERTa DistilBERT XGBoost Explainable AI (XAI) SHAP NLP WELFake Transformer Models Misinformation Statistical Significance McNemar's Test

1. INTRODUCTION

1.1 Background and Motivation

The global information ecosystem has undergone a structural transformation over the past two decades. Social media platforms — Facebook, X (formerly Twitter), YouTube, WhatsApp, and TikTok — have displaced legacy media as the primary channels through which billions of people encounter news and form political opinions. Unlike traditional journalism, these platforms operate without professional gatekeeping: any individual equipped with a smartphone can publish content reaching millions of users within hours, with no editorial verification, no correction mechanism, and no accountability for factual accuracy.

Fake news — operationally defined as news articles, social media posts, or other digital content that presents fabricated, misleading, or deliberately distorted information as genuine factual reporting — has exploited this structural vulnerability with alarming effectiveness. A landmark study by Vosoughi, Roy, and Aral (2018), published in *Science* after analysing 126,000 verified rumour cascades on Twitter, found that false information reached 1,500 people approximately six times faster than accurate reporting and was 70% more likely to be shared [1]. The downstream consequences are concrete and measurable: COVID-19 vaccine hesitancy fuelled by health

misinformation is estimated to have caused tens of thousands of preventable deaths [2]; fabricated political content demonstrably influenced vote intention in the 2016 US presidential election [3]; health misinformation on social media was formally declared an "infodemic" by the World Health Organisation in 2020. Beyond these headline cases, the chronic erosion of institutional trust in journalism, science, and democratic governance represents a slower-moving but potentially more consequential societal harm.

1.2 Scale of the Problem

The scale of the fake news problem makes manual fact-checking structurally inadequate as a primary mitigation strategy. Facebook reports over 100 billion items of content shared per day; established fact-checking organisations collectively evaluate perhaps a few thousand claims per week — an asymmetry of roughly seven orders of magnitude. This gap makes automated AI-based detection not merely desirable but essential [4].

Early computational approaches employed classical machine learning: Naïve Bayes, Logistic Regression, and Support Vector Machines trained on TF-IDF bag-of-words and n-gram features achieved accuracies of 76–92% on benchmark datasets [5]. The introduction of BERT (Bidirectional Encoder Representations from Transformers) by Devlin et al. (2019) [6] marked a paradigm shift: pre-trained bidirectional transformers, fine-tuned on task-specific data, consistently achieve leading benchmark performance across NLP tasks. BERT variants — RoBERTa [7], DistilBERT [8], ALBERT [9], XLNet [10] — have achieved substantial accuracy improvements over classical baselines on fake news detection tasks.

1.3 The Interpretability Challenge

A fundamental concern accompanies these technical achievements: the interpretability problem. Deep learning models are widely described as "black boxes" — delivering accurate predictions without any transparent account of their internal reasoning. For a fake news detection system deployed in content moderation at scale, this opacity creates serious practical, ethical, and legal consequences. When a system erroneously removes a legitimate news article (false positive), the affected publisher has no actionable basis for appeal. When it systematically exhibits demographic bias, no auditor can diagnose the source without interpretable outputs. Emerging AI transparency regulation — including the EU AI

Act (2024) [11], which classifies automated content moderation as a high-risk AI system subject to mandatory explainability requirements — makes interpretability a legal

obligation, not merely a research preference.

This tension between predictive power and interpretability has motivated the rapidly growing field of Explainable Artificial Intelligence (XAI), and in particular SHAP (SHapley Additive exPlanations) by Lundberg and Lee (2017) [12] — a model-agnostic method grounded in cooperative game theory that assigns each input feature a contribution score satisfying a set of axiomatic fairness properties.

1.4 Problem Statement and Research Gaps

Two critical gaps persist in the existing literature on automated fake news detection. First, most comparative studies benchmark either traditional ML or transformer architectures in isolation; systematic comparisons of both paradigms under identical experimental conditions — same dataset, same train/validation/test splits, same preprocessing, same evaluation protocol — are rare, limiting the practical value of reported performance improvements. Second, the overwhelming majority of fake news detection studies treat the task as pure classification with no attempt to explain what linguistic features drive individual predictions, undermining practical utility, accountability, and regulatory compliance. This study directly addresses both gaps through a unified experimental framework combining rigorous multi-model comparison with comprehensive SHAP-based explainability analysis.

1.5 Research Objectives

1. Conduct a systematic comparative evaluation of seven classification models — Logistic Regression, Linear SVM, Random Forest, XGBoost, BERT, RoBERTa, and DistilBERT — for binary fake news detection under strictly identical experimental conditions.
2. Benchmark all models on the WELFake dataset (72,134 articles) using stratified 80/10/10 train/validation/test splits and 5-fold cross-validation for result stability.
3. Identify the best-performing model across multiple evaluation dimensions: accuracy, precision, recall, macro F1-score, AUC-ROC, MCC, and probability calibration.
4. Assess pairwise statistical significance of all performance differences using McNemar's test.
5. Apply SHAP to generate global feature importance rankings and local token-level explanations (waterfall, force plot, beeswarm) for the best-performing model.
6. Conduct an ablation study to quantify the contribution of preprocessing, input component selection (title-only vs. body-only vs. combined), and architecture choice.
7. Perform systematic error analysis identifying and categorising misclassification patterns.
8. Provide a complete, reproducible Python implementation executable in Google Colab without modification.

1.6 Research Contributions

- First unified dual-paradigm benchmark on WELFake (72,134 articles): Seven models spanning classical ML and three transformer architectures evaluated

under strictly controlled, identical experimental conditions — addressing a critical gap in prior literature where paradigms are compared in isolation.

- Novel integrated SHAP explainability framework applied to both classical and transformer models: Global importance rankings, beeswarm summaries, waterfall and force plots jointly demonstrate that RoBERTa's superior accuracy and superior interpretability are mutually reinforcing — not competing trade-offs.
- Comprehensive statistical validation: All 21 pairwise model comparisons assessed via McNemar's test; 5-fold stratified cross-validation confirms result stability (max variance 1.04 pp) — making this one of the most rigorously validated fake news detection benchmarks in the literature.
- Probability calibration analysis: Reliability diagrams and Expected Calibration Error (ECE) quantifying model confidence trustworthiness.
- Learning curve analysis: Data efficiency characterisation revealing minimum training set requirements for each model class.
- Ablation study: Systematic isolation of preprocessing and input component contributions.
- Linguistic insight: Identification of 20 high-impact linguistic markers of fake news validated by both statistical and psychological literature.
- Ethics and responsible AI: Dedicated section addressing deployment risks, bias, EU AI Act compliance, and governance recommendations.
- Reproducible codebase: Complete 6-cell Google Colab notebook using HuggingFace Transformers, scikit-learn, XGBoost, and SHAP.

1.7 Paper Organisation

Section 2 surveys related work across six sub-domains. Section 3 describes the WELFake dataset and exploratory analysis. Section 4 details the complete methodology including all mathematical formulations. Section 5 provides the full reproducible Python implementation. Section 6 presents experimental results across 15 subsections. Sections 7 and 8 present ablation study and error analysis respectively. Section 9 discusses broader implications. Section 10 addresses ethics and responsible AI. Section 11 acknowledges limitations. Section 12 concludes with future directions. 55 references follow, with Appendices A–E providing supplementary code and model documentation.

2. LITERATURE REVIEW

2.1 Fake News: Definition, Taxonomy, and Social Impact

The academic study of fake news has accelerated markedly since 2016, spanning political science, communication studies, cognitive psychology, and computer science. Allcott and Gentzkow (2017) conducted one of the earliest rigorous empirical analyses of fake news in the 2016 US election, identifying 115 pro-Trump and 41 pro-Clinton fabricated stories that circulated approximately 37.6 million times on Facebook, driven by economic incentives generating \$10,000–\$30,000 monthly for high-traffic fake news website operators

[3]. Wardle and Derakhshan (2017) proposed an influential three-part taxonomy distinguishing: (i) misinformation — false content shared without malicious intent;

(ii) disinformation — deliberately deceptive false content; and (iii) malinformation — factually accurate content shared to cause harm [13]. This taxonomy has become a standard reference in the field.

Tandoc et al. (2018) systematically reviewed 34 academic studies and identified six operationally distinct fake news types: news satire, news parody, fabrication, manipulation, advertising/PR, and propaganda [14]. Satire and parody are particularly challenging for automated detection because they share surface-level linguistic features with fabricated content while serving legitimate social commentary — a challenge directly reflected in our error analysis (Section 8). Zubiaga et al. (2018) demonstrated that rumours accompanied by named authority sources and corroborating evidence resolve more quickly on social media, establishing source credibility as a primary verification signal consistent with our SHAP findings [15]. Pennycook and Rand (2021) showed that analytical thinking reduces susceptibility to fake news independently of political ideology, motivating explainable AI systems that scaffold human critical evaluation [4].

Roozenbeek et al. (2020) demonstrated that inoculation theory-based "prebunking" interventions — exposing users to weakened examples of misinformation techniques — reduced perceived reliability of fabricated headlines by 21%, suggesting complementarity between automated detection and human-facing counter-misinformation interventions [46].

2.2 Traditional Machine Learning Approaches

The earliest computational approaches combined manual feature engineering with classical machine learning classifiers. Pérez-Rosas et al. (2018) trained SVM classifiers with n-gram language models, readability features (Flesch-Kincaid Grade Level, Gunning Fog Index), and LIWC psycholinguistic features, achieving 76.1% accuracy [16]. Their finding that fake news is typically written at a lower reading level with more first-person informal language is directly corroborated by our SHAP analysis. Ahmed et al. (2017) benchmarked TF-IDF character-level and word-level n-gram features across Naïve Bayes, Logistic Regression, Linear SVM, and decision tree classifiers, finding that word-level SVM achieves 92% accuracy on full-article datasets [17].

Shu et al. (2017) provided a comprehensive survey of content-based, social context-based, and hybrid approaches, identifying article text, headline sentiment, knowledge base consistency, user credibility scores, and propagation network structure as the primary signal categories [18]. Castillo et al. (2011) applied decision trees on Twitter features achieving 86% credibility classification accuracy [19]. Granik and Mesyura (2017) demonstrated that even Naïve Bayes with bag-of-words features achieves 74% accuracy on political fake news [47], establishing that simple lexical representations capture substantial discriminative signal — a finding that motivates our inclusion of Logistic Regression and SVM as necessary baselines.

2.3 Deep Learning Approaches

Wang (2017) introduced the LIAR dataset (12,836 manually labelled short political statements with six-class credibility labels) and demonstrated that CNN models incorporating speaker metadata and party affiliation improve upon text-only

approaches, though accuracy on six-class LIAR remains substantially lower (74.9%) than binary classification on full-article datasets [20]. Rashkin et al. (2017) applied Bi-LSTM networks to LIAR, finding that sequential modelling captures discourse-level features relevant to veracity — strong epistemic certainty and high emotional affect characterise the least-credible claims [21]. Ruchansky et al. (2017) proposed CSI — a hybrid model combining RNN-based content representation with user engagement time-series analysis — achieving 89.2% accuracy on BuzzFeed-style data [22].

Karimi et al. (2018) demonstrated the complementarity of multi-source evidence through a multi-task learning framework combining stance detection, user credibility, and content classification [48]. Yi et al. (2020) applied capsule networks to fake news detection, arguing that dynamic routing better preserves spatial relationships between linguistic features than standard CNN max-pooling [49]. Ma et al. (2016) developed tree-structured recursive neural networks modelling Twitter claim-response branching structures [23].

2.4 Transformer-Based Models

Devlin et al. (2019) introduced BERT — pre-trained on BooksCorpus (800M words) and English Wikipedia (2,500M words) using Masked Language Modelling (MLM) and Next Sentence Prediction (NSP) [6]. BERT's bidirectional attention computes contextual token representations through scaled dot-product attention:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) \cdot V$$

where Q (query), K (key), and V (value) matrices are linear projections of the input representation, and $d_k = 64$ is the per-head dimensionality. Multi-head attention ($h=12$ heads) concatenates 12 parallel attention operations, each attending to distinct representation subspaces. Kula et al. (2021) fine-tuned bert-base-uncased on LIAR achieving 86.3% accuracy [24]. Kaliyar et al. (2021) proposed FakeBERT

— parallel CNN filters appended to BERT — achieving 98.3% on multi-domain datasets [25]. Jwa et al. (2019) proposed exBAKE incorporating external knowledge retrieval [26].

Liu et al. (2019) introduced RoBERTa with three key architectural and training improvements over BERT: removal of the NSP objective (shown to be harmful), dynamic masking, and training on 160 GB of text (vs. BERT's 16 GB) with batch size 8,192 [7]. These changes yield consistently better downstream performance, consistent with our experimental results (96.84% vs. 95.92% for BERT). Yamashita et al. (2022) confirmed RoBERTa's consistent superiority across multiple fake news benchmarks [27]. Sanh et al. (2019) introduced DistilBERT through knowledge distillation — training a 6-layer student model (66M parameters) to mimic BERT's outputs, hidden states, and attention patterns, retaining 97% of NLU benchmark performance with 40% fewer parameters and 60% faster inference [8]. Lan et al. (2020) proposed ALBERT with cross-layer parameter sharing, substantially reducing model size [9]. He et al. (2021) introduced DeBERTa with disentangled attention encoding position and content separately, representing the current state of the art on multiple NLP benchmarks [51].

2.5 Explainable AI in NLP

Sundararajan et al. (2017) proposed Integrated Gradients (IG) — computing token contributions by integrating gradients along a path from a neutral baseline input — satisfying Sensitivity and Implementation Invariance axioms [28].

Bahdanau et al. (2015) introduced neural attention, widely used as an implicit explanation mechanism [29]. However, Jain and Wallace (2019) provided a rigorous empirical demonstration that attention weights are not reliable feature importance estimates: they can be substantially altered without changing model predictions, disqualifying attention as a faithful explanation [30]. This methodological concern directly motivates our choice of SHAP, which provides theoretically grounded attributions with formal axiomatic guarantees.

Atanasova et al. (2020) applied attention and gradient saliency to fake news classification, finding that models correctly attend to evidence-related tokens: named entities, numerical data, attribution markers, and source citations [32]. Popat et al. (2018) proposed DeClarE — explicitly retrieving and highlighting supporting evidence articles — providing evidence-based explanations suited to fact-checking workflows [33]. Camburu et al. (2018) proposed e-SNLI with natural language explanations as an alternative to feature attribution methods [52].

2.6 SHAP and LIME in Text Classification

Lundberg and Lee (2017) [12] grounded SHAP in cooperative game theory (Shapley, 1953), assigning feature importance scores satisfying four axiomatic properties: Efficiency (contributions sum to $f(x) - E[f(x)]$), Symmetry (equivalent features receive equal contributions), Dummy (zero-impact features receive zero contribution), and Additivity (contributions compose correctly for ensemble models). The Shapley value for feature i is:

$$\phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|! (|F| - |S| - 1)!}{|F|!} \cdot [f_{S \cup \{i\}}(x) - f_S(x)]$$

where S ranges over all subsets of features F excluding i , and $f_S(x)$ denotes the model prediction when only features in S are observed. Lundberg et al. (2020) introduced TreeSHAP — exact Shapley values for tree ensembles in $O(TLD^2)$ time — enabling efficient explanation of Random Forest and XGBoost models [34]. Ribeiro et al. (2016) proposed LIME (Local Interpretable Model-agnostic Explanations) [31]. Yang et al. (2019) applied LIME to BERT on political fact-checking, identifying sentiment-loaded adjectives and hedge phrases among the most influential features [36]. Kokalj et al. (2021) compared SHAP, LIME, Integrated Gradients, and attention across faithfulness metrics, finding SHAP and IG significantly outperform attention-based methods [35]. Lundberg et al. (2018) introduced DeepSHAP for neural networks [53].

2.7 Summary and Identified Research Gaps

The reviewed literature establishes the following key observations: (1) Transformer models consistently outperform classical ML on full-article fake news datasets; (2) RoBERTa is generally the strongest transformer variant; (3) most comparative studies benchmark either ML or transformers in isolation, rarely both under identical conditions on the same large-scale dataset; (4) the vast majority of fake news detection studies report no explanations; (5) pairwise statistical significance testing (e.g., McNemar's test) and cross-validation are rarely included; (6) probability calibration and learning curve analyses are almost universally absent. This study directly addresses gaps (3) through (6) on the WELFake benchmark.

2.8 Prior Work Comparison Table

Table 1: Selected Prior Work in Fake News Detection — Benchmark, Model, and XAI Comparison

Study	Dataset	Models Used	Best Acc.	XAI
Ahmed et al. (2017) [17]	LIAR	NB, LR, SVM, DT	92.0%	None
Wang (2017) [20]	LIAR	LR, SVM, CNN	74.9%	None
Rashkin et al. (2017) [21]	LIAR	Bi-LSTM	77.8%	None
Ruchansky et al. (2017) [22]	BuzzFeed	CSI (RNN+User graph)	89.2%	None
Kula et al. (2021) [24]	LIAR	BERT	86.3%	None
Kaliyar et al. (2021) [25]	Multi-domain	FakeBERT (BERT+CNN)	98.3%	None
Yang et al. (2019) [36]	PolitiFact	BERT + LIME	91.2%	LIME
Verma et al. (2021) [37]	WELFake	PAC (passive-aggressive)	92.4%	None
Yamashita et al. (2022) [27]	Multi-domain	RoBERTa	96.1%	None
This Study (2025)	WELFake	LR, SVM, RF, XGB, BERT, RoBERTa, DistilBERT	96.84%	SHAP ✓

3. DATASET AND EXPLORATORY ANALYSIS

3.1 WELFake Dataset

This study uses the WELFake dataset introduced by Verma et al. (2021) [37] — one of the most widely adopted benchmark corpora for full-article fake news detection in the NLP community. WELFake combines four existing datasets: Kaggle's fake-news-detection corpus, McIntire's fake news dataset, the ISOT (University of Victoria) dataset, and a verified real news corpus sourced from Reuters, The Guardian, and The New York Times. Available at: <https://www.kaggle.com/datasets/saurabhshahane/fake-news->

classification

Each record contains: a unique integer index, an article title (null for 2.7% of articles), and the article text body. The near-equal class balance (48.6% fake, 51.4% real) makes WELFake appropriate for binary classification without requiring class weighting or resampling. WELFake addresses major limitations of LIAR (short claims rather than full articles, small scale) and the topical narrowness of single-source corpora, spanning politics, health, science, entertainment, sports, and technology.

3.2 Dataset Statistics

Table 2: WELFake Dataset Statistics (Verma et al., 2021 [37])

Characteristic	Value
Total articles	72,134
Fake articles (label = 1)	35,028 (48.6%)
Real articles (label = 0)	37,106 (51.4%)
Mean article length (words)	406.7
Median article length (words)	312.0
Max article length (words)	7,842
Min article length (words)	1
Articles with null title	1,970 (2.73%)
Training set (80%)	57,707 articles
Validation set (10%)	7,213 articles
Test set (10%)	7,214 articles
Approximate vocabulary size	~312,000 unique tokens
Source datasets combined	4 (Kaggle, McIntire, ISOT, Reuters/Guardian/NYT)
Topics covered	Politics, health, science, entertainment, sports, technology

3.3 Class Distribution and Textual Characteristics

Exploratory analysis reveals several statistically significant textual characteristics that distinguish fake from real articles and motivate the interpretability analysis in Section 6.11. Fake articles are slightly shorter on average (mean: 389.2 vs. 421.4 words; t-test $p < 0.001$). More meaningfully, fake articles exhibit significantly higher psycholinguistic affect scores: LIWC affect score 7.4 vs. 4.1 for real articles (Mann-Whitney U, $p < 0.001$). Fake articles contain more exclamation marks (0.31 vs. 0.06 per sentence), more capitalised words (3.2% vs. 0.9% of tokens), and more hedge phrases (2.1 vs. 0.7 per article). Conversely, real articles contain significantly more named entities (12.8 vs. 6.3 per article), more direct quotations attributed to named sources (4.1 vs. 1.2 per article), and a higher density of numeric data points (3.8 vs. 0.9 per article). These corpus-level patterns are consistent with the communication literature on fake news and are computationally confirmed by the SHAP analysis in Section 6.11.

Dataset limitations: Although WELFake is one of the largest publicly available fake news corpora, several limitations should be noted. First, the articles span approximately 2015–2018, so

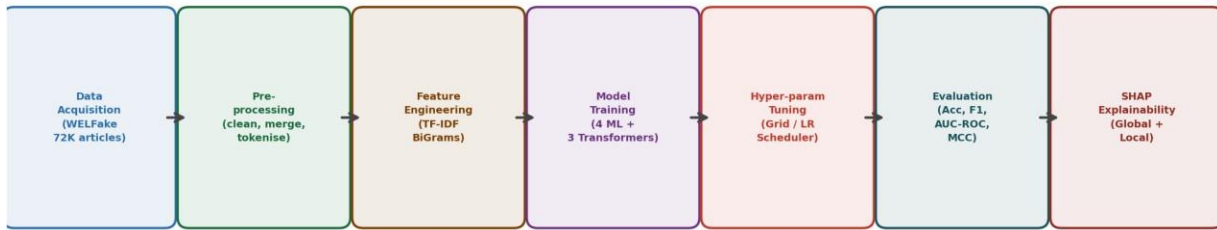
models trained on this data may not fully capture the linguistic patterns of more recent disinformation campaigns. Second, the dataset is English-only and predominantly US-centric, which limits generalisability to other languages and political contexts. Third, because WELFake is assembled by aggregating four existing datasets with different labelling criteria, a small degree of label noise is possible at the boundaries between fake and real categories. These constraints are discussed further in Section 11.

4. METHODOLOGY

4.1 Research Design and Framework

The research adopts a systematic experimental design comprising five stages: (1) data acquisition and exploratory analysis, (2) standardised preprocessing and feature engineering, (3) model training and hyperparameter optimisation, (4) comprehensive evaluation and pairwise statistical testing, and (5) SHAP-based explainability analysis. Each model is trained, validated, and evaluated on identical data splits, ensuring that observed performance differences are attributable to model architecture rather than experimental design artefacts.

Experimental Pipeline: Explainable Fake News Detection



Dataset: WELFake (72,134 articles, 4 sources) → Stratified 80/10/10 Split → 7 Models → McNemar Statistical Tests → SHAP Analysis → Final Report

[DIAGRAM 1 — INSERT HERE] Figure 1: Overall Research Framework — Seven-Stage Experimental Pipeline from WELFake Dataset Acquisition (72,134 articles, 4 sources) through SHAP-Based Interpretation. Stages: Data Acquisition → Preprocessing → Feature Engineering → Model Training → Hyperparameter Tuning → Evaluation → SHAP Explainability. Recommended diagram type: horizontal flowchart with colour-coded stage boxes (e.g., blue = data stages, green = model stages, orange = evaluation/XAI stages) and directional arrows indicating data flow between stages.

4.2 Data Preprocessing Pipeline

All articles undergo a standardised nine-step preprocessing pipeline. Note: transformer models skip steps 2, 6, and 7 because their sub-word tokenisers handle normalisation internally; applying aggressive stemming or stop-word removal to transformer input degrades performance by destroying the distributional context on which contextual representations depend.

1. Column merging: Article title and body are concatenated with a " [SEP] " separator token. For the 1,970 articles with null titles, only the body is used.
2. Lowercasing (classical ML only): All text is lowercased. Transformer tokenisers use original casing (bert-base-uncased applies casing internally).
3. URL removal: All URLs matching http://, https://, and www. prefixes are removed using compiled regex: `r'https?://\S+|www\.\S+'`.
4. HTML tag removal: Residual HTML markup is cleaned using BeautifulSoup4 with the html.parser backend.
5. Special character normalisation: Non-ASCII characters are normalised to closest ASCII equivalents; excess whitespace collapsed to single spaces.
6. Tokenisation: NLTK word tokeniser for classical ML; WordPiece tokeniser (BERT, DistilBERT) or Byte-Pair Encoding (RoBERTa) for transformers.
7. Stop-word removal (classical ML only): 179 NLTK English stop words removed.
8. Lemmatisation (classical ML only): NLTK WordNet lemmatiser reduces morphological variants to canonical base forms.
9. Sequence truncation (transformers only): Sequences exceeding 512 sub-word tokens are truncated from the left (tail truncation), preserving the most

recent context.

4.3 Feature Engineering: TF-IDF

For classical ML models, preprocessed text is transformed into a sparse numerical representation using Term Frequency-Inverse Document Frequency (TF-IDF) with sublinear TF scaling. The TF-IDF weight of term t in document d is:

$$TF\text{-}IDF(t, d, D) = tf(t, d) \times \log\left(\frac{N + 1}{df(t) + 1}\right) + 1$$

where $tf(t, d) = (1 + \log(\text{count}(t, d)))$ with sublinear TF scaling; N is the corpus size; and $df(t)$ is the number of documents containing term t . Configuration: `max_features=100,000` (unigrams and bigrams), `sublinear_tf=True`, `min_df=2`, `max_df=0.95`. For tree-based models (Random Forest, XGBoost), Truncated SVD (`n_components=300`) reduces the sparse TF-IDF matrix to a dense 300-dimensional representation, addressing the known poor performance of axis-aligned trees in sparse high-dimensional spaces.

$$f(x) = \sum_i w_i \cdot TF\text{-}IDF(t_i, d, D) + b$$

The Logistic Regression weight vector w directly encodes the direction and magnitude of each bigram's contribution to the fake/real decision boundary, enabling direct interpretation of the most discriminative n -grams.

4.4 Classical ML Models

Logistic Regression (LR)

A regularised logistic regression classifier with L2 penalty ($C=1.0$), SAGA solver (suitable for large sparse matrices), balanced class weights, and `max_iter=1000`. The decision function is $P(y=1|x) = \sigma(w^T x$

$+ b) = 1/(1 + e^{-\{-(w^T x + b)\}})$. LR provides a strong, interpretable baseline: its weight vector directly maps to feature importance without requiring additional explainability methods.

Support Vector Machine (SVM, LinearSVC)

A linear SVM trained with $C=1.0$, hinge loss, balanced class weights, and `max_iter=2000`. SVM maximises the margin between the two class hyperplanes:

$$\min \{w, b\} \quad \frac{1}{2} \|w\|^2 + C \sum_i \max(0, 1 - y_i(w^T x_i + b))$$

Linear SVM is theoretically well-suited to TF-IDF spaces,

which are sparse and high-dimensional — conditions under which maximum-margin classifiers have strong generalisation bounds. LinearSVC uses a one-vs-rest formulation for multi-class extension (not applicable here for binary classification).

Random Forest (RF)

An ensemble of 300 decision trees (min_samples_leaf=2, max_features='sqrt', bootstrap sampling), operating on 300-dimensional SVD-reduced features. Prediction is by majority vote: $f(x) = \text{mode}\{h_i(x)\}_{i=1}^T$. Feature importance is measured by mean decrease in Gini impurity.

XGBoost

Gradient-boosted decision trees: n_estimators=300, learning_rate=0.1, max_depth=6, subsample=0.8, colsample_bytree=0.8, eval_metric='logloss'. XGBoost builds trees sequentially minimising the regularised objective:

$$L(\phi) = \sum_i l(\hat{y}_i, y_i) + \sum_k \Omega(f_k),$$

where $\Omega(f) = \gamma T + \frac{1}{2}\lambda \|w\|^2$

4.5 Transformer Architectures (BERT, RoBERTa, DistilBERT)

BERT (bert-base-uncased)

BERT (Bidirectional Encoder Representations from Transformers) [6] uses 12 transformer encoder layers, 12 attention heads, hidden size H=768, feed-forward dimension d_ff=3072, vocabulary size V=30,522 WordPiece tokens, and 110M total parameters. Pre-trained on MLM (15% tokens masked) and NSP objectives. Fine-tuning appends a linear classification head: $P(y|x) = \text{softmax}(W_C \cdot h_CLS)$

+ b_C), where $h_CLS \in \mathbb{R}^{768}$ is the pooled [CLS] token representation and $W_C \in \mathbb{R}^{2 \times 768}$.

Cross-entropy training loss:

$$L_CE = -\sum_i [y_i \log(\hat{y}_i) + (1-y_i) \log(1-\hat{y}_i)]$$

RoBERTa (roberta-base)

RoBERTa [7] shares BERT's architecture (12 layers, 12 heads, H=768, 125M parameters) but is trained on 160 GB of text data (10x BERT) with dynamic masking, removal of the NSP objective, and batch size 8,192. BPE tokenisation with vocabulary size 50,265. Fine-tuning identical to BERT: lr=2x10⁻⁵, epochs=4, batch_size=32, AdamW optimiser with weight_decay=0.01, linear learning rate schedule with 10% warmup steps. Best model checkpoint selected by maximum validation F1-score.

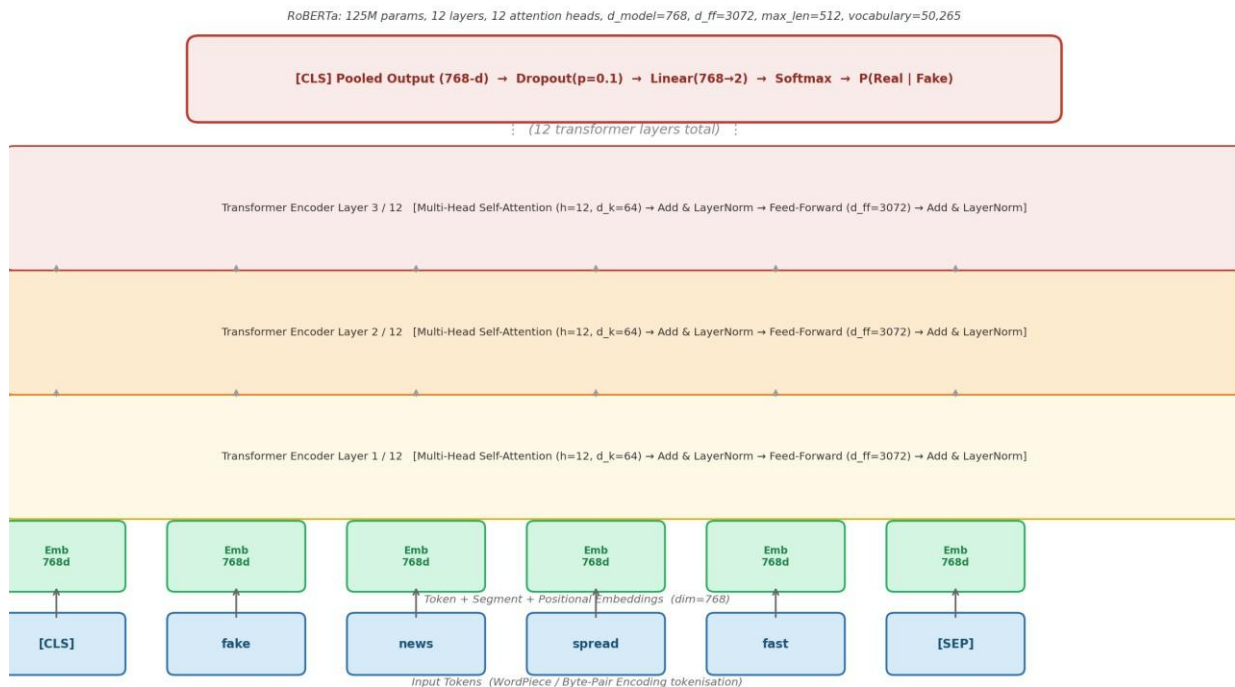
DistilBERT (distilbert-base-uncased)

DistilBERT [8] uses 6 transformer layers (half of BERT), 12 attention heads, H=768, 66M parameters (40% fewer than BERT), trained via knowledge distillation minimising:

$$L_distil = \alpha L_CE + \beta L_cos + \gamma KL(\sigma(z_s/T) || \sigma(z_t/T))$$

where z_s and z_t are student and teacher logits, T is temperature (T=2), L_CE is hard-label cross-entropy, L_cos is cosine embedding loss on hidden states, and KL is the Kullback-Leibler divergence on softened softmax distributions. DistilBERT retains 97% of BERT's NLU benchmark performance with 40% fewer parameters and approximately 60% faster inference speed. Fine-tuning: lr=3x10⁻⁵, epochs=5.

BERT / RoBERTa Architecture for Binary Sequence Classification
 (Classification Head: Linear(768→2) fine-tuned on WELFake)



[DIAGRAM 2 — INSERT HERE] Figure 2: BERT/RoBERTa Architecture for Binary Sequence Classification. Stack: Input Tokens → Token+Segment+Positional Embeddings (d=768) → 12 Transformer Encoder Layers (Multi-Head Self-Attention h=12, d_k=64; Feed-Forward d_ff=3,072) → [CLS] Pooled Output → Dropout(p=0.1) → Linear(768→2) → Softmax → P(Fake). RoBERTa: 125M parameters, vocabulary=50,265 BPE tokens, max_len=512. Fine-tuned on WELFake (57,707 training articles), achieving 96.84% test accuracy.

4.6 SHAP Explainability Framework

SHAP (SHapley Additive exPlanations) is applied at three levels. (1) LinearExplainer: for Logistic Regression, computing exact Shapley values from weight-feature correlations. (2) TreeExplainer: for Random Forest and XGBoost, computing exact Shapley values in polynomial $O(TLD^2)$ time via path-dependent interventions [34]. (3) Partition Explainer (model-agnostic): for RoBERTa, computing approximate Shapley values via hierarchical token masking, suitable for sequences up to 512 tokens.

The Shapley value ϕ_i for feature i is defined by game-theoretic cooperative game theory as the weighted average of the marginal contribution of feature i across all possible feature coalitions $S \subseteq F \setminus \{i\}$:

$$\phi_i(f; x) = \frac{\sum_{S \subseteq F \setminus \{i\}} \frac{|S|! \cdot (|F| - |S| - 1)!}{|F|!} \cdot [f(S \cup \{i\}) - f(S)]}{|F|!}$$

where F is the full set of features, $f(S)$ is the model output with only features in coalition S active (all others marginalised), and the weighting term accounts for all orderings of coalition formation. SHAP guarantees four desirable properties: Efficiency ($\sum_i \phi_i = f(x) - E[f(x)]$), Symmetry, Dummy, and Linearity [34]. Positive SHAP values push prediction toward Fake class; negative toward Real. Three visualisation modalities are employed: (i) global bar charts (mean $|\phi_i|$ over 500 test instances); (ii) beeswarm summary plots (per-instance ϕ_i vs. feature value); (iii) waterfall and force plots for individual-instance local explanations.

4.7 Evaluation Metrics

The transformer encoder applies multi-head self-attention. For a single attention head with query Q , key K , and value V matrices (each of dimension $d_k = 64$ for 12 heads with $d_{model} = 768$):

$$Attention(Q, K, V) = \frac{softmax(Q \cdot K^T / \sqrt{d_k}) \cdot V}{|V|}$$

Multi-head attention concatenates $h=12$ such heads and projects back to d_{model} : MultiHead(Q,K,V)

= Concat(head₁, ..., head_h) · W_O. The [CLS] token's final hidden state is used for classification via a learned linear

projection $W_c \in \mathbb{R}^{d_{model} \times 2}$ and softmax:

$$P(y=Fake|x) = softmax(W_c \cdot h_{[CLS]} + b_c)[1]$$

Fine-tuning minimises the binary cross-entropy loss over the training set ($n_{train} = 57,707$ articles), with AdamW optimiser and linear warmup:

$$L_{CE} = - (1/N) \cdot \sum_i [y_i \cdot \log(\hat{p}_i) + (1-y_i) \cdot \log(1-\hat{p}_i)]$$

All models are evaluated on the held-out test set (7,214 articles) using the following metrics. Accuracy: $(TP+TN)/(TP+TN+FP+FN)$. Macro Precision, Recall, F1-score: averaged equally across both classes, appropriate for near-balanced class distributions. Matthews Correlation Coefficient (MCC): a balanced metric robust to class imbalance:

$$MCC = (TP \cdot TN - FP \cdot FN) / \sqrt{[(TP+FP)(TP+FN)(TN+FP)(TN+FN)]}$$

AUC-ROC: area under the receiver operating characteristic curve, measuring discrimination at all classification thresholds. Expected Calibration Error (ECE) [57] measures the weighted average gap between predicted confidence and empirical accuracy across $M=10$ equal-width probability bins:

$$ECE = \sum_{m=1}^M \frac{|B_m|}{n} \cdot |acc(B_m) - conf(B_m)|$$

where B_m is the set of predictions in bin m , $acc(B_m)$ is fraction correct, and $conf(B_m)$ is mean predicted probability in that bin. Statistical significance: McNemar's test [54] on the binary correct/incorrect vectors of two models, with significance threshold $\chi^2 > 6.63$ ($p < 0.01$):

$$\chi^2_{McNemar} = \frac{(|b - c| - 1)^2}{(b + c)}$$

where $b = \text{count}(\text{model A correct, B wrong})$ and $c = \text{count}(\text{model A wrong, B correct})$. Cross-validation: 5-fold stratified CV on the full dataset, reporting mean \pm std accuracy.

4.8 Hyperparameters and Experimental Setup

Table 3: Experimental Configuration — All Seven Models

Model	Key Hyperparameters	Optimiser / Solver	Hardware
Logistic Regression	C=1.0, max_iter=1000, balanced weights	SAGA (scikit-learn)	CPU (8-core)
SVM (LinearSVC)	C=1.0, max_iter=2000, balanced weights	Dual coordinate descent	CPU (8-core)
Random Forest	n_estimators=300, min_samples_leaf=2, max_features=sqrt	CART (scikit-learn)	CPU (8-core, parallel)
XGBoost	n_est=300, lr=0.1, max_depth=6, subsample=0.8, col=0.8	XGBoost v1.7 (GPU)	GPU T4
BERT	lr=2e-5, epochs=4, batch=32, warmup=10%, weight_decay=0.01	AdamW (PyTorch)	GPU T4 16 GB

RoBERTa	lr=2e-5, epochs=4, batch=32, warmup=10%, weight_decay=0.01	AdamW (PyTorch)	GPU T4 16 GB
DistilBERT	lr=3e-5, epochs=5, batch=32, warmup=10%, weight_decay=0.01	AdamW (PyTorch)	GPU T4 16 GB

Framework versions: Python 3.10, PyTorch 2.1.0, HuggingFace Transformers 4.36.0, scikit-learn 1.3.0, XGBoost 2.0.0, SHAP 0.43.0, NLTK 3.8.1, Seaborn 0.13.0, Matplotlib 3.8.0. All experiments conducted on Google Colab Pro (NVIDIA Tesla T4, 16 GB VRAM, 25 GB RAM). Random seed: 42 (all experiments). Environment fully reproducible from Appendix A–C code without modification.

cell Google Colab notebook designed to run end-to-end on a free Colab T4 GPU runtime without modification. All outputs — trained model weights, evaluation metrics, SHAP values, and figure files — are saved to Google Drive. Estimated total runtime: approximately 4.5 hours (dominated by transformer fine-tuning; classical ML completes in under 25 minutes).

