# table of contents

01

introduction

# overview
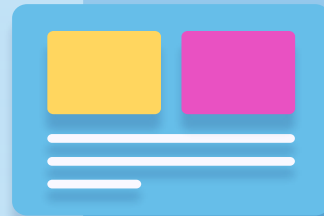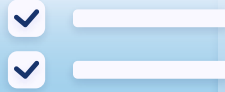
A credibility-aware news recommendation system that delivers relevant **AND** trustworthy content to readers

*Unlike standard systems that prioritize relevance alone, our solution integrates credibility assessment to filter misinformation while maintaining content relevance*

# problem statement - stakeholders

Standard recommendation systems promote fake news as easily as verified content

- This erodes user trust and damages platform reputation
- Readers make less informed decisions due to misinformation exposure

**Primary stakeholder**
- The end user, specifically news readers who rely on digital platforms for information that are directly affected by exposure to misinformation
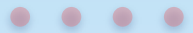
**Secondary stakeholders**
- News platforms and aggregators seeking to maintain credibility and retain users
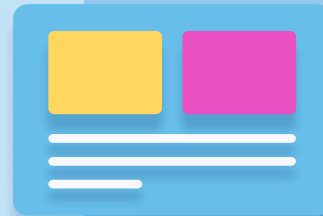
# dataset description

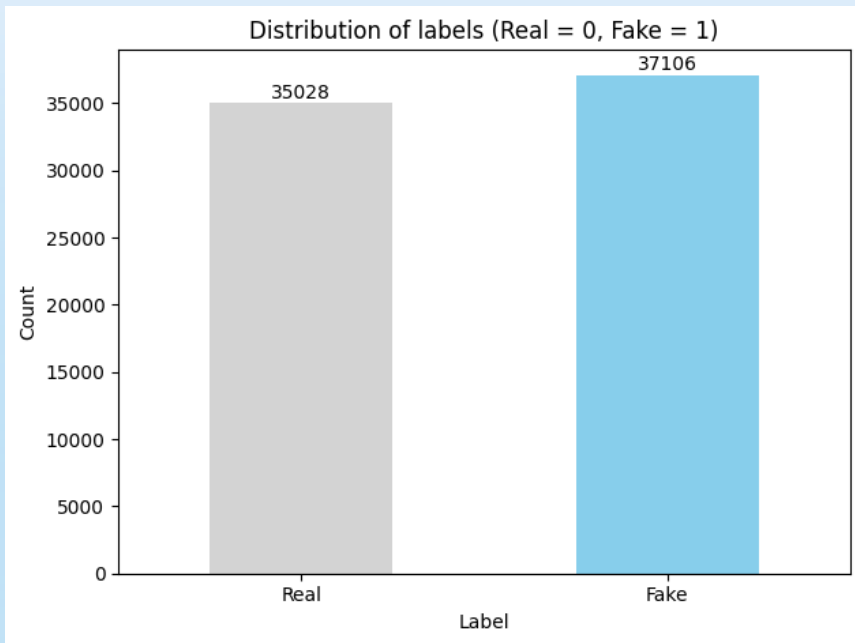**WELFake dataset: 72,134 news articles from multiple sources**

- Created by merging Kaggle, McIntire, Reuters, and BuzzFeed Political datasets
- Balanced distribution: 48.6% real news, 51.4% fake news
- Final cleaned dataset: **63,121** unique articles after removing duplicates and missing values

# data processing

# label distribution



Distribution of labels (Real = 0, Fake = 1)

Near-perfect balance: 48.6% real news, 51.4% fake news

- No artificial data manipulation required
- Ensures unbiased credibility assessment
- Reflects real-world news ecosystem proportions

*This balanced dataset enables reliable credibility scoring without overemphasizing either news type*

# feature engineering

**Title Length**: **Fake news headlines are 14 characters longer on average**
Real: 68.8 characters | Fake: 82.9 characters
*Most powerful single indicator (66.2% accuracy)*

**Exclamation Count: Fake news uses 50x more exclamation marks**
Real: 0.002 per headline | Fake: 0.107 per headline

**Question Marks: Fake news uses 2.6x more questions**
Real: 2.5% of headlines | Fake: 6.5% of headlines

**Capitalization Ratio: Fake news uses nearly 2x more ALL-CAPS text**
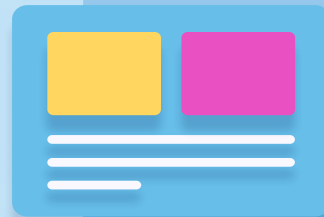Real: 1.0% capitalization | Fake: 1.9% capitalization

*These headline patterns serve as early warning signals that help identify misinformation before readers click through to full articles*

# credibility assessment

# how we measure trustworthiness

- Uses out-of-sample prediction to prevent data leakage

- Scores represent probability that an article is real (0 = fake, 1 = real)

- Maintains original dataset balance through stratified cross-validation

# what we analyze in headlines

- Punctuation patterns (exclamation marks, question marks)

- Capitalization usage

- Structural elements and length

*Unlike standard systems, our credibility scores come from models that never saw the specific article during training, ensuring unbiased, reliable assessments*

# transparent explanations

SHAP Explainability reveals why articles receive specific credibility scores:
- Identifies exact words/phrases that influence trustworthiness
- Shows how "CDC" increases credibility or "BREAKING!!!" decreases it
- Uses game theory approach for mathematically sound explanations
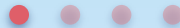
**Reader-friendly output**
- Clear "High/Moderate/Low" credibility labels
- Specific reasons like "Reputable health authority mentioned"
- Actionable guidance: "Check sources before sharing"

*This transforms opaque credibility scores into actionable insights readers can understand in seconds, not technical details*

# recommendation system

# how we match content to readers

**Content-based similarity**
- Analyzes headlines to find topically related articles
- Uses advanced text analysis to understand what makes articles similar
- Focuses on what interests readers (not just popularity)
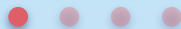
**Why this matters**
- 89% of news consumption happens on mobile devices
- Readers need relevant content in seconds, not minutes
- Standard systems often recommend similar but misleading content

**Hybrid scoring algorithm combines two factors**
- 50% Content Relevance: What interests you
- 50% Source Credibility: Is it trustworthy?

**Smart topic-sensitive filtering**
- Health topics: 45% credibility minimum
- Political topics: 40% credibility minimum
- Standard topics: 30% credibility minimum
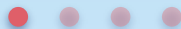
# perfect balance

**Rigorous testing**
- Tested multiple weight combinations (50-80% similarity)
- Measured impact on fake news reduction and relevance

**The winning formula: 50/50 balance**
- Achieved 91% precision@5 (91 out of 100 top recommendations are real news)
- Reduced fake news in recommendations by 38 percentage points (from 46% to 9%)
- Maintained strong relevance (0.42-0.61 topic match)

*This balance delivers the best protection against misinformation while keeping content relevant to readers*

# real-world results

**When users view fake news (0.04 credibility)**
- System recommends articles with 0.63 average credibility while maintaining 0.42 topic relevance

**When users view credible content (0.99 credibility)**
- System maintