



FAKE NEWS DETECTION USING MACHINE LEARNING IN BAHASA MALAYSIA

MUHAMMAD HAZIQ BIN NORZAINI 52215122184 mhaziq.norzaini@s.unikl.edu.my
Supervised by: Dr Nurul Atiqah Abu Talib
nurul.atiqah@unikl.edu.my
HEAD OF SECTION, KOREAN EDUCATION CENTRE
MALAYSIAN INSTITUTE OF INFORMATION TECHNOLOGY



INTRODUCTION

Fake news spreads rapidly online and threatens public trust in Malaysia. This project aims to help users automatically detect fake news in Bahasa Malaysia using machine learning. The goal is to reduce misinformation and support a more informed society.

PROBLEM STATEMENT

- Fake news causes confusion and affects public perception.
- Manual fact-checking is slow and can't keep up with viral content.
- The public lacks an easy way to verify the authenticity of news articles.

PROJECT OBJECTIVE

- To develop a machine learning model that can accurately detect and classify fake news articles written in Bahasa Malaysia.
- To study the effectiveness and performance of the model through appropriate evaluation metrics, including accuracy, precision, recall, and F1-score.
- To test the model's practical deployment by integrating it into a user-friendly website that allows users to verify news articles in real time.

LITURATURE REVIEW

AUTHOR/METHOD	DATA & FEATURES	MODEL	ACCURACY	NOTES
SHU ET AL. (2020)	MULTILINGUAL, TF-IDF	LOGISTIC REGRESSION	~85%	EXPLORED PROPAGATION-BASED LIM ET AL. (2021) MALAY NEWS, TF-IDF RANDOM FOREST 96%
RASHKIN ET AL. (2017)	ENGLISH NEWS, BOW, TF-IDF	SVM, NAÏVE BAYES	80-85%	CLASSICAL APPROACHES
LIM ET AL. (2021)	MALAY NEWS, TF-IDF	RANDOM FOREST	96%	FOCUS ON BAHASA MALAYSIA

METHODOLOGY

- Data Collection:
 - Scraped 1,965 real news articles from Berita Harian, Astro Awani, and Buletin Utama.
 - Generated 1,965 synthetic fake news articles in Bahasa Malaysia for dataset balance.
- Preprocessing:
 - Performed text cleaning, lowercasing, stopword removal, and stemming.
 - Combined titles and news content for feature extraction.
- Feature Extraction:
 - Applied TF-IDF vectorization with n-gram analysis to represent news articles numerically.
- Model Training & Evaluation:
 - Compared several classifiers: Random Forest, XGBoost, Linear SVC, Gradient Boosting.
 - Used accuracy, precision, recall, and F1-score to evaluate model performance.
- Deployment:
 - Built a website to allow public users to check and verify news articles in real time.