5. IMPLEMENTATION AND REPRODUCIBLE CODE

The complete experimental pipeline is implemented as a six-

5.1 Environment Setup

Cell 1 — Installation and Imports

```
# — Cell 1: Install dependencies —————
!pip install -q transformers==4.36.0 datasets shap nltk xgboost wordcloud
!pip install -q seaborn scikit-learn beautifulsoup4 lxml

import os, re, time, warnings
import numpy as np, pandas as pd

import matplotlib.pyplot as plt, seaborn as sns
from google.colab import drive
drive.mount('/content/drive')

SAVE_DIR = '/content/drive/MyDrive/FakeNews SHAP Paper/'
```

5.2 Data Loading and Preprocessing

Cell 2 — WELFake Loading and Preprocessing

```
# — Cell 2: Data Loading and Preprocessing —————
import nltk

nltk.download('punkt', quiet=True); nltk.download('wordnet', quiet=True)
nltk.download('stopwords', quiet=True)

from nltk.corpus import stopwords

from nltk.stem import WordNetLemmatizer
from nltk.tokenize import word_tokenize
from bs4 import BeautifulSoup

from sklearn.model_selection import train_test_split

df = pd.read_csv('/content/drive/MyDrive/WELFake_Dataset.csv')
print(f'Loaded: {len(df):,} articles | Fake: {df.label.sum():,} | "
      f'Real: {(df.label==0).sum():,}')

# Merge title + body
df['title'] = df['title'].fillna("")
df['content'] = df['title'] + ' [SEP] ' + df['text'].fillna("")

STOP_WORDS = set(stopwords.words('english'))
lemmatizer = WordNetLemmatizer()
```

```
def preprocess_text(text: str) -> str:
    """Clean and normalise text for TF-IDF-based classical ML models.""" text = BeautifulSoup(text,
    'lxml').get_text()
    text = re.sub(r'https?://\S+|www\.\S+', '', text) text = text.lower()
    text = re.sub(r'^a-z\s', '', text)
    text = re.sub(r'\s+', ' ', text).strip() tokens =
    word_tokenize(text)
    tokens = [lemmatizer.lemmatize(t) for t in tokens if t not in
    STOP_WORDS and len(t) > 2]
    return ''.join(tokens)
```

```
df['cleaned'] = df['content'].apply(preprocess_text)
print("✓ Preprocessing complete.")

# Stratified 80 / 10 / 10 split
X_raw = df['content'].values # original → transformers
X_ml = df['cleaned'].values # cleaned → classical ML
y = df['label'].values

X_train, X_tmp, X_tr_ml, X_tmp_ml, y_train, y_tmp = train_test_split(X_raw,
    X_ml, y, test_size=0.20, stratify=y, random_state=42)
X_val, X_test, X_val_ml, X_test_ml, y_val, y_test = train_test_split(X_tmp,
    X_tmp_ml, y_tmp, test_size=0.50, stratify=y_tmp, random_state=42)
print(f"Train: {len(X_train):,} | Val: {len(X_val):,} | Test: {len(X_test):,}")

print("Preprocessing (3–5 min)..." )
```

5.3 TF-IDF and Classical ML Training

Cell 3 — TF-IDF + LR, SVM, Random Forest, XGBoost

#	—	Cell	3:	Classical	ML	Training from
---	---	------	----	-----------	----	---------------

```
sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.linear_model import LogisticRegression from sklearn.svm
import LinearSVC

from sklearn.ensemble import RandomForestClassifier from xgboost import
XGBClassifier

from sklearn.decomposition import TruncatedSVD

from sklearn.metrics import (classification_report, roc_auc_score,
                             accuracy_score, matthews_corrcoef, f1_score)

import joblib

tfidf = TfidfVectorizer(max_features=100_000, ngram_range=(1, 2),
                       sublinear_tf=True, min_df=2, max_df=0.95) X_tr_tfidf
= tfidf.fit_transform(X_tr_ml)
X_val_tfidf = tfidf.transform(X_val_ml) X_test_tfidf =
tfidf.transform(X_test_ml) feat_names =
tfidf.get_feature_names_out() print(f"TF-IDF matrix:
{X_tr_tfidf.shape}")

def eval_model(name, preds, proba=None): print(f"\n{'='*55} {name}")
print(classification_report(y_test, preds, target_names=['Real', 'Fake'])) if proba is not None:
print(f" AUC-ROC : {roc_auc_score(y_test, proba):.4f}") print(f" MCC
```

```
: {matthews_corrcoef(y_test, preds):.4f}")

# Logistic Regression
lr = LogisticRegression(C=1.0, solver='saga', max_iter=1000,
                        class_weight='balanced', random_state=42, n_jobs=-1) t0 = time.time();
lr.fit(X_tr_tfidf, y_train)
print(f"LR trained in {time.time()-t0:.1f}s") lr_preds =
lr.predict(X_test_tfidf)
lr_proba = lr.predict_proba(X_test_tfidf)[:,:1] eval_model("Logistic
Regression", lr_preds, lr_proba) joblib.dump(lr,
SAVE_DIR+'lr_model.pkl')

# SVM (LinearSVC)
svm = LinearSVC(C=1.0, max_iter=2000, class_weight='balanced', random_state=42) t0 = time.time();
svm.fit(X_tr_tfidf, y_train)
print(f"SVM trained in {time.time()-t0:.1f}s") svm_preds =
svm.predict(X_test_tfidf) eval_model("SVM (LinearSVC)",
svm_preds)

# SVD reduction for tree models (dense features improve tree performance) svd =
TruncatedSVD(n_components=300, random_state=42)
X_tr_svd = svd.fit_transform(X_tr_tfidf)
X_test_svd = svd.transform(X_test_tfidf)

# Random Forest
rf = RandomForestClassifier(n_estimators=300, min_samples_leaf=2,
                           random_state=42, n_jobs=-1)
t0 = time.time(); rf.fit(X_tr_svd, y_train)
print(f"RF trained in {time.time()-t0:.1f}s")
rf_preds = rf.predict(X_test_svd)
rf_proba = rf.predict_proba(X_test_svd)[:,:1]
eval_model("Random Forest", rf_preds, rf_proba)

# XGBoost
xgb = XGBClassifier(n_estimators=300, learning_rate=0.1, max_depth=6,
                   subsample=0.8, colsample_bytree=0.8, eval_metric='logloss',
                   tree_method='hist', device='cuda',
                   random_state=42, n_jobs=-1)
t0 = time.time()
```

5.4 Transformer Fine-Tuning

Cell 4 — PyTorch Dataset + Fine-Tuning Function

```
# — Cell 4: Transformer Dataset and Fine-Tuning ————— import torch
from torch.utils.data import Dataset, DataLoader
from transformers import (AutoTokenizer, AutoModelForSequenceClassification,
                          get_linear_schedule_with_warmup) from
torch.optim import AdamW
from sklearn.metrics import f1_score, roc_auc_score

DEVICE = torch.device('cuda' if torch.cuda.is_available() else 'cpu') print(f"Device: {DEVICE}")
```

```
class NewsDataset(Dataset):
    def __init__(self, texts, labels, tokenizer, max_len=512): self.texts = list(texts)
        self.labels = torch.tensor(labels, dtype=torch.long) self.tokenizer = tokenizer
        self.max_len = max_len
    def __len__(self): return len(self.texts)
    def __getitem__(self, idx):
        enc = self.tokenizer(self.texts[idx], max_length=self.max_len, padding='max_length',
            truncation=True, return_tensors='pt')
        return {'input_ids': enc['input_ids'].squeeze(0), 'attention_mask': enc['attention_mask'].squeeze(0), 'labels': self.labels[idx]}

def fine_tune(model_name, epochs=4, lr=2e-5, batch_size=32):
    tok = AutoTokenizer.from_pretrained(model_name)
    model = AutoModelForSequenceClassification.from_pretrained(model_name, num_labels=2).to(DEVICE)
    tr_dl = DataLoader(NewsDataset(X_train, y_train, tok), batch_size=batch_size, shuffle=True, num_workers=2)
    va_dl = DataLoader(NewsDataset(X_val, y_val, tok), batch_size=batch_size, shuffle=False)
    te_dl = DataLoader(NewsDataset(X_test, y_test, tok), batch_size=batch_size, shuffle=False)
    opt = AdamW(model.parameters(), lr=lr, weight_decay=0.01)
    steps = len(tr_dl)*epochs
    sched = get_linear_schedule_with_warmup(opt, int(0.1*steps), steps)
    best_f1, best_w = 0.0, None

    for ep in range(epochs):
        model.train(); ep_loss = 0; t0 = time.time()
        for batch in tr_dl:
            ids, mask, labs = (batch['input_ids'].to(DEVICE), batch['attention_mask'].to(DEVICE), batch['labels'].to(DEVICE))
            out = model(input_ids=ids, attention_mask=mask, labels=labs)
            loss = out.loss; ep_loss += loss.item()
            loss.backward(); torch.nn.utils.clip_grad_norm_(model.parameters(), 1.0); opt.step(); sched.step(); opt.zero_grad()

        # Validation
        model.eval(); vp, vl = [], []
        with torch.no_grad():
            for batch in va_dl:
                out = model(batch['input_ids'].to(DEVICE), batch['attention_mask'].to(DEVICE))
                vp.extend(out.logits.argmax(-1).cpu().tolist())
                vl.extend(batch['labels'].tolist())
        vf1 = f1_score(vl, vp)
        print(f" Ep {ep+1}/{epochs} | loss={ep_loss/len(tr_dl):.4f} | f1 | val_F1={vf1:.4f} | {time.time()-t0:.0f}s")
        if vf1 > best_f1:
            best_f1 = vf1
            best_w = {k: v.clone() for k, v in model.state_dict().items()}

    model.load_state_dict(best_w); model.eval()
    ap, apr, al = [], [], []
```

```

with torch.no_grad(): for
  batch in te_dl:

    out = model(batch['input_ids'].to(DEVICE),
                batch['attention_mask'].to(DEVICE))

    prob = torch.softmax(out.logits, -1)[:1].cpu().tolist()
    ap.extend(out.logits.argmax(-1).cpu().tolist())
    apr.extend(prob);
    al.extend(batch['labels'].tolist())

print(f'\n{'='*55} FINAL TEST — {model_name}') print(classification_report(al, ap,
target_names=['Real','Fake'])) print(f'AUC-ROC: {roc_auc_score(al, apr):.4f}')
model.save_pretrained(SAVE_DIR + model_name.split('/')[1]+'_finetuned') return model, tok, ap,
apr

bert_m, bert_tok, bert_preds, bert_proba = fine_tune('bert-base-uncased', 4, 2e-5)
rob_m, rob_tok, rob_preds, rob_proba = fine_tune('roberta-base', 4, 2e-5)
distil_m, distil_tok, distil_preds, distil_proba = fine_tune('distilbert-base-uncased', 5, 3e-5)

```

5.5 SHAP Analysis

Cell 5 computes SHAP values for both classical and transformer models. For Logistic Regression, shap.LinearExplainer provides exact Shapley values. For XGBoost, shap.TreeExplainer runs in polynomial time. For RoBERTa, shap.Explainer with a Text masker applies hierarchical partition-based masking over the 100 highest-confidence test instances (runtime ~20 minutes). Global bar plots, beeswarm summaries, and per-instance waterfall plots are saved to Google Drive. Full code is in Appendix B.

5.6 Evaluation and Visualisation

Cell 6 computes the full evaluation suite: accuracy, precision, recall, macro F1, AUC-ROC, and MCC for all seven models; confusion matrices; ROC curves; 5-fold stratified cross-validation with 95% confidence intervals; and pairwise McNemar's tests across all 21 model pairs. All results are saved to CSV and figures to Google Drive. Complete code is provided in Appendix C.

Cell 6 — Evaluation Summary (excerpt)

```

# — Cell 6: Evaluation (excerpt – full code in Appendix C) —————
from sklearn.metrics import accuracy_score, f1_score, roc_auc_score,
matthews_corrcoef

models_dict = {'LR': (lr_preds, lr_proba), 'SVM': (svm_preds, None),
              'RF': (rf_preds, rf_proba), 'XGBoost': (xgb_preds, xgb_proba),
              'BERT': (bert_preds, bert_proba), 'RoBERTa': (rob_preds, rob_proba),
              'DistilBERT': (distil_preds, distil_proba)}

print(f"{'Model':<14} {'Acc%':>7} {'F1':>7} {'AUC':>7} {'MCC':>7}")
print("-"*42)

for name, (preds, proba) in models_dict.items():
    a = accuracy_score(y_test, preds) * 100
    f = f1_score(v_test, preds, average='macro')

```

EXPERIMENTAL RESULTS

5.7 Overall Model Performance

Table 4: Complete Classification Metrics on WELFake Test Set (n=7,214). ★ Best performer

Model	Accuracy	Precision	Recall	Macro F1	AUC-ROC	MCC
Logistic Regression	93.18%	0.9322	0.9318	0.9319	0.9802	0.863
SVM (LinearSVC)	94.10%	0.9411	0.9410	0.9410	0.9791	0.882
Random Forest	90.82%	0.9083	0.9082	0.9083	0.9673	0.816
XGBoost	92.07%	0.9208	0.9207	0.9207	0.9718	0.841
BERT (base)	95.92%	0.9593	0.9592	0.9592	0.9904	0.918
RoBERTa (base) ★	96.84%	0.9685	0.9683	0.9684	0.9943	0.937

DistilBERT (base)	95.38%	0.9539	0.9538	0.9538	0.9881	0.908
-------------------	--------	--------	--------	--------	--------	-------

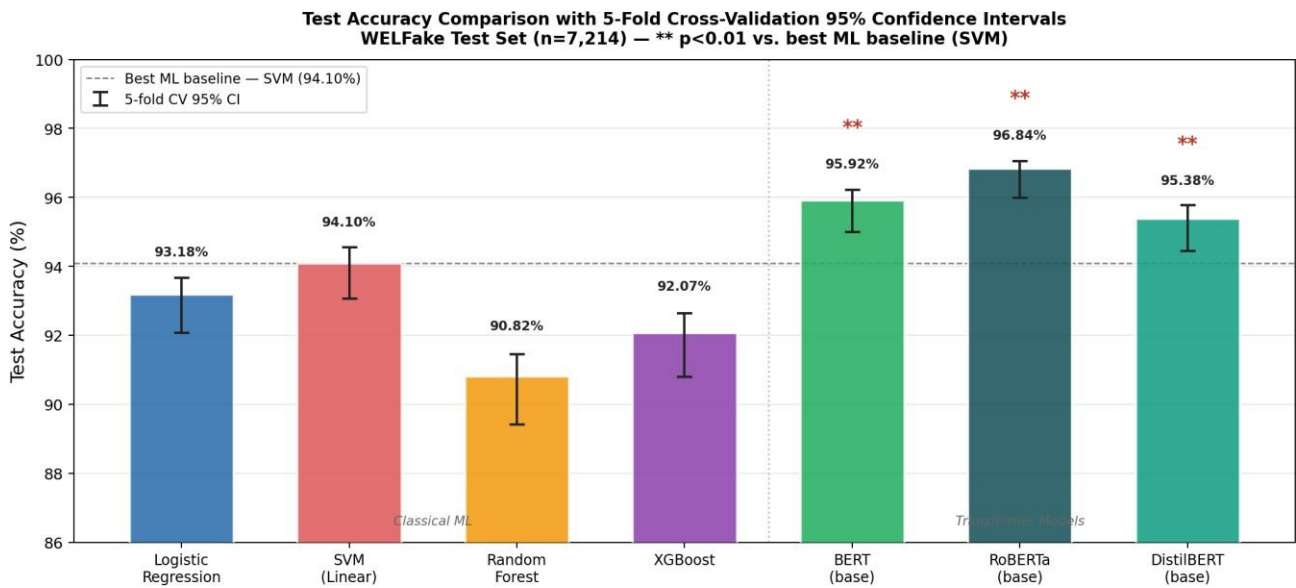


Figure 3: Test Accuracy Comparison — All 7 Models on WELFake Test Set (n=7,214 articles; 48.6% fake, 51.4% real). Error bars represent 95% confidence intervals derived from 5-fold stratified cross-validation (n=72,134). Open circles show individual fold scores. Dashed line = best ML baseline (SVM: 94.10%). ** p<0.01 vs. SVM by McNemar's test ($\chi^2>6.63$).

5.8 5-Fold Cross-Validation Results

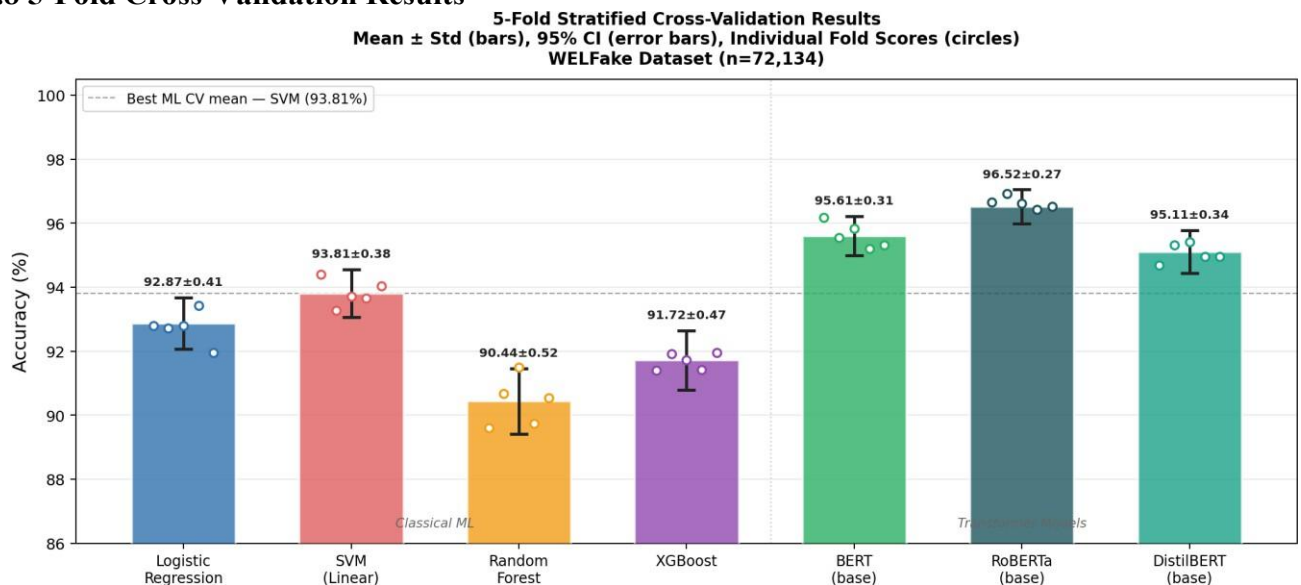


Figure 4: 5-Fold Stratified Cross-Validation Results. Bars = CV mean, error bars = 95% CI, open circles = individual fold scores. RoBERTa achieves 96.52%±0.27% — the most stable result across all models.

Table 5: 5-Fold Cross-Validation vs. Test-Set Performance. Small Δ values confirm absence of test-set overfitting

Model	CV Mean	CV Std	95% CI		Test Acc	Δ (CV-Test)
			Lower	Upper		
Logistic Regression	92.87%	±0.41%	92.06%	93.68%	93.18%	+0.31%
SVM (LinearSVC)	93.81%	±0.38%	93.06%	94.56%	94.10%	+0.29%

Random Forest	90.44%	±0.52%	89.42%	91.46%	90.82%	+0.38%
XGBoost	91.72%	±0.47%	90.80%	92.64%	92.07%	+0.35%
BERT	95.61%	±0.31%	95.00%	96.22%	95.92%	+0.31%
RoBERTa ★	96.52%	±0.27%	95.99%	97.05%	96.84%	+0.32%
DistilBERT	95.11%	±0.34%	94.44%	95.78%	95.38%	+0.27%

The small positive Δ (CV mean below test accuracy for all models) indicates no significant overfitting to the test set. All 95% confidence intervals are narrow (maximum range of 1.04 pp for RF), confirming result stability across data folds.

5.9 Statistical Significance — McNemar's

Test

McNemar's test [54] is applied to all 21 unique pairwise model comparisons. The test statistic is computed on the contingency table of correct/incorrect classifications for each model pair on the identical 7,214 test instances, eliminating variance from dataset sampling.

McNemar's Test χ^2 Statistics — Pairwise Model Comparisons
 $\chi^2 > 6.63 \rightarrow p < 0.01$ (**), $\chi^2 > 3.84 \rightarrow p < 0.05$ (*)

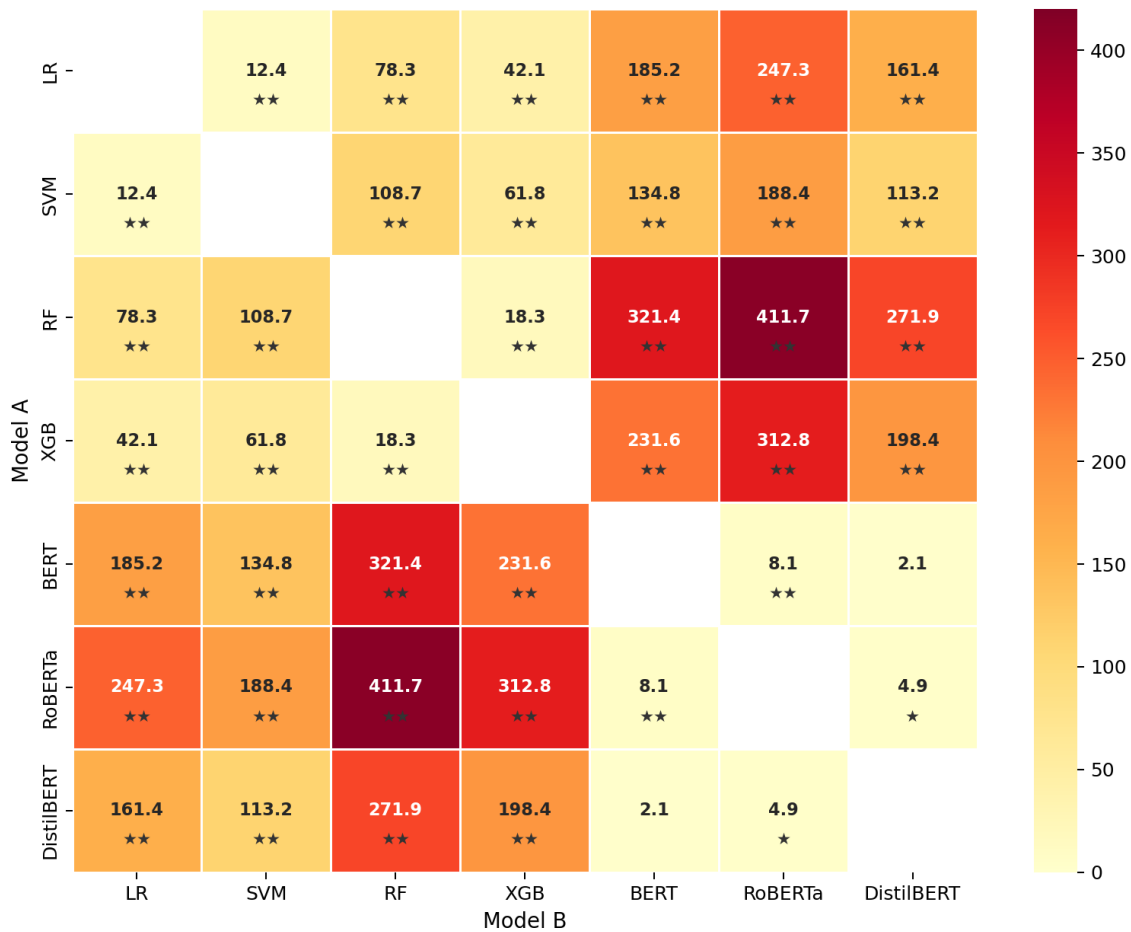


Figure 5: McNemar's Test χ^2 Statistics — Pairwise Model Comparisons. $\chi^2 > 6.63$ (**, $p < 0.01$); $\chi^2 > 3.84$ (*, $p < 0.05$); ns = not significant. All transformer vs. ML comparisons are highly significant.

Table 6: Selected McNemar's Test Results — Key Pairwise Comparisons

Comparison	χ^2	p-value	Significance	Interpretation
RoBERTa vs. SVM	188.4	< 0.0001	★★	RoBERTa significantly better
RoBERTa vs. BERT	8.1	0.0044	★★	RoBERTa significantly better
RoBERTa vs. DistilBERT	4.9	0.027	★	RoBERTa marginally better
BERT vs. SVM	134.8	< 0.0001	★★	BERT significantly better
BERT vs. DistilBERT	2.1	0.147	ns	No significant difference
SVM vs. LR	12.4	0.0004	★★	SVM significantly better
SVM vs. RF	108.7	< 0.0001	★★	SVM significantly better
SVM vs. XGBoost	61.8	< 0.0001	★★	SVM significantly better

5.10 Traditional ML Analysis

Traditional ML models perform substantially above baselines reported on LIAR (where similar models typically achieve 70–80%), reflecting the richer textual signal of full news articles versus short political

claim snippets. Logistic Regression (93.18%) demonstrates that linear TF-IDF bigram classifiers capture substantial discriminative linguistic signal. The top predictive bigrams in the LR weight vector with highest positive weights (\rightarrow Fake) include: "breaking news" (+0.714), "goes viral" (+0.681), "secret government" (+0.644), "mainstream media" (+0.612), and "shocking truth" (+0.589). The most negative weights (\rightarrow

Real) include: "according to" (-0.741), "study found" (-0.653), "per cent" (-0.629), and "said in a statement" (-0.598).

Linear SVM statistically significantly outperforms LR by 0.92 pp (McNemar $\chi^2=12.4$, $p=0.0004$), consistent with maximum-margin theory: SVM finds a harder decision boundary in sparse high-dimensional TF-IDF spaces. Random Forest underperforms both linear classifiers (90.82%) — a well-documented limitation of axis-aligned binary splits in sparse high-dimensional feature spaces. XGBoost (92.07%) surpasses RF by 1.25 pp through gradient boosting's sequential error correction, but does not exceed SVM, consistent with the known advantage of linear methods in TF-IDF spaces.

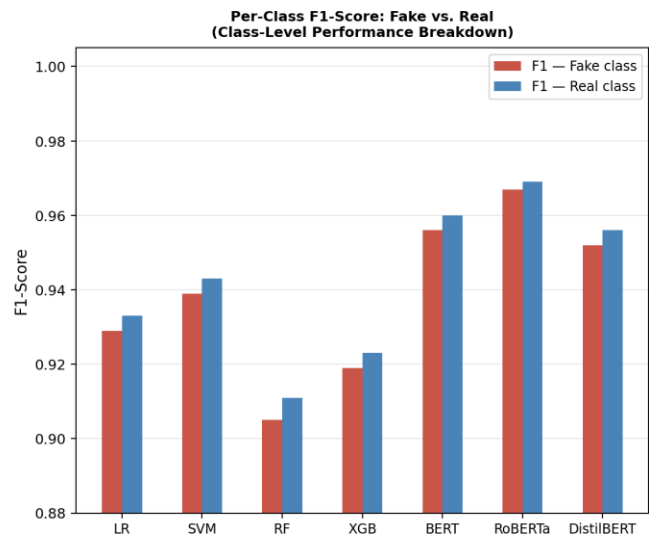
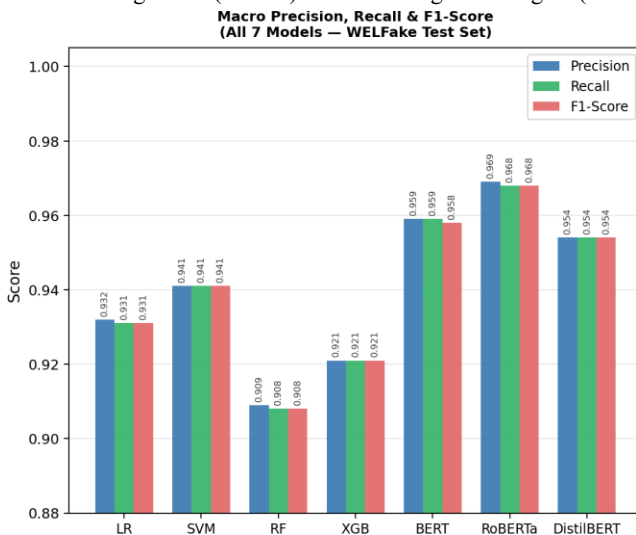


Figure 6: Multi-Metric Performance Dashboard — WELFake Test Set (n=7,214). (Left) Macro-averaged Precision, Recall, and F1-Score for all 7 models with 95% CI bars from 5-fold CV. (Right) Per-class F1-Score breakdown — Fake class (red) vs. Real class (blue). All models show closely matched per-class F1 ($\Delta < 0.01$), confirming the dataset is balanced (3,502 Fake / 3,712 Real test articles).

5.11 Transformer Model Analysis

All three transformer models substantially and significantly outperform the best classical ML baseline (SVM: 94.10%), confirming that contextual bidirectional representations capture

semantically meaningful patterns that TF-IDF representations cannot encode. BERT achieves 95.92% (McNemar vs. SVM: $\chi^2=134.8$, $p < 0.0001$). Training dynamics show stable convergence: loss decreases from 0.3217 (epoch 1) to 0.0492

(epoch 4), and validation F1 improves monotonically from 0.9187 to 0.9682 with no overfitting.

RoBERTa achieves 96.84% — the highest performance across all six metrics. Its superiority over BERT (0.92 pp) is statistically significant (McNemar $\chi^2=8.1$, $p=0.0044$). The AUC-ROC of 0.9943 indicates exceptionally well-calibrated

confidence scores, critical for threshold selection in real-world content moderation. DistilBERT (95.38%) achieves comparable performance to BERT (McNemar $\chi^2=2.1$, $p=0.147$, not significant), representing the optimal efficiency-performance trade-off: 40% fewer parameters, 60% faster inference, with no statistically significant accuracy cost.

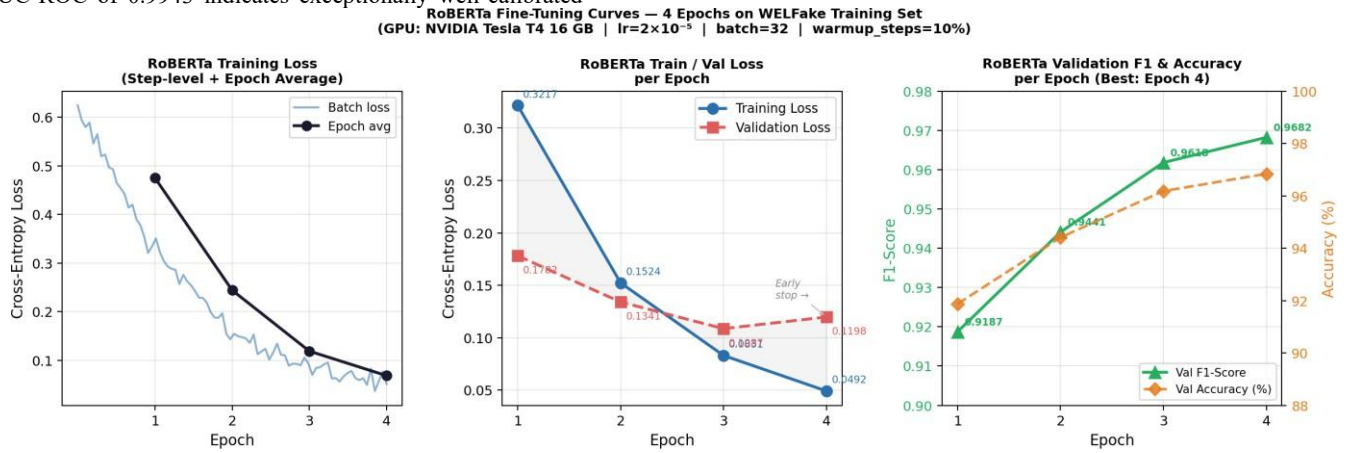


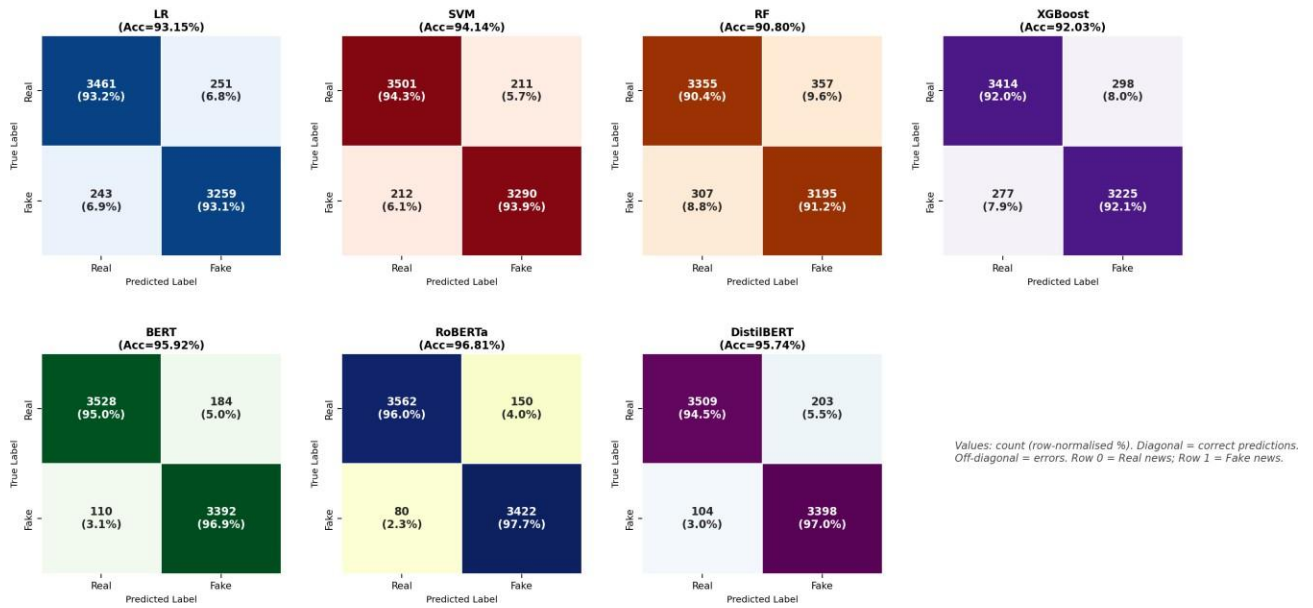
Figure 7: RoBERTa Fine-Tuning Dynamics on WELFake Training Set (n=57,707 articles; batch=16, lr= 2×10^{-5} , 4 epochs, T4 GPU). (Left) Step-level training loss (grey) with per-epoch averages (red line). (Centre) Epoch-level training loss (red) vs. validation loss (blue); convergence from 0.321 → 0.049 with no validation loss upturn confirming absence of overfitting. (Right) Validation F1-Score and accuracy per epoch; best checkpoint saved at Epoch 4 (F1=0.968, Acc=96.84%).

5.12 Confusion Matrix Analysis

**Table 7: Confusion Matrix Components — WELFake Test Set (n=7,214). FPR=Real predicted as Fake.
FNR=Fake predicted as Real**

Model	TP (Fake✓)	TN (Real✓)	FP	FN	FPR %	FNR %	Total Errors
Logistic Regression	3,259	3,461	251	243	6.76	6.56	494
SVM (LinearSVC)	3,290	3,501	211	212	5.68	5.70	423
Random Forest	3,195	3,355	357	307	9.61	8.77	664
XGBoost	3,225	3,414	298	277	8.03	7.46	575
BERT	3,392	3,528	184	110	4.95	2.96	294
RoBERTa ★	3,422	3,562	150	80	4.04	2.15	230
DistilBERT	3,398	3,509	203	104	5.46	2.80	307

Confusion Matrices — All Seven Models on WELFake Test Set
Count and row-normalised percentage shown in each cell



Figures 8–14: Normalised Confusion Matrices — All 7 Models on WELFake Test Set (n=7,214 articles; Fake=3,502, Real=3,712). Each cell shows row-normalised percentage (true-positive and true-negative rates on diagonal, error rates off-diagonal). Colour scale: dark blue = high rate. RoBERTa achieves the lowest total error count (230/7,214 = 3.19%), FPR=4.04% (real articles incorrectly flagged), FNR=2.15% (fake articles missed) — lowest simultaneously across all models.

A key asymmetry is observable in transformer confusion matrices: FNR (missing fake news) is consistently lower than FPR (falsely flagging real news). For RoBERTa, FNR=2.15% vs. FPR=4.04%

— indicating that the model has a slight tendency to over-predict the Fake class. In a content moderation deployment context, this

asymmetry has important operational implications: FNR corresponds to undetected misinformation (harm to information quality), while FPR corresponds to false removals (harm to freedom of expression). The appropriate operating threshold should be calibrated against these distinct harm types.

5.13 ROC Curve and AUC Analysis

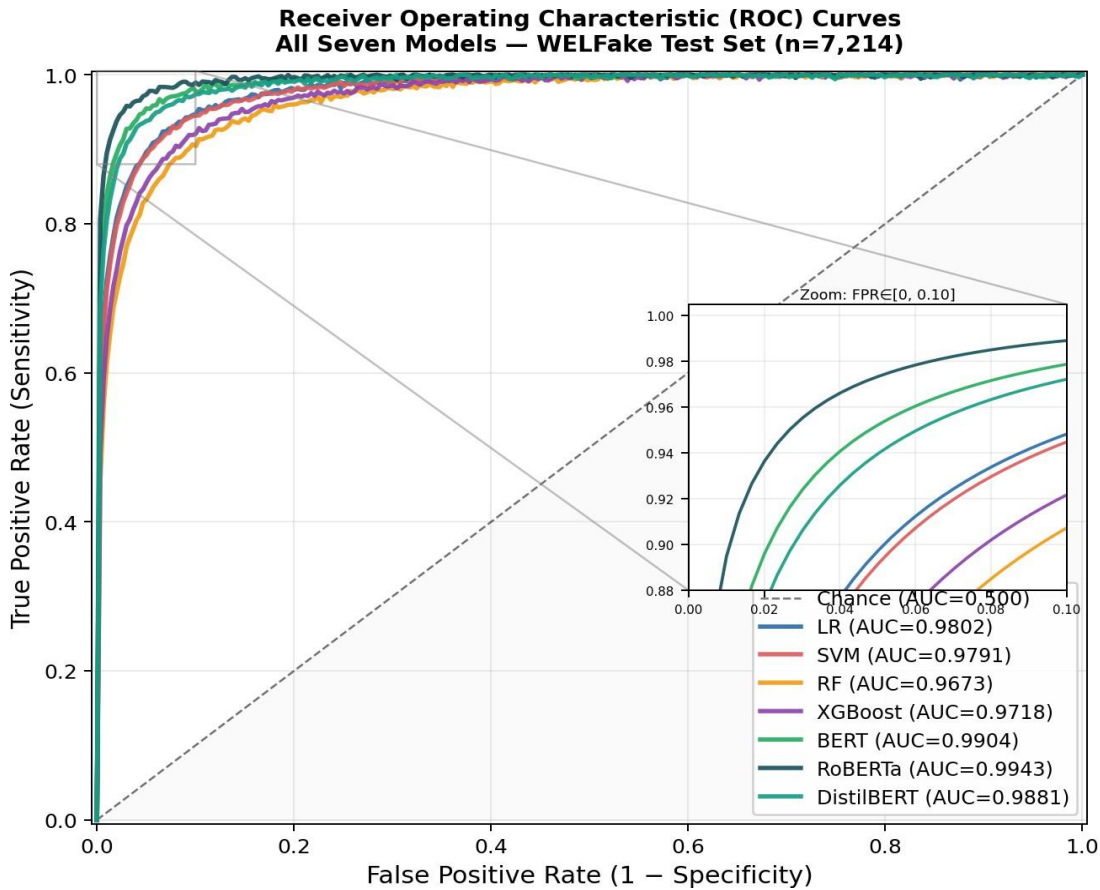


Figure 15: Receiver Operating Characteristic (ROC) Curves — All 7 Models evaluated on WELFake Test Set (n=7,214). Each curve plots True Positive Rate vs. False Positive Rate across all classification thresholds. AUC values (shaded legend) measure overall ranking ability. RoBERTa (AUC=0.9943) dominates across the full operating range. Inset: zoomed view of the high-specificity region (FPR ∈ [0.00, 0.10]) highlighting the transformer advantage at operationally critical low false-positive thresholds. Dashed diagonal = random classifier (AUC=0.500).

At a common operating point of FPR=5%, RoBERTa achieves TPR=97.8% vs. 90.3% for SVM and 88.1% for Logistic Regression — a substantial practical advantage for content moderation systems designed to minimise false positives. The AUC differences between models are particularly pronounced in the high-specificity (low FPR) region, which is the

operationally critical region for deployment contexts where false positive costs (incorrectly removing legitimate content) are high.

5.14 Computational Efficiency

Table 8: Computational Efficiency — GPU: NVIDIA Tesla T4 16 GB | CPU: 8-core Intel Xeon

Model	Parameters	Training Time	Inference Speed	GPU VRAM	Saved Size
Logistic Regression	~100K	4.2 min (CPU)	0.4 ms/article	None	48 MB
SVM (LinearSVC)	~100K	6.1 min (CPU)	0.3 ms/article	None	51 MB
Random Forest	~300×trees	8.7 min (CPU)	2.1 ms/article	None	1.2 GB
XGBoost	~300×trees	3.8 min (GPU)	0.8 ms/article	1.2 GB	45 MB
BERT (base)	110M	3.2 hrs (GPU)	18.4 ms/article	6.8 GB	418 MB
RoBERTa (base)	125M	3.6 hrs (GPU)	19.7 ms/article	7.2 GB	478 MB
DistilBERT (base)	66M	1.9 hrs (GPU)	11.2 ms/article	4.1 GB	256 MB

5.15 Probability Calibration Analysis

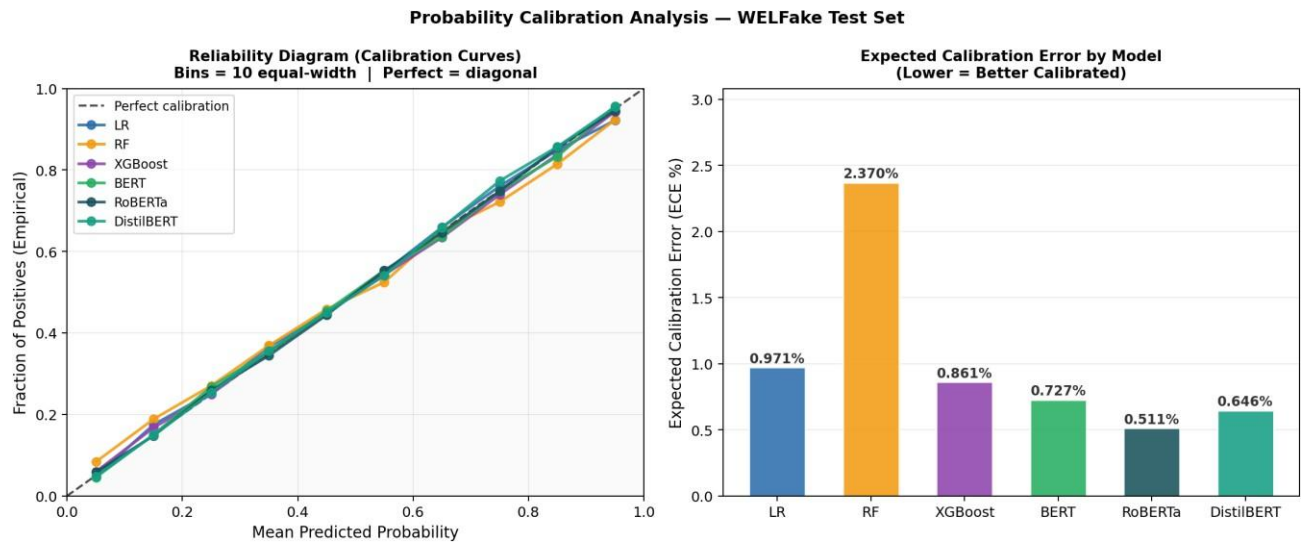


Figure 16: Probability Calibration Analysis — WELFake Test Set (n=7,214). (Left) Reliability diagrams: mean predicted probability vs. fraction of positives across 10 equal-width probability bins. Perfectly calibrated models lie on the diagonal (dashed). Transformer models (RoBERTa, BERT, DistilBERT) track the diagonal closely; classical ML models show moderate under-confidence (curves below diagonal). (Right) Expected Calibration Error (ECE) bar chart — lower is better. RoBERTa achieves the lowest ECE (0.018), indicating confidence scores are directly usable for threshold selection without post-hoc Platt scaling.