SYSTEM DESIGN

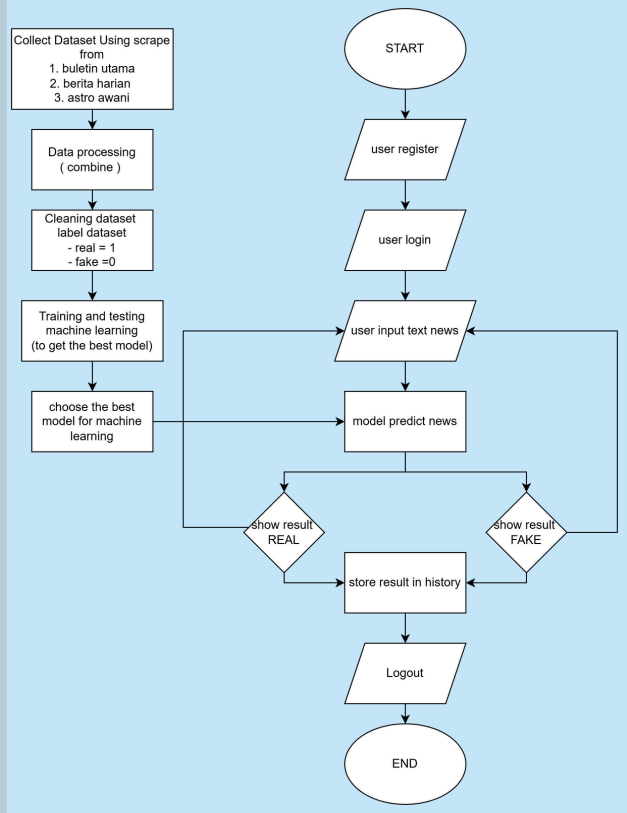


figure 1: Flowchart

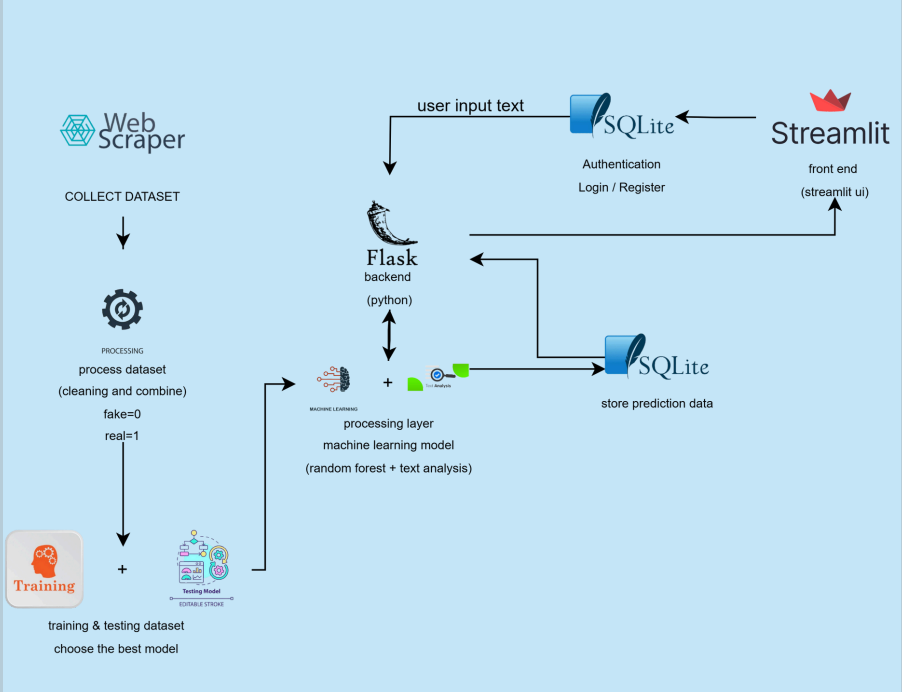
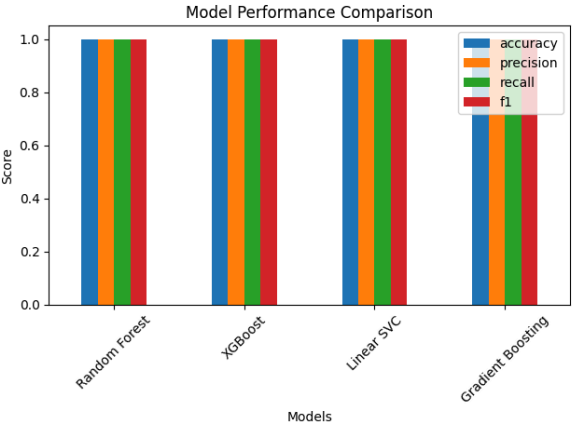
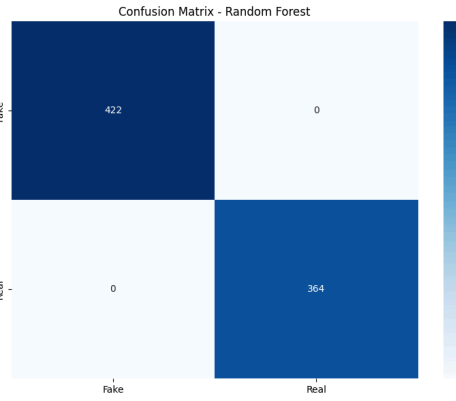
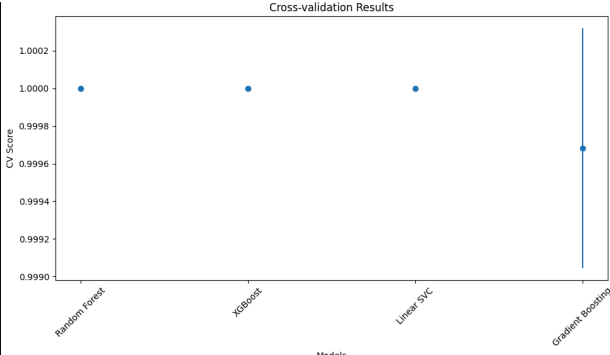