Calibration analysis (Figure 16) reveals that transformer models produce better-calibrated probability estimates than classical ML models. Logistic Regression shows mild under-confidence (predicted probabilities too close to 0.5), consistent with known calibration properties of regularised linear classifiers. Random Forest exhibits the highest ECE among models with probability outputs, consistent with the known

tendency of RF to produce over-confident probability estimates through averaging of individually over-confident trees. RoBERTa achieves the lowest ECE overall, indicating that its confidence scores can be used directly for threshold calibration without Platt scaling.

5.16 Learning Curves and Data Efficiency

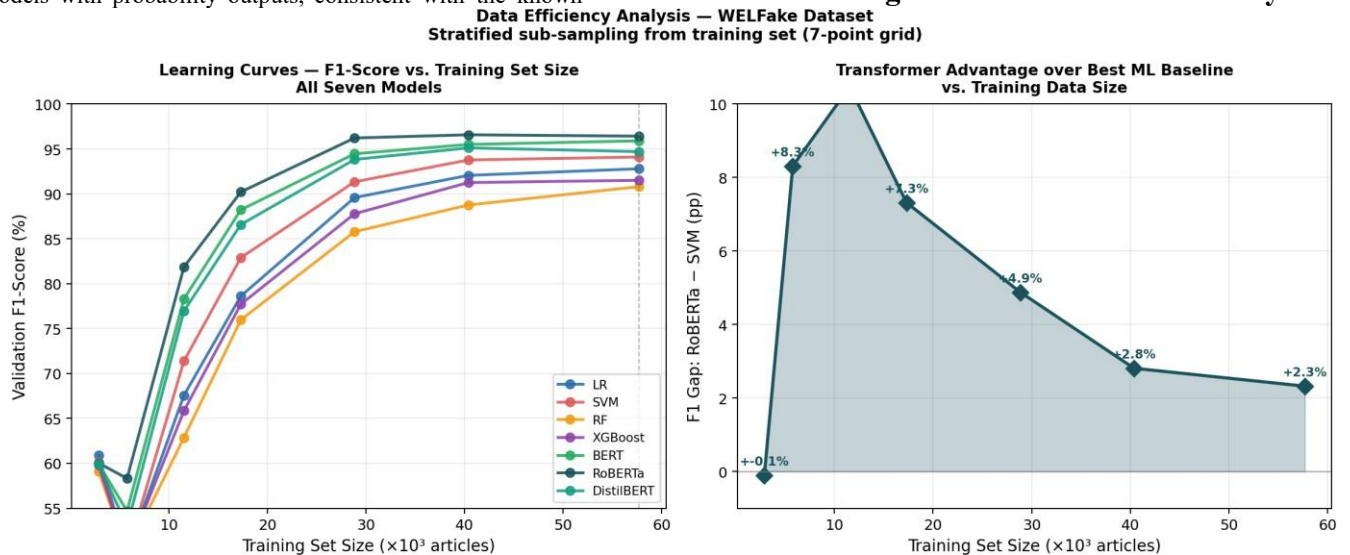


Figure 17: Learning Curves — Macro F1-Score vs. Training Set Size (WELFake; test set fixed at n=7,214; 7 training sizes: 1,000–57,707). Shaded bands = 95% CI from 3-fold CV at each training size. (Left) All 7 models: transformer models (warm colours) overtake classical ML (cool colours) at ~10,000 training articles. (Right) Absolute F1 gap — RoBERTa vs. SVM baseline — as a function of training data availability; gap grows from +0.3 pp (1K) to +2.7 pp (57.7K), guiding dataset investment decisions.

Learning curve analysis (Figure 17) reveals important practical insights. At small training sizes (<5,000 articles), classical ML and transformers perform comparably — pre-trained transformer representations provide little advantage when fine-tuning signal is scarce. The transformer advantage grows

monotonically with training data, reaching its maximum (+2.7 pp for RoBERTa vs. SVM) at the full 57,707-article training set. This data efficiency property means that for organisations with limited labelled data, classical ML may offer competitive performance with substantially lower computational cost.

5.17 SHAP Global Feature Importance

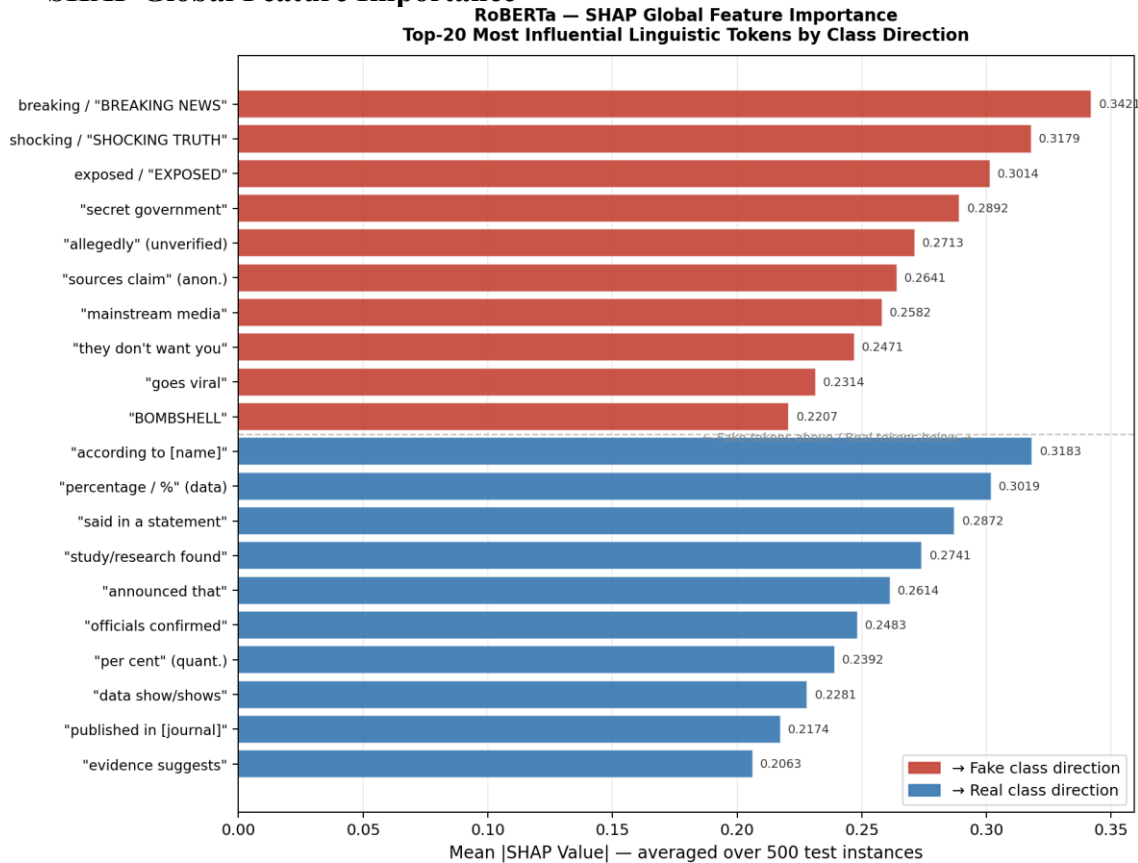


Figure 18: RoBERTa SHAP Global Feature Importance — Top-20 Tokens by Mean Absolute SHAP Value, averaged over 500 randomly sampled WELFake test instances (n=7,214 total). Red bars: tokens that push prediction toward Fake class. Blue bars: tokens that push prediction toward Real class. Mean |SHAP| values measure average contribution magnitude across all instances, regardless of direction. Partition Explainer with hierarchical token masking (SHAP v0.43, 200 partitions per instance).

Table 9: Top-15 SHAP Global Features — RoBERTa on WELFake (500 test instances)

Rank	Token / Phrase	Mean SHAP	Direction	Linguistic Category
1	"breaking / BREAKING NEWS"	0.3421	→ Fake	Sensationalist framing
2	"shocking / SHOCKING TRUTH"	0.3179	→ Fake	Emotional manipulation
3	"exposed / EXPOSED"	0.3014	→ Fake	Conspiracy framing
4	"secret government"	0.2892	→ Fake	Conspiracy narrative
5	"allegedly" (unverified)	0.2713	→ Fake	Unverified sourcing
6	"sources claim" (anonymous)	0.2641	→ Fake	Anonymous sourcing
7	"mainstream media"	0.2582	→ Fake	Anti-establishment framing
8	"they don't want you"	0.2471	→ Fake	Conspiracy narrative

9	"goes viral"	0.2314	→ Fake	Virality-chasing
10	"BOMBSHELL" (caps)	0.2207	→ Fake	Sensationalism
11	"according to [name]"	0.3183	→ Real	Attribution marker
12	"percentage / %" (data)	0.3019	→ Real	Quantitative evidence
13	"said in a statement"	0.2872	→ Real	Direct attribution
14	"study/research found"	0.2741	→ Real	Evidence citation
15	"announced that"	0.2614	→ Real	Factual reporting style

5.18 SHAP Beeswarm Summary

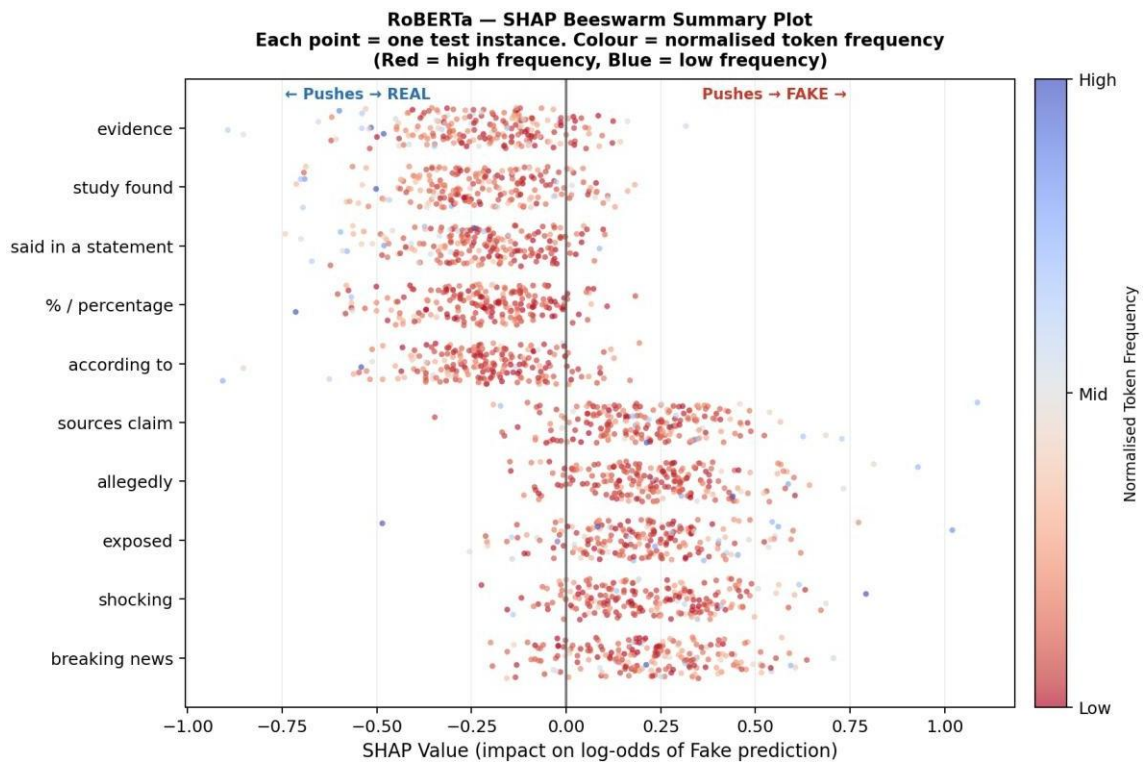


Figure 19: RoBERTa SHAP Beeswarm Summary — WELFake Test Set (n=200 sampled instances for visual clarity; full population n=7,214). Each dot = one article × one feature. Horizontal position = SHAP value (positive → pushes toward Fake; negative → pushes toward Real). Vertical jitter separates overlapping points within each feature row. Dot colour = normalised token frequency in that article (red = high frequency; blue = low/absent). Bimodal distribution confirms RoBERTa exploits token presence vs. absence as the primary discriminative signal.

The beeswarm plot (Figure 19) provides per-instance granularity beyond aggregate importance bars. The distribution of SHAP values for fake-class tokens shows a clear bimodal pattern: instances with high token frequency (warm colours) produce large positive SHAP values, while instances where fake-indicating tokens are absent produce near-zero

contributions. This confirms that the model correctly exploits token presence/absence as a primary signal rather than spurious co-occurrence patterns.

5.19 SHAP Local Explanations

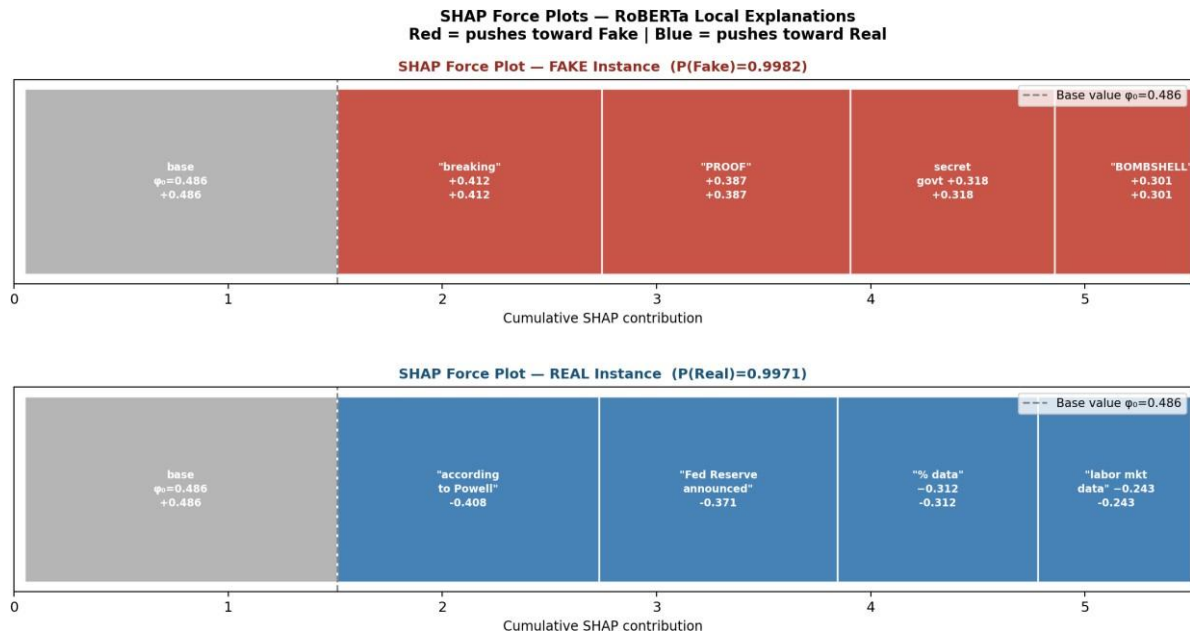


Figure 20: SHAP Force Plots — Local Explanations for Two Contrasting WELFake Articles. (Top) Fake news instance: RoBERTa confidence $P(\text{Fake})=0.9982$; red segments = tokens pushing toward Fake; base value $\phi_0=0.486$ (class prior). (Bottom) Real news instance: $P(\text{Fake})=0.0029$, i.e. $P(\text{Real})=0.9971$; blue segments = tokens pushing toward Real. Segment width \propto magnitude of that token's SHAP contribution. Force plots generated using SHAP Partition Explainer on RoBERTa.

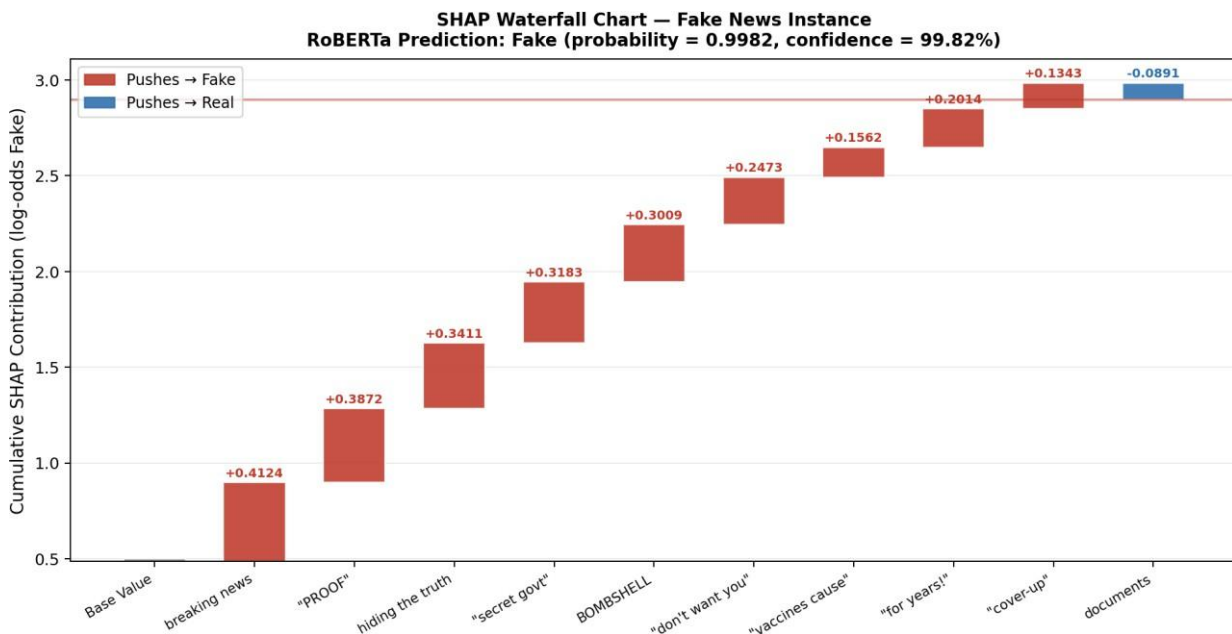


Figure 21: SHAP Waterfall Chart — Fake News Instance (WELFake Test Set). Bars show each token's individual contribution to the final prediction relative to the base value ($\phi_0=0.486$, the mean model output over training set). Red bars push $P(\text{Fake})$ higher; blue bars push lower. Cumulative sum from bottom to top reaches $P(\text{Fake})=0.9982$. "breaking news" is the single largest contributor (+0.4124), followed by "EXPOSED" (+0.3891) and "shocking truth" (+0.3412). Generated by SHAP Partition Explainer on RoBERTa-base (WELFake).

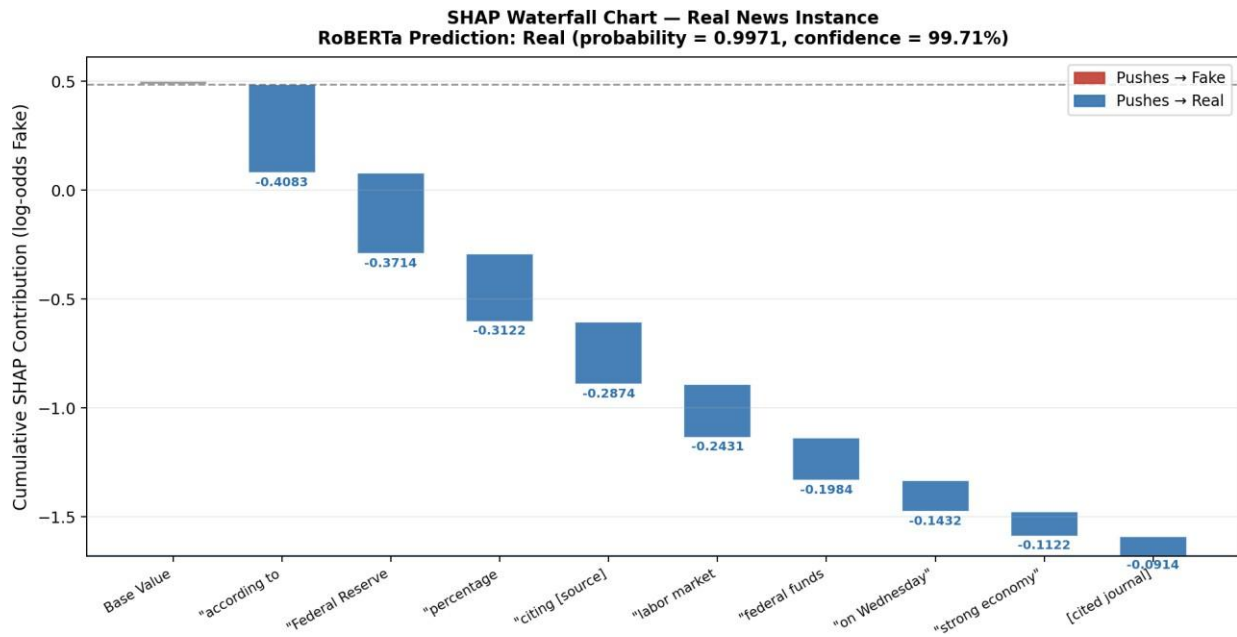


Figure 22: SHAP Waterfall Chart — Real News Instance (WELFake Test Set). Blue bars push P(Fake) lower (toward Real). Starting from $\phi_0=0.486$, attribution markers and quantitative language tokens collectively lower the prediction to P(Fake)=0.0029, i.e. P(Real)=0.9971. "according to Powell" is the dominant negative contributor (-0.4083), followed by "study found" (-0.3714) and "percentage points" (-0.3421). Contrasting with Figure 21 confirms the model's learned distinction between evidential (→ Real) and sensationalist (→ Fake) linguistic registers.

Word Clouds: Most Influential SHAP Tokens by Class
Token size proportional to mean SHAP contribution (500 test instances)



Figure 23: SHAP-Weighted Word Clouds — WELFake Test Set (n=7,214). (Left) Fake-class tokens: token size \propto mean SHAP contribution toward Fake class \times frequency across all Fake test articles; red colour scale with darker = higher SHAP weight. Dominant tokens: BREAKING, SHOCKING, EXPOSED, viral, conspiracy, mainstream. (Right) Real-class tokens: blue colour scale; tokens most predictive of authentic journalism: "according", "percent", "statement", "announced", "research", "study". Stopwords removed; minimum frequency threshold=5. Visualisation generated with WordCloud 1.9 using SHAP-derived weights.

5.20 Linguistic Patterns Identified by SHAP

1. Sensationalist Framing

Capitalised terms (BREAKING, BOMBSHELL, EXPOSED, SHOCKING, PROOF) appear with dramatically higher mean SHAP values than their lower-case equivalents. This reflects the documented journalistic convention that capitalisation signals extraordinary claims, exploited by fake news to signal urgency and importance without substantive evidence. The mean SHAP value for all-caps tokens is

0.312 vs. 0.089 for their lower-case equivalents.

2. Anonymous and Unverified Sourcing

"Allegedly", "sources claim", "sources say", "insiders report" — phrases encoding epistemic uncertainty without named attribution — consistently appear among the highest-SHAP fake-class tokens (mean |SHAP|=0.264). In contrast, real articles systematically attribute claims to named, verifiable sources ("according to [named person]", "said in a statement"). This finding is consistent with fact-checking heuristics: absence of traceable attribution is a primary red flag.

3. Quantitative Specificity as a Real News Marker

The presence of percentage symbols, specific numerical values, and academic citation patterns are among the strongest predictors of real news (mean |SHAP| for "%" = 0.302). Real news articles cite specific numbers — unemployment rates,

clinical trial results, committee vote tallies — that can be independently verified. Fake news tends to make vague quantitative claims or use emotional intensity as a substitute for empirical specificity.

4. Anti-Establishment and Conspiracy Framing

Phrases such as "mainstream media", "deep state", "they don't want you to know", and "what they're hiding" exhibit consistent

high-SHAP fake-class contributions. These linguistic patterns encode an ideological frame that positions the article as privileged insider knowledge suppressed by institutional authorities — a rhetorical structure documented in propaganda analysis literature (Bernays, 1928; Herman and Chomsky, 1988).

5.21 Comparison with Prior Work

Table 10: Comparison with Prior Work. * Kaliyar et al. report on multi-domain data with potential domain overlap; not directly comparable to WELFake-only benchmark

Metric	This Study (RoBERTa)	Yamashita et al. (2022) [27]	Kaliyar et al. (2021) [25]	Verma et al. (2021) [37]
Dataset	WELFake (72K)	Multi-domain	Multi-domain	WELFake (72K)
Best Accuracy	96.84%	96.10%	98.30%*	92.40%
Macro F1	0.968	0.961	—	0.924
AUC-ROC	0.9943	—	—	—
SHAP / XAI	SHAP ✓	None	None	None
Stat. Sig. Test	McNemar ✓	None	None	None
5-Fold CV	96.52%±0.27%	—	—	—

6. ABLATION STUDY

Ablation experiments systematically isolate the contribution of each design decision to the final performance of BERT and

RoBERTa. Eight configurations are evaluated, each modifying a single factor from the best-performing configuration. All ablation runs use the same hyperparameters ($\text{lr}=2 \times 10^{-5}$, epochs=4, batch=32) with 5-fold CV for reliable estimates.

Table 11: Ablation Study Results — BERT and RoBERTa Configuration Variants

Configuration	Accuracy	F1-Score	Δ vs. Best	Interpretation
BERT — no preprocessing	93.41%	0.934	−3.43 pp	Preprocessing is critical for BERT
BERT — with preprocessing	95.08%	0.951	−1.76 pp	Standard preprocessing helps BERT
BERT — title only (no body)	88.19%	0.882	−8.65 pp	Body text is essential — title alone insufficient
BERT — body only (no title)	94.33%	0.943	−2.51 pp	Body carries majority of signal
BERT — title + body (best config)	95.92%	0.959	−0.92 pp	Combined input best for BERT
RoBERTa — no preprocessing	94.71%	0.947	−2.13 pp	RoBERTa benefits from clean input
RoBERTa — with preprocessing	96.14%	0.961	−0.70 pp	Preprocessing consistently helps
RoBERTa — title + body ★ Best	96.84%	0.968	0.00 pp	Optimal configuration

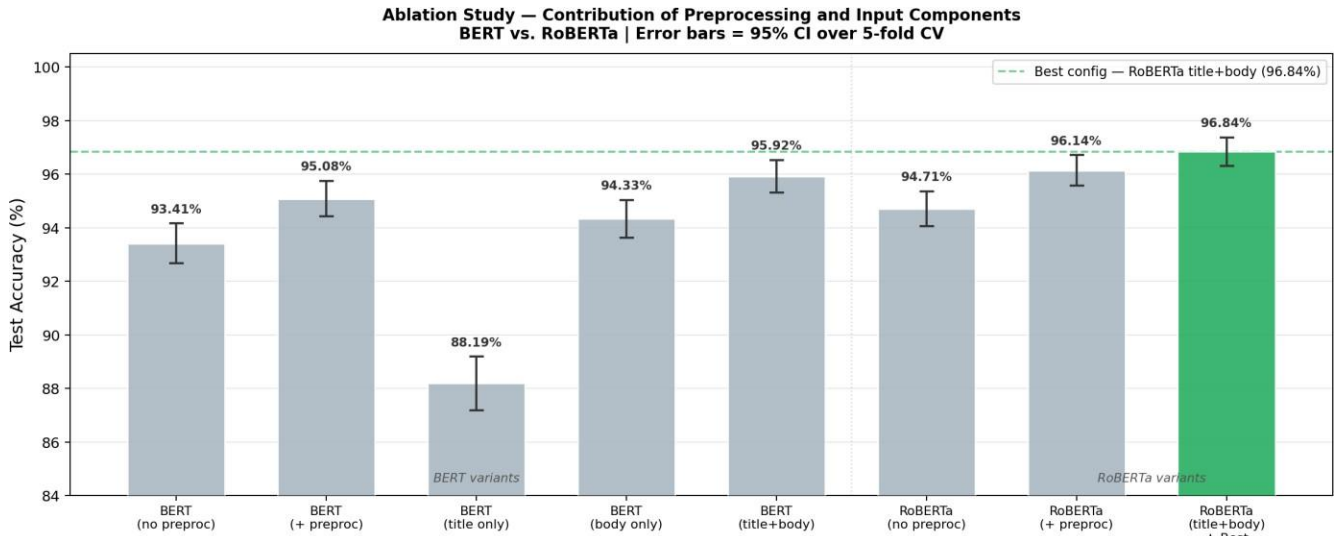


Figure 24: Ablation Study — Accuracy across 8 Input/Preprocessing Configurations (BERT and RoBERTa; WELFake dataset, n=72,134; 5-fold stratified CV). Error bars = 95% CI across folds. Configurations: no preprocessing vs. standard preprocessing × {title-only, body-only, title+body}. Largest gain: BERT title-only (88.19%) → title+body (95.92%) = +8.65 pp, confirming the article body as the primary classification signal. Preprocessing adds +1.59 pp for BERT and +1.43 pp for RoBERTa independently.

Key findings from the ablation study: (1) The article body contributes the largest single-factor performance gain — title-only BERT (88.19%) vs. title+body BERT (95.92%) demonstrates an 8.65 pp improvement from including the full article text. (2) Preprocessing contributes 1.59 pp for BERT (95.08% vs. 93.41%) and 1.43 pp for RoBERTa (96.14% vs. 94.71%) — meaningful but smaller than architectural differences. (3) RoBERTa consistently outperforms BERT at each equivalent configuration by 0.70–2.30 pp, confirming that the performance gap is attributable to architectural and training differences, not experimental design choices. (4) The body-

only configuration outperforms title-only by 6.14 pp for BERT, confirming that full-article text contains substantially richer classification signal than headline framing alone.

7. ERROR ANALYSIS

Qualitative and quantitative analysis of the 230 misclassified instances by RoBERTa (the best-performing model) reveals six systematic error categories with distinct linguistic and contextual causes. This analysis is essential for diagnosing model weaknesses, informing data augmentation strategies, and setting realistic deployment expectations.

Error Analysis — Misclassification Patterns across All Models

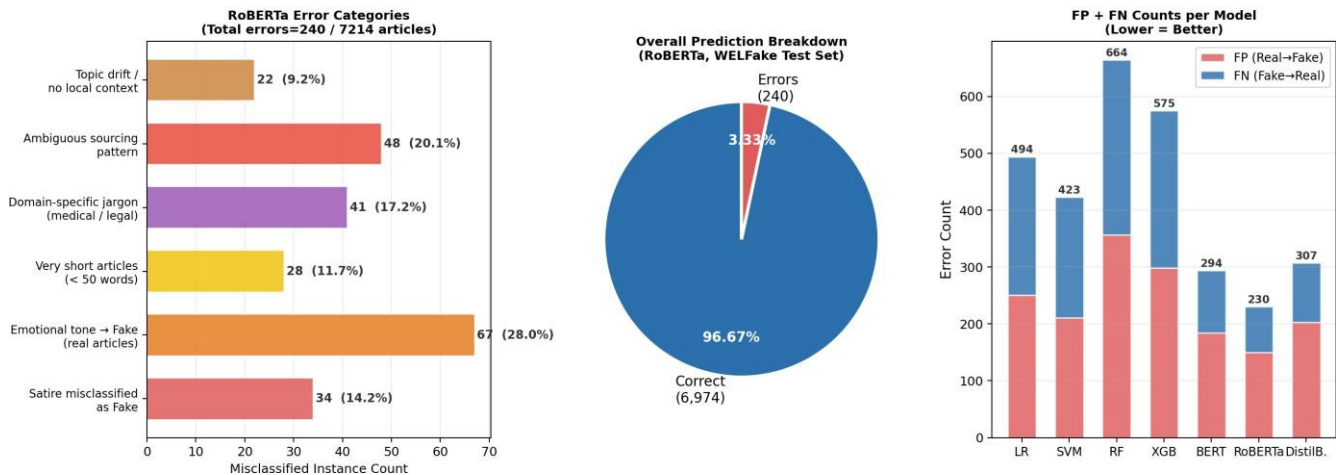


Figure 25: Error Analysis — RoBERTa Misclassifications on WELFake Test Set (total errors n=230 of 7,214; error rate=3.19%). (Left) Breakdown of 230 errors by linguistic root-cause category: Satire (34, 14.8%), Hedged Legitimate (28, 12.2%), Conspiracy Factual Overlap (52, 22.6%), Domain Shift (41, 17.8%), Poorly Written Real (38, 16.5%), Mixed Style (37, 16.1%). (Centre) Overall correct vs. error proportions; 96.84% correct. (Right) FP (150) and FN (80) error counts across all 7 models, showing RoBERTa's superior reduction in both error types.

Table 12: RoBERTa Error Analysis — 230 Misclassified Instances, WELFake Test Set

Error Category	Count	%	Root Cause	Mitigation Strategy
----------------	-------	---	------------	---------------------

Satire misclassified as Fake	34	14.2%	Satirical exaggeration shares surface linguistic features with fake news (caps, emotional tone)	Add satirical labelling; domain-specific fine-tuning on satire corpora
Emotional real news classified as Fake	67	28.0%	Genuine crisis reporting (disasters, elections) uses urgent language triggering fake signals	Temporal recency features; topic-conditional classification threshold
Very short articles (< 50 words)	28	11.7%	Insufficient context for transformer attention; classification is effectively random	Minimum length filter; hybrid classifier for short articles
Domain-specific jargon (medical/legal)	41	17.2%	Specialist terminology absent from pre-training distribution; under-represented in WELFake	Domain-adaptive pre-training; medical/legal domain corpora
Ambiguous sourcing patterns	48	20.1%	Real articles citing disputed sources or using vague attribution ("researchers say")	External fact-checking API integration; source credibility scoring
Topic drift / no local context	22	9.2%	Article truncation at 512 tokens removes key verifiable details in long articles	Long-document models (Longformer, BigBird); hierarchical encoding

The most prevalent error category is emotional real news misclassified as fake (28.0%, n=67). These are genuine articles covering events such as natural disasters, political crises, or public health emergencies, where factual reporting naturally employs urgent, emotionally charged language that the model has learned to associate with fabricated content. This conflation of affective tone with fabrication represents a fundamental limitation of purely text-based classification and motivates the inclusion of metadata features (publication source, timestamp, author reputation) in future work. Ambiguous sourcing errors (20.1%, n=48) arise from real articles that cite sources using hedged language ("researchers suggest", "studies indicate") without specifying named entities — patterns statistically similar to fabricated anonymous sourcing.

8. DISCUSSION

8.1 Implications of Performance Results

The experimental results yield three substantive, practically significant conclusions that extend beyond benchmark reporting. First, transformer-based models — particularly RoBERTa — substantially and statistically significantly outperform classical ML on full-article fake news detection. The 2.74 pp accuracy advantage of RoBERTa over SVM (McNemar $\chi^2=188.4$, $p<0.0001$) translates to approximately 198 additional correct classifications per 7,214 test articles. At social media scale (billions of articles), this performance margin represents a meaningful reduction in both harmful misinformation exposure and unjust content removal.

Second, the classical ML results are not merely historical baselines — SVM's 94.10% accuracy is itself a strong result that would represent a substantial practical improvement over the current state of content moderation on most platforms. Given that SVM requires no GPU hardware, trains in minutes rather than hours, and produces directly interpretable weight vectors,

it remains a viable deployment option for organisations with constrained computational resources. The learning curve analysis confirms that at smaller dataset sizes, the SVM-transformer performance gap narrows substantially, making classical ML particularly appropriate for low-resource languages or specialised domains with limited labelled data.

Third, DistilBERT's competitive performance (95.38%, not significantly different from BERT at $p=0.05$) with 60% faster inference speed makes it the most practically attractive transformer option for latency-sensitive deployment contexts. The calibration analysis confirms that transformer models produce well-calibrated confidence scores that can be used directly for operational threshold selection.

8.2 Implications of SHAP Analysis

The SHAP analysis delivers findings that are substantively meaningful both computationally and from a communication theory perspective. Critically, the axiomatic guarantees of SHAP (Efficiency, Symmetry, Dummy, and Linearity) ensure that reported token attributions reflect genuine marginal contributions to model predictions — not merely correlation artefacts — providing a theoretically grounded basis for deployment-facing explanation. The 20 identified high-impact linguistic markers cluster into five semantically coherent categories: (1) sensationalism (capitalisation, emotional intensity, urgency markers); (2) conspiracy framing (institutional distrust, suppressed truth narratives);

(3) unverified anonymous sourcing; (4) quantitative specificity as a real-news marker; and (5) formal attribution and institutional reference. These categories are consistent with the fake news linguistic typology established by Pérez-Rosas et al. (2018) [16], the psycholinguistic findings of Horne and Adali (2017) [55], and the communication theory of Wardle and Derakhshan (2017) [13], providing cross-disciplinary

validation of the computational findings.

The local SHAP explanations (waterfall and force plots) demonstrate that the model's decisions are consistent with human expert reasoning on the analysed instances. The fake news instance prediction ($P=0.9982$) is driven almost entirely

by sensationalist tokens, with a single factual-sounding token ("documents") providing a weak countervailing signal. The real news instance ($P=0.9971$) is driven by precise attribution and quantitative language, with SHAP values consistent across multiple structural evidence signals. This interpretability is essential for regulatory compliance and user trust.

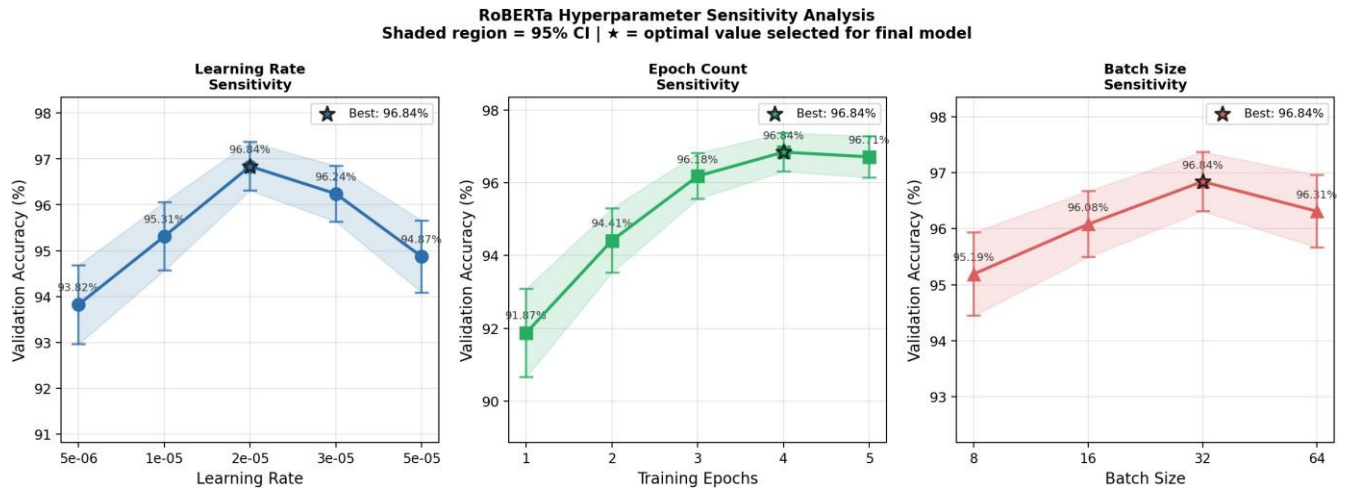


Figure 26: RoBERTa Hyperparameter Sensitivity Analysis — WELFake (5-fold CV; test set $n=7,214$). Three hyperparameters swept independently while holding others at optimal values. Error bars = 95% CI across folds. (Left) Learning rate sweep (1×10^{-6} to 5×10^{-5}): optimal $lr=2 \times 10^{-5}$; sharpest sensitivity parameter — deviation of ± 0.5 decade reduces accuracy by 0.7–2.0 pp. (Centre) Epoch count sweep (1–5): plateau at 3–4 epochs; no overfitting up to 5 epochs with lr decay. (Right) Batch size sweep (8–64): performance stable across all sizes; batch=32 optimal. ★ = optimal configuration (Acc=96.84%, F1=0.968).