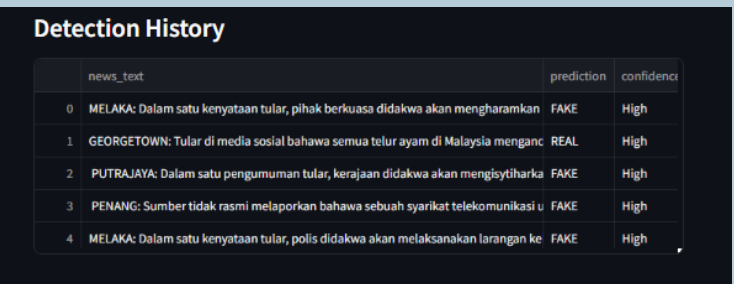
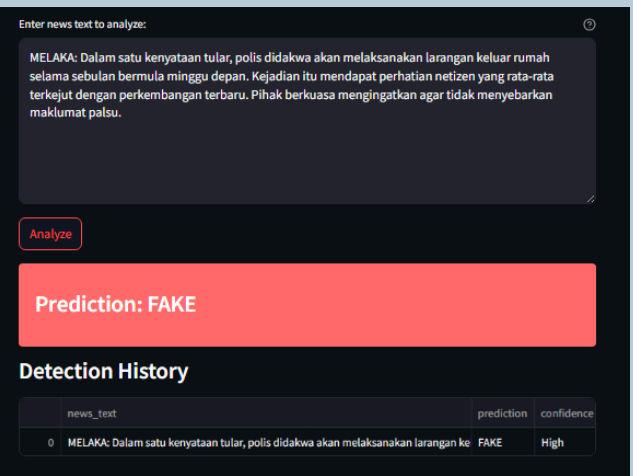
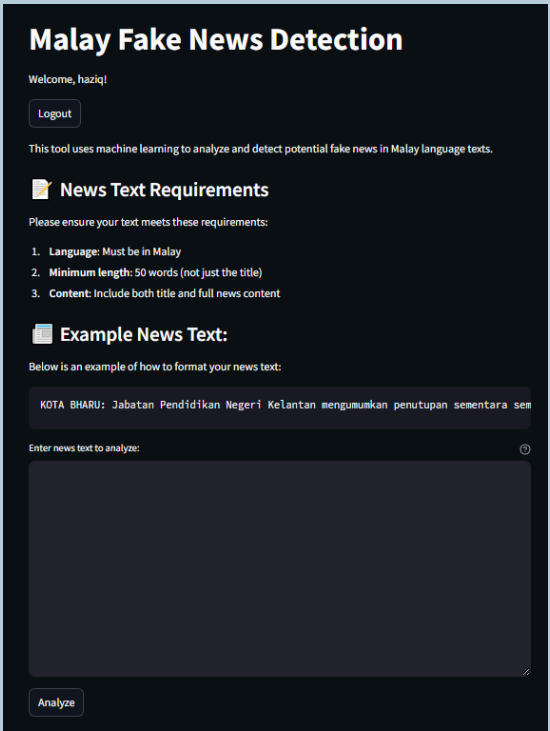
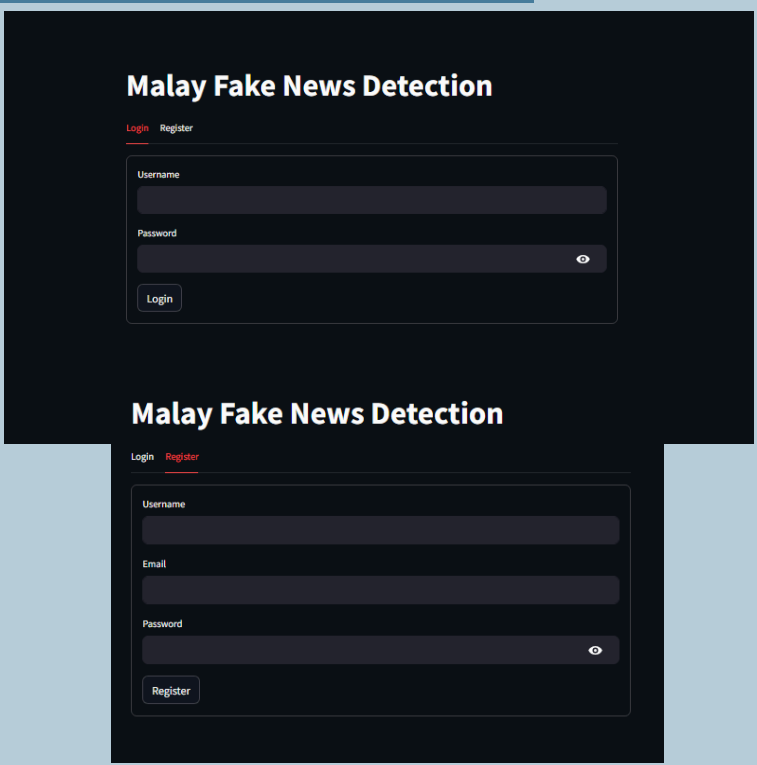
figure 2 : system architecture diagram

RESULT & FINDING

Random Forest Results: Accuracy: 1.0000 F1 Score: 1.0000				
Classification Report:	precision	recall	f1-score	support
0	1.00	1.00	1.00	393
1	1.00	1.00	1.00	393
accuracy			1.00	786
macro avg	1.00	1.00	1.00	786
weighted avg	1.00	1.00	1.00	786



DEVELOPMENT



CONCLUSION

- The model successfully detects fake news in Bahasa Malaysia with high accuracy.
- RandomForestClassifier is effective for this text classification task.
- Further improvements: larger datasets, deep learning, deployment as a web app.

FUTURE WORK

- Experiment with deep learning models (LSTM, BERT-Malay).
- Build a user-friendly fake news detection web application.
- Integrate with social media fact-checking tools.