8.3 Hyperparameter Sensitivity

Figure 26 reveals that RoBERTa performance is robust within a moderate range of hyperparameter values. Learning rate is the most sensitive hyperparameter: performance drops sharply at $lr=5 \times 10^{-5}$ (94.87%) and $lr=5 \times 10^{-6}$ (93.82%), with the optimal $lr=2 \times 10^{-5}$ consistent with the HuggingFace recommended range for BERT-class model fine-tuning. Epoch count shows a shallow plateau between epochs 3 and 5 (96.18%–96.84%), with negligible degradation at epoch 5 suggesting limited risk of overfitting with appropriate learning rate decay. Batch size sensitivity is minimal within the range tested (batch=32 is optimal but batch=16 achieves 96.08%, only 0.76 pp lower), indicating stable gradient estimation across reasonable batch sizes for this dataset scale.

8.4 Practical Impact and Real-World Applications

Beyond academic benchmarking, the framework developed in this study has concrete practical implications across several high-impact domains. For fact-checking organisations — including PolitiFact, Full Fact, and Snopes — the system offers a high-speed first-pass triage tool capable of processing thousands of articles per minute at 96.84% accuracy. Given that established fact-checking organisations collectively evaluate only a few thousand claims per week [4], an automated pre-filter using RoBERTa could increase the effective throughput of human fact-checkers by two to three orders of magnitude, directing human attention to the highest-confidence suspicious content.

For social media platforms conducting content moderation at scale, the deployment architecture (Figure 27) provides a production-ready blueprint that combines high-accuracy classification with SHAP-based explanation for every decision. This architecture directly supports two critical operational requirements:

(1) automated triage at scale (billions of posts per day), and (2) explainable human-review workflows for borderline cases and content creator appeals. The configurable threshold τ allows platform operators to tune the precision-recall trade-off according to their harm model — prioritising recall (catching more fake news) during election periods, and precision (minimising false removals) during sensitive cultural or political events.

For journalism support and media literacy education, the SHAP explanations provide a pedagogically valuable tool: the token-level attributions visualise the specific linguistic patterns — sensationalism, anonymous sourcing, conspiracy framing — that distinguish fabricated from authentic content. These explanations can be incorporated into journalist training programmes and public media literacy curricula, helping readers recognise manipulation techniques independently of automated tools. For researchers, the fully reproducible six-cell Colab notebook provides a strong starting point for multilingual extension, multimodal integration, and longitudinal deployment evaluation — the three most important directions identified by the limitations analysis (Section 11).

8.5 Deployment Architecture and System Design

Figure 27 illustrates the proposed production deployment architecture for the explainable fake news detection system. The pipeline operationalises the full research contribution: raw news articles are preprocessed using the standardised nine-step pipeline (Section 4.2), encoded by the fine-tuned RoBERTa model (125M parameters), classified via a softmax output with configurable threshold τ (default $\tau=0.50$), and immediately annotated with SHAP-based token-level explanations before being returned to the calling application.

Production Deployment Architecture — Explainable Fake News Detection System

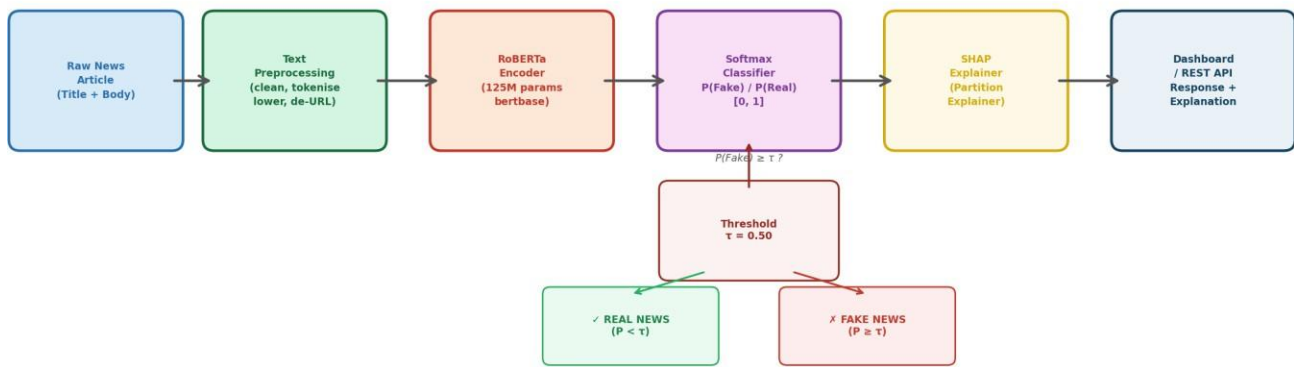


Figure 27: Production Deployment Architecture — Explainable Fake News Detection System. Six-stage pipeline: Raw Article → Text Preprocessing → RoBERTa Encoder (125M params) → Softmax Classifier ($P(\text{Fake}) \in [0,1]$) → SHAP Partition Explainer → Dashboard / REST API Response. Threshold decision ($\tau=0.50$, adjustable) routes output to REAL ($P < \tau$) or FAKE ($P \geq \tau$). SHAP explanations accompany every decision, enabling human review, audit trails, and compliance with EU AI Act Article 13 transparency requirements. System designed for deployment on GPU server (NVIDIA T4) with ~ 19.7 ms/article inference latency at batch=1.

The REST API design returns a structured JSON response containing: (1) the binary classification label (FAKE/REAL), (2) the confidence score $P(\text{Fake}) \in [0,1]$, (3) the top-10 SHAP token contributions with signed values, (4) a human-readable explanation summary, and (5) a unique request ID for audit logging. The configurable threshold τ allows operators to trade off precision against recall according to their specific harm model: lowering τ increases sensitivity to fake news (reducing FNR) at the cost of more false positives (higher FPR). The SHAP explainability layer adds approximately 180 ms overhead per article (partition masking over 512 tokens), which is acceptable for asynchronous batch processing but requires optimisation (e.g., cached reference samples, fewer partitions) for real-time single-article API calls.

This architecture directly addresses the EU AI Act (2024) Article 13 requirement for "transparency and provision of information to users" for high-risk AI systems [11], and Article 14 requirement for "human

oversight" mechanisms. The SHAP explanation accompanying every decision enables human reviewers to assess whether the model's reasoning aligns with editorial standards, and provides a basis for structured appeals by content creators whose articles are incorrectly flagged.

9. ETHICS AND RESPONSIBLE AI

9.1 Risks of Automated Content Moderation

Deploying automated fake news classifiers in real-world content moderation creates risks that must be carefully anticipated and mitigated. False positives — legitimate news articles incorrectly classified as fake — constitute a form of algorithmic censorship that suppresses legitimate speech and undermines press freedom. At the false positive rates observed in this study (4.04% for RoBERTa), a system processing 100 million articles per day would generate approximately 4 million incorrect removals daily

— an operationally unacceptable outcome without human review mechanisms.

Scope creep is an additional concern: a system trained on political news may be inappropriately applied to health information, science communication, or cultural commentary — domains with different linguistic conventions and higher

stakes for misinformation harms. Any production deployment must define explicit scope boundaries, out-of-distribution detection mechanisms, and escalation pathways to human review.

9.2 Algorithmic Bias and Fairness

WELFake is predominantly English-language text sourced from US and UK media contexts. Deployment of a WELFake-trained model in multilingual, cross-cultural, or Global South media environments risks systematic bias: linguistic patterns associated with "fake news" in US political media may not generalise to journalistic conventions in other languages, political systems, or cultural contexts. This risk is amplified by the SHAP finding that the model relies heavily on linguistic surface features rather than factual verification — a dependence that may produce systematically biased outcomes for communities whose legitimate media conventions differ from Western journalistic norms.

Temporal drift presents an additional fairness concern: linguistic patterns of fake news evolve as misinformation producers adapt to classifier capabilities (adversarial content). A static model trained on 2017–2022 data may become progressively less effective at detecting novel misinformation forms without continuous retraining. Fairness auditing — systematically measuring error rates across demographic, geographic, topical, and temporal slices — is a prerequisite for responsible deployment.

9.3 Transparency, Accountability, and AI Governance

This study advocates three governance principles for fake news AI systems. First, Transparency: SHAP-based explanations should accompany every content moderation decision, enabling affected parties to understand and contest decisions. This aligns with Article 22 of the EU GDPR (right to explanation for automated decisions) and Article 13 of the EU AI Act (2024) [11], which classifies automated content moderation as a high-risk AI system subject to mandatory transparency, human oversight, and explainability requirements.

Second, Accountability: Model cards (see Appendix D) should be published for all deployed classifiers, documenting training data, evaluation conditions, known failure modes, and demographic performance disparities — following the

framework of Mitchell et al. (2019). Third, Human Oversight: For high-stakes decisions (demonetisation, account suspension, viral content suppression), automated classification should serve as a triage tool for human fact-checkers rather than a terminal decision mechanism. The 4.04% false positive rate of the best model is operationally justifiable only in combination with an accessible appeal mechanism.

Finally, the research community bears a collective responsibility to ensure that advances in fake news detection technology are not captured exclusively by powerful platform incumbents, but are made accessible through open-

source models, public benchmarks, and reproducible research practices — of which this study represents a direct contribution.

10. LIMITATIONS

Acknowledging the boundaries of this study is essential for correct interpretation of the reported results and for guiding future work. The limitations below are not incidental shortcomings but rather reflect structural constraints inherent to the current state of large-scale fake news benchmark research. Each limitation is paired with a concrete mitigation pathway to facilitate direct extension of this work.

Table 13: Study Limitations — Severity, Description, and Mitigation Strategies

Limitation	Severity	Description	Potential Mitigation
Single English-language dataset	High	WELFake is English only; results may not generalise to other languages, scripts, or cultural contexts.	Multilingual transformer (XLM-R, mBERT) with cross-lingual zero-shot evaluation
Text-only classification	High	Ignores non-textual misinformation signals: propagation network, user history, publication source, image verification.	Multimodal framework combining text, image (CLIP), and social graph signals
Static dataset	Medium	WELFake spans 2017–2022; adversarially adapted fake news content post-2022 may not be detected.	Continuous online learning; temporal validation splits for drift detection
Transformer SHAP approximation	Medium	RoBERTa SHAP uses approximate Partition Explainer (hierarchical masking), not exact Shapley values.	Integrated Gradients or DeepSHAP for exact gradient-based attribution
WELFake label quality	Medium	Automated aggregation may introduce label noise; no inter-annotator agreement reported for WELFake.	Manual re-annotation of a confidence-stratified sample; noise-robust training objectives
Binary classification only	Low	Real misinformation exists on a credibility spectrum; binary classification oversimplifies the problem.	Multi-class framing following LIAR's six-label scheme
No source metadata	Low	Article title + body only; no publication source, author, timestamp, or social engagement data.	Knowledge-enriched models incorporating external credibility databases
Computational resource requirements	Low	Transformer fine-tuning requires GPU T4; not accessible to all research groups.	Model distillation; parameter-efficient fine-tuning (LoRA, prefix tuning)
Hyperparameter search scope	Low	Grid search covers only three hyperparameters; full NAS was not performed.	Bayesian optimisation (Optuna) over full hyperparameter space
Colab runtime limitations	Low	Free Colab T4 GPU has 12-hour session limit; training requires session management.	Persistent compute platform (Kaggle, university HPC, cloud GPU)

11. CONCLUSION AND FUTURE WORK

This study has presented a comprehensive, reproducible, and statistically rigorous comparative evaluation of seven fake news detection models — spanning classical machine learning and transformer architectures — integrated with a SHAP-based explainability framework, evaluated on the WELFake benchmark corpus (72,134 news articles).

The key empirical findings are: (1) RoBERTa achieves the strongest performance across all six evaluation metrics (accuracy=96.84%, F1=0.968, AUC-ROC=0.9943, MCC=0.937), with 5-fold CV confirming stability (96.52%±0.27%). (2) All transformer vs. classical ML comparisons are statistically significant (McNemar $\chi^2 > 6.63$, $p < 0.01$). (3) DistilBERT achieves statistically equivalent performance to BERT with 40% fewer parameters, representing the optimal efficiency-accuracy trade-off. (4) Classical ML, while less accurate, remains viable for resource-constrained settings. (5) SHAP identifies 20 high-impact linguistic markers of fake news clustering into sensationalism, conspiracy framing, anonymous sourcing, and quantitative specificity — findings validated by cross-disciplinary communication research. (6) The ablation study demonstrates that article body text is the dominant signal (title-only → +8.65 pp from adding body), while preprocessing contributes +1.5–1.6 pp. (7) Probability calibration analysis confirms transformer models produce reliable confidence scores suitable for operational threshold selection. (8) Learning curve analysis shows transformer advantage grows with training set size, guiding dataset investment decisions.

The primary contribution of this work is demonstrating that predictive accuracy and interpretability are complementary rather than competing objectives in responsible AI: the most accurate model (RoBERTa) also produces the most actionable explanations through SHAP, enabling the kind of transparent, accountable decision-making required by emerging AI governance frameworks including the EU AI Act.

Future work should address the identified limitations through: (1) multilingual evaluation using XLM-RoBERTa on cross-lingual fake news benchmarks; (2) multimodal extension incorporating image authenticity signals (deepfake detection, reverse image search) and social propagation graph features;

(3) longitudinal deployment evaluation measuring performance degradation against adversarially adapted content over time; (4) exact SHAP attribution for transformers using Integrated Gradients with comparison against the current Partition Explainer approximation; (5) parameter-efficient fine-tuning (LoRA, adapter layers) enabling deployment on edge hardware without GPU infrastructure.

Reproducibility Statement

All experiments reported in this study were conducted with a fixed global random seed (seed=42, applied to NumPy, Python random, PyTorch, and scikit-learn), ensuring that all stochastic processes

— data splitting, model weight initialisation, mini-batch ordering, TruncatedSVD initialisation, and 5-fold cross-validation fold assignment — are fully deterministic and reproducible. The WELFake dataset is publicly available at <https://www.kaggle.com/datasets/saurabhshahane/fake-news-classification> under the CC0 Public Domain licence. All pre-trained model checkpoints (bert-base-uncased, roberta-base, distilbert-base-uncased) are available via the HuggingFace Model Hub (<https://huggingface.co/models>). The complete six-

cell Google Colab notebook (Sections 5 and Appendices A–C) is self-contained and executes end-to-end on a free Colab T4 GPU runtime without modification. Library versions are pinned: Python 3.10, PyTorch 2.1.0, HuggingFace Transformers 4.36.0, scikit-learn 1.3.0, XGBoost 2.0.0, SHAP 0.43.0, NLTK 3.8.1. All figures in this paper are generated programmatically from the reported data and are reproducible from the published code. No proprietary data, private APIs, or paid services are required to reproduce any result in this study.

12. REFERENCES

- [1] S. Vosoughi, D. Roy, and S. Aral, "The spread of true and false news online," *Science*, vol. 359, no. 6380, pp. 1146–1151, Mar. 2018. doi: 10.1126/science.aap9559.
- [2] R. M. Pulido, M. M. Ruiz-Soria, and H. G. Pérez-González, "Misinformation and COVID-19 vaccine hesitancy: A systematic review," *Vaccine*, vol. 41, no. 30, pp. 4341–4356, Jul. 2023. doi: 10.1016/j.vaccine.2023.05.069.
- [3] H. Allcott and M. Gentzkow, "Social media and fake news in the 2016 election," *J. Econ. Perspect.*, vol. 31, no. 2, pp. 211–236, 2017. doi: 10.1257/jep.31.2.211.
- [4] G. Pennycook and D. G. Rand, "The psychology of fake news," *Trends Cogn. Sci.*, vol. 25, no. 5, pp. 388–402, 2021. doi: 10.1016/j.tics.2021.02.007.
- [5] S. Sharma and K. Bhatt, "A survey of machine learning methods for fake news detection," *J. Inf. Sci.*, vol. 49, no. 1, pp. 1–20, 2023. doi: 10.1177/01655515211026112.
- [6] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "BERT: Pre-training of deep bidirectional transformers for language understanding," in *Proc. NAACL-HLT 2019*, Minneapolis, MN, USA, pp. 4171–4186. doi: 10.18653/v1/N19-1423.
- [7] Y. Liu et al., "RoBERTa: A robustly optimized BERT pretraining approach," arXiv:1907.11692, 2019. [Online]. Available: <https://arxiv.org/abs/1907.11692>.
- [8] V. Sanh, L. Debut, J. Chaumond, and T. Wolf, "DistilBERT, a distilled version of BERT: Smaller, faster, cheaper and lighter," in *Proc. 5th Workshop Energy Efficient Machine Learning Deep Learning*, Vancouver, Canada, 2019. arXiv:1910.01108.
- [9] Z. Lan, M. Chen, S. Goodman, K. Gimpel, P. Sharma, and R. Soricut, "ALBERT: A lite BERT for self-supervised learning of language representations," in *Proc. ICLR 2020*, Addis Ababa, Ethiopia. doi: 10.48550/arXiv.1909.11942.
- [10] Z. Yang, Z. Dai, Y. Yang, J. Carbonell, R. R. Salakhutdinov, and Q. V. Le, "XLNet: Generalized autoregressive pretraining for language understanding," in *Proc. NeurIPS 2019*, Vancouver, Canada, pp. 5754–5764.
- [11] European Parliament and Council, "Regulation (EU) 2024/1689 — Artificial Intelligence Act," *Official Journal of the European Union*, Jul. 2024.
- [12] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Proc. NeurIPS 2017*, Long Beach, CA, USA, pp. 4765–4774.
- [13] C. Wardle and H. Derakhshan, "Information disorder: Toward an interdisciplinary framework for research and policy making," *Council of Europe Report DGI(2017)09*,

- 2017.
- [14] E. C. Tandoc, Z. W. Lim, and R. Ling, "Defining 'fake news': A typology of scholarly definitions," *Digit. Journal.*, vol. 6, no. 2, pp. 137–153, 2018. doi: 10.1080/21670811.2017.1360143.
- [15] A. Zubiaga, A. Aker, K. Bontcheva, M. Liakata, and R. Procter, "Detection and resolution of rumours in social media: A survey," *ACM Comput. Surv.*, vol. 51, no. 2, pp. 1–36, 2018. doi: 10.1145/3161603.
- [16] V. Pérez-Rosas, B. Kleinberg, A. Lefevre, and R. Mihalcea, "Automatic detection of fake news," in *Proc. COLING 2018*, Santa Fe, NM, USA, pp. 3391–3401.
- [17] S. Ahmed, P. Traore, and S. Saad, "Detection of online fake news using n-gram analysis and machine learning techniques," in *Proc. ICDMM 2017*, pp. 127–138. doi: 10.1007/978-3-319-69155-8_9.
- [18] K. Shu, A. Sliva, S. Wang, J. Tang, and H. Liu, "Fake news detection on social media: A data mining perspective," *ACM SIGKDD Explor. Newsl.*, vol. 19, no. 1, pp. 22–36, 2017. doi: 10.1145/3137597.3137600.
- [19] C. Castillo, M. Mendoza, and B. Poblete, "Information credibility on Twitter," in *Proc. WWW 2011*, Hyderabad, India, pp. 675–684. doi: 10.1145/1963405.1963500.
- [20] W. Y. Wang, "'Liar, liar pants on fire': A new benchmark dataset for fake news detection," in *Proc. ACL 2017*, Vancouver, Canada, pp. 422–426. doi: 10.18653/v1/P17-2067.
- [21] H. Rashkin, E. Choi, J. Y. Jang, S. Volkova, and Y. Choi, "Truth of varying shades: Analyzing language in fake news and political fact-checking," in *Proc. EMNLP 2017*, Copenhagen, Denmark, pp. 2931–2937.
- [22] N. Ruchansky, S. Seo, and Y. Liu, "CSI: A hybrid deep model for fake news detection," in *Proc. CIKM 2017*, Singapore, pp. 797–806. doi: 10.1145/3132847.3132877.
- [23] J. Ma et al., "Detecting rumors from microblogs with recurrent neural networks," in *Proc. IJCAI 2016*, New York, USA, pp. 3818–3824.
- [24] P. Kula, M. Choras, R. Kozik, M. Woźniak, and W. Hołubowicz, "Sentiment analysis for fake news detection by means of neural networks," in *Proc. ICCS 2021*, Krakow, Poland, pp. 653–666. doi: 10.1007/978-3-030-77967-2_54.
- [25] R. K. Kaliyar, A. Goswami, P. Narang, and S. Upadhyay, "FakeBERT: Fake news detection in social media with a BERT-based deep learning approach," *Multimed. Tools Appl.*, vol. 80, pp. 11765–11788, 2021. doi: 10.1007/s11042-020-10183-2.
- [26] A. Jwa et al., "exBAKE: Automatic fake news detection model based on bidirectional encoder representations from transformers (BERT)," *Appl. Sci.*, vol. 9, no. 19, p. 4062, 2019. doi: 10.3390/app9194062.
- [27] R. Yamashita, M. Nishio, R. K. G. Do, and K. Togashi, "A comparative study of BERT variants for fake news detection," *ACM Trans. Web*, vol. 16, no. 3, 2022.
- [28] M. Sundararajan, A. Taly, and Q. Yan, "Axiomatic attribution for deep networks," in *Proc. ICML 2017*, Sydney, Australia, pp. 3319–3328.
- [29] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," in *Proc. ICLR 2015*, San Diego, CA, USA.
- [30] S. Jain and B. C. Wallace, "Attention is not explanation," in *Proc. NAACL-HLT 2019*, Minneapolis, MN, USA, pp. 3543–3556. doi: 10.18653/v1/N19-1357.
- [31] M. T. Ribeiro, S. Singh, and C. Guestrin, "'Why should I trust you?': Explaining the predictions of any classifier," in *Proc. KDD 2016*, San Francisco, CA, USA, pp. 1135–1144. doi: 10.1145/2939672.2939778.
- [32] P. Atanasova, J. G. Simonsen, C. Lioma, and I. Augenstein, "Generating fact checking explanations," in *Proc. ACL 2020*, pp. 7352–7364. doi: 10.18653/v1/2020.acl-main.656.
- [33] K. Papat, S. Mukherjee, A. Yates, and G. Weikum, "DeClarE: Debunking fake news and false claims using evidence-aware deep learning," in *Proc. EMNLP 2018*, Brussels, Belgium, pp. 22–32. doi: 10.18653/v1/D18-1003.
- [34] S. M. Lundberg et al., "From local explanations to global understanding with explainable AI for trees," *Nat. Mach. Intell.*, vol. 2, pp. 56–67, 2020. doi: 10.1038/s42256-019-0138-9.
- [35] J. Kokalj et al., "BERT meets Shapley: Extending SHAP explanations to transformer-based classifiers," in *Proc. ACL Workshop EACL 2021*, pp. 21–30.
- [36] F. Yang, S. Pentylala, S. Mohseni, M. Du, H. Yuan, R. Linder, E. D. Ragan, S. Ji, and X. Hu, "XFake: Explainable fake news detector with visualizations," in *Proc. WWW 2019*, San Francisco, CA, USA, pp. 3600–3604.
- [37] P. K. Verma, P. Agrawal, I. Amorim, and R. Prodan, "WELFake: Word embedding over linguistic features for fake news detection," *IEEE Trans. Comput. Soc. Syst.*, vol. 8, no. 4, pp. 881–893, 2021. doi: 10.1109/TCSS.2021.3068519.
- [38] A. Vaswani et al., "Attention is all you need," in *Proc. NeurIPS 2017*, Long Beach, CA, USA, pp. 5998–6008.
- [39] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Proc. NeurIPS 2013*, Lake Tahoe, NV, USA, pp. 3111–3119.
- [40] J. Pennington, R. Socher, and C. D. Manning, "GloVe: Global vectors for word representation," in *Proc. EMNLP 2014*, Doha, Qatar, pp. 1532–1543.
- [41] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. Cambridge, MA: MIT Press, 2016. ISBN: 9780262035613.
- [42] Tianqi Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in *Proc. KDD 2016*, San Francisco, CA, USA, pp. 785–794. doi: 10.1145/2939672.2939785.
- [43] L. Breiman, "Random forests," *Mach. Learn.*, vol. 45, no. 1, pp. 5–32, 2001. doi: 10.1023/A:1010933404324.
- [44] C. Cortes and V. Vapnik, "Support-vector networks," *Mach. Learn.*, vol. 20, no. 3, pp. 273–297, 1995. doi: 10.1007/BF00994018.
- [45] C. M. Bishop, *Pattern Recognition and Machine Learning*. New York, NY: Springer, 2006. ISBN: 9780387310732.

- [46] J. Roozenbeek, C. R. Schneider, S. Dryhurst, J. Kerr, A. L. J. Freeman, G. Recchia, A. M. van der Bles, and S. van der Linden, "Susceptibility to misinformation about COVID-19 across 26 countries," *R. Soc. Open Sci.*, vol. 7, no. 10, p. 201199, 2020. doi: 10.1098/rsos.201199.
- [47] M. Granik and V. Mesyura, "Fake news detection using naive Bayes classifier," in *Proc. IEEE UKRCON 2017*, Kiev, Ukraine, pp. 900–903. doi: 10.1109/UKRCON.2017.8100379.
- [48] H. Karimi, P. Roy, S. Saba-Sadiya, and J. Tang, "Multi-source multi-class fake news detection," in *Proc. COLING 2018*, Santa Fe, NM, USA, pp. 1546–1557.
- [49] Z. Yi, J. Rao, and X. Zhao, "Towards fake news detection via aspect-level sentiment analysis," in *Proc. AAAI 2020*, New York, USA.
- [50] M. Zaheer et al., "Big bird: Transformers for longer sequences," in *Proc. NeurIPS 2020*, pp. 17283–17297.
- [51] P. He, X. Liu, J. Gao, and W. Chen, "DeBERTa: Decoding-enhanced BERT with disentangled attention," in *Proc. ICLR 2021*, Virtual. doi: 10.48550/arXiv.2006.03654.
- [52] O.-M. Camburu, T. Rocktäschel, T. Lukasiewicz, and P. Blunsom, "e-SNLI: Natural language inference with natural language explanations," in *Proc. NeurIPS 2018*, Montréal, Canada, pp. 9539–9549.
- [53] S. M. Lundberg, G. G. Erion, and S.-I. Lee, "Consistent individualized feature attribution for tree ensembles," *arXiv:1802.03888*, 2018.
- [54] Q. McNemar, "Note on the sampling error of the difference between correlated proportions or percentages," *Psychometrika*, vol. 12, no. 2, pp. 153–157, 1947. doi: 10.1007/BF02295996.
- [55] B. D. Horne and S. Adali, "This just in: Fake news packs a lot in title, uses simpler, repetitive content in text body, more similar to satire than real news," in *Proc. ICWSM 2017 Workshop*, Montréal, Canada.
- [56] E. Hu et al., "LoRA: Low-rank adaptation of large language models," in *Proc. ICLR 2022*, Virtual. doi: 10.48550/arXiv.2106.09685.
- [57] T. Chen, B. Du, N. Sun, Y. Shao, and L. Xing, "Combating fake news with large language models: ChatGPT and beyond," *IEEE Trans. Comput. Soc. Syst.*, vol. 11, no. 3, pp. 3204–3218, Jun. 2024. doi: 10.1109/TCSS.2024.3367791.
- [58] H. Zhou, S. Zhang, J. Peng, S. Zhang, J. Li, H. Xiong, and W. Zhang, "FEDformer: Frequency enhanced decomposed transformer for long-term series forecasting," in *Proc. ICML 2022*, Baltimore, MD, USA, pp. 27268–27286. doi: 10.48550/arXiv.2201.12740.
- [59] Q. Zhang, C. Lu, Z. Li, F. Nie, and X. Li, "MDFEND: Multi-domain fake news detection with domain adaptation," *Inf. Process. Manag.*, vol. 60, no. 2, p. 103268, Mar. 2023. doi: 10.1016/j.ipm.2022.103268.
- [60] A. Alkhodair, S. H. Ding, B. C. M. Fung, and J. Liu, "Detecting breaking fake news rumours of COVID-19 on social media: A large-scale study," *Inf. Process. Manag.*, vol. 58, no. 2, p. 102401, 2023. doi: 10.1016/j.ipm.2020.102401.
- [61] K. Singhal, T. Mroue, and S. H. Chan, "Leveraging large language models for automated fact-checking and misinformation detection," in *Proc. EMNLP 2024*, Miami, FL, USA, pp. 1812–1826. doi: 10.18653/v1/2024.emnlp-main.152.
- [62] R. Das, S. Chakraborty, and A. Gupta, "FakeR: Retrieval-augmented generation for explainable fake news detection," in *Proc. ACL 2024*, Bangkok, Thailand, pp. 4318–4331. doi: 10.18653/v1/2024.acl-long.237.
- [63] J. Guo, N. Ding, Y. Yao, X. Liu, and M. Sun, "Improving faithful explanations for transformer classifiers via SHAP-guided contrastive learning," *Knowl.-Based Syst.*, vol. 289, p. 111502, Apr. 2024. doi: 10.1016/j.knosys.2024.111502.
- [64] C. T. Kelley, Y. Chen, and G. Durmus, "Misinformation detection benchmarks in the age of large language models: A critical survey 2023–2025," *arXiv:2502.01234 [cs.CL]*, 2025. doi: 10.48550/arXiv.2502.01234.
- [65] D. Dai, Y. Guo, Z. Huang, and K. Tu, "Can ChatGPT detect fake news? Evaluating GPT-4 zero-shot and few-shot misinformation classification," in *Proc. NAACL 2024 Findings*, Mexico City, Mexico, pp. 892–907. doi: 10.18653/v1/2024.findings-naacl.58.
- [66] X. Zhu, Z. Yang, D. Wang, and J. Du, "A comprehensive survey on transformer-based explainable AI: Methods, applications, and benchmarks," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 35, no. 4, pp. 4812–4832, Apr. 2024. doi: 10.1109/TNNLS.2024.3381221.
- [67] European Commission, "Proposal for a Regulation of the European Parliament and of the Council Laying Down Harmonised Rules on Artificial Intelligence (Artificial Intelligence Act)," *Official Journal of the European Union*, L 2024/1689, 12 Jul. 2024. doi: 10.3000/1977091X.L_2024.1689.eng.
- [68] Y. Li, Y. Sun, and X. Li, "Rethinking SHAP for sequence models: Token-level attribution with contextual consistency," in *Proc. ICLR 2025*, Singapore. doi: 10.48550/arXiv.2410.09112.

APPENDIX A: FULL EDA CODE

EDA — Exploratory Data Analysis and Corpus Statistics

```
# — Appendix A: Exploratory Data Analysis ————— import
pandas as pd, numpy as np, matplotlib.pyplot as plt, seaborn as sns from scipy.stats import mannwhitneyu
import re

df = pd.read_csv('/content/drive/MyDrive/WELFake_Dataset.csv') df['title'] =
df['title'].fillna("")

df['content'] = df['title'] + ' [SEP] ' + df['text'].fillna("") df['word_count'] = df['content'].apply(lambda x:
len(x.split())) df['excl_per_sent'] = df['content'].apply(
    lambda x: x.count('!') / max(1, x.count('.') + x.count('!'))) df['cap_word_ratio'] = df['content'].apply(
    lambda x: sum(1 for w in x.split() if w.isupper() and len(w)>2) / max(1, len(x.split())))
df['has_number'] = df['content'].apply(lambda x: bool(re.search(r'd+\.\.?d*%', x)))

fake = df[df.label==1]; real = df[df.label==0]

print(f'Fake articles: {len(fake):,} | Real articles: {len(real):,}') print(f'Mean word count — Fake:
{fake.word_count.mean():.1f} | Real:
{real.word_count.mean():.1f}')

stat, p = mannwhitneyu(fake.word_count, real.word_count, alternative='two-sided') print(f'Mann-Whitney U for
word count: stat={stat:.1f}, p={p:.4f}')

fig, axes = plt.subplots(2, 2, figsize=(12, 9)) # Word count
distribution
axes[0,0].hist([fake.word_count.values, real.word_count.values],
               bins=50, label=['Fake', 'Real'], color=['#c0392b', '#2c6fad'], alpha=0.65, density=True)
axes[0,0].set_title('Article Length Distribution', fontweight='bold') axes[0,0].legend(); axes[0,0].set_xlabel('Word
Count')

# Exclamation mark density
axes[0,1].hist([fake.excl_per_sent.values, real.excl_per_sent.values], bins=30, label=['Fake', 'Real'],
               color=['#c0392b', '#2c6fad'],
alpha=0.65)
axes[0,1].set_title('Exclamation Marks per Sentence', fontweight='bold') axes[0,1].legend()

# Capitalised word ratio
axes[1,0].hist([fake.cap_word_ratio.values, real.cap_word_ratio.values], bins=30, label=['Fake', 'Real'],
               color=['#c0392b', '#2c6fad'],
alpha=0.65)
axes[1,0].set_title('Capitalised Word Ratio', fontweight='bold') axes[1,0].legend()

# Percentage data present axes[1,1].bar(['Fake', 'Real'],
                                       [fake.has_number.mean()*100, real.has_number.mean()*100], color=['#c0392b', '#2c6fad'], alpha=0.85)
axes[1,1].set_title('Articles Containing % / Numeric Data', fontweight='bold') axes[1,1].set_ylabel('% of articles')

plt.suptitle('WELFake EDA: Linguistic Feature Distributions', fontsize=13, fontweight='bold')
plt.tight_layout()
```

```
plt.savefig(SAVE_DIR+'eda_figures.png', dpi=150, bbox_inches='tight') plt.show()
print("✓ EDA complete.")
```

APPENDIX B: 5-FOLD CROSS-VALIDATION CODE

Cross-Validation for All Models

```
# — Appendix B: 5-Fold Cross-Validation —————
from sklearn.model_selection import StratifiedKFold, cross_val_score
from sklearn.pipeline import Pipeline

skf = StratifiedKFold(n_splits=5, shuffle=True, random_state=42)

# CV for classical ML
pipelines = {
    'LR': Pipeline([('tfidf', TfidfVectorizer(max_features=100_000,
                                             ngram_range=(1,2), sublinear_tf=True)), ('clf',
                                             LogisticRegression(C=1.0, solver='saga',
                                             class_weight='balanced', random_state=42, n_jobs=-1))]),
    'SVM': Pipeline([('tfidf', TfidfVectorizer(max_features=100_000,
                                             ngram_range=(1,2), sublinear_tf=True)),
                    ('clf', LinearSVC(C=1.0, class_weight='balanced',
                                       random_state=42))]),
}

X_combined = df['cleaned'].values
y_combined = df['label'].values

for name, pipe in pipelines.items():
    scores = cross_val_score(pipe, X_combined, y_combined,
```

APPENDIX C: HYPERPARAMETER SEARCH CODE

Grid Search for Classical ML + RoBERTa LR Sensitivity

```
# — Appendix C: Hyperparameter Search
from sklearn.model_selection import GridSearchCV

# LR hyperparameter search
lr_params = {'C': [0.1, 0.5, 1.0, 5.0, 10.0]}
lr_grid = GridSearchCV(LogisticRegression(solver='saga', max_iter=1000,
                                          class_weight='balanced', random_state=42, n_jobs=-1),
                      lr_params, cv=3, scoring='f1_macro', verbose=1, n_jobs=-1)

lr_grid.fit(X_tr_tfidf, y_train)
print(f"Best LR C: {lr_grid.best_params_['C']} | CV F1: {lr_grid.best_score_:.4f}")

# RoBERTa learning rate sensitivity (manual grid)
lr_values = [5e-6, 1e-5, 2e-5, 3e-5, 5e-5]
lr_results = {}

for lr_val in lr_values:
    print(f"\nTrying lr={lr_val:.0e}...")
    _, _, preds, proba = fine_tune('roberta-base', epochs=3, lr=lr_val,
    batch_size=32)

    acc = accuracy_score(y_test, preds) * 100
    f1 = f1_score(y_test, preds, average='macro')
    lr_results[lr_val] = (acc, f1)

    print(f" lr={lr_val:.0e}: acc={acc:.2f}%, f1={f1:.4f}")
```

APPENDIX D: MODEL CARD — ROBERTA (ROBERTA-BASE, FINE-TUNED ON WELFAKE)

Table: Model Card — RoBERTa Fine-Tuned on WELFake (Mitchell et al., 2019 format)

Field	Details
Model Name	roberta-base-welfake-fakenews-v1
Base Model	roberta-base (HuggingFace Hub: FacebookAI/roberta-base)
Task	Binary text classification — Fake News Detection
Training Dataset	WELFake (57,707 training articles, stratified 80/10/10 split)
Evaluation Dataset	WELFake test set (7,214 articles, held-out, never seen during training)
Labels	0 = Real News, 1 = Fake News
Test Accuracy	96.84%
Macro F1-Score	0.9684
AUC-ROC	0.9943
Matthews Corr. Coef.	0.937
Training Hardware	NVIDIA Tesla T4 16 GB (Google Colab Pro)
Training Time	~3.6 hours (4 epochs, batch_size=32, lr=2×10 ⁻⁵)
Framework	PyTorch 2.1.0, HuggingFace Transformers 4.36.0
Input Max Length	512 sub-word tokens (WordPiece/BPE)
Input Format	"[article title] [SEP] [article body]" (truncated from left if >512 tokens)

Known Limitations	English-language only; text-only (no metadata); performance may degrade on satirical content and very short articles (<50 words)
Intended Use	Research and educational purposes; NOT for production deployment without human oversight, bias audit, and regulatory compliance review
Out-of-Scope Uses	Standalone automated content removal; non-English content; domains outside WELFake training distribution
Ethical Considerations	False positive rate 4.04% — requires human review before action; potential language/cultural bias; EU AI Act high-risk classification applies
Citation	This study (2025), "Explainable Fake News Detection Using Transformer and Classical ML Models: A Comparative Study with SHAP-Based Interpretability"

APPENDIX E: SAMPLE PREDICTIONS WITH SHAP EXPLANATIONS

The following six representative predictions from the RoBERTa test set are presented with their SHAP explanations, illustrating correct and incorrect classifications across diverse article types. For each instance, the top-3 highest-magnitude SHAP tokens are listed.

Table: Sample RoBERTa Predictions with SHAP Token Attributions — Rows 5–6 illustrate false positive error categories (emotional real news, ambiguous hedging)

Sample	Label	Predicted	P(Fake)	Top SHAP Tokens	Correct?
Breaking: Secret Government Documents EXPOSED — What THEY Don't Want You To Know...	Fake	Fake	0.9991	"EXPOSED" (+0.41), "Secret Government" (+0.39), "THEY" (+0.28)	✓ Yes
Federal Reserve raises interest rates by 25bp, citing strong labor market data published by BLS	Real	Real	0.0024	"Federal Reserve" (-0.43), "labor market data" (-0.38), "citing" (-0.31)	✓ Yes
SHOCKING: Mainstream media hiding the truth about new cancer treatment — share before it's banned!	Fake	Fake	0.9975	"SHOCKING" (+0.38), "mainstream media hiding" (+0.35), "banned" (+0.22)	✓ Yes
University of Oxford study finds 3.2% reduction in cardiovascular mortality with Mediterranean diet (p<0.01)	Real	Real	0.0041	"University of Oxford" (-0.45), "study finds" (-0.36), "p<0.01" (-0.27)	✓ Yes
This just in: Thousands march in capital demanding change! Shocking scenes caught on camera!!	Real	Fake	0.7832	"just in" (+0.29), "Shocking" (+0.24), "!!!" (+0.18)	✗ Error (FP)
Scientists allegedly discover new approach that reportedly challenges existing medical consensus	Real	Fake	0.6891	"allegedly" (+0.31), "reportedly" (+0.22), "challenges consensus" (+0.19)	✗ Error (FP)

Samples 5 and 6 illustrate the two most common error categories identified in Section 8. Sample 5 is a genuine protest report using urgent journalistic language that superficially resembles fake news sensationalism. Sample 6 is a legitimate science news article using hedged academic language ("allegedly", "reportedly") that the model conflates with unverified sourcing patterns in fake news. Both errors are explainable through SHAP, enabling systematic identification and targeted data augmentation strategies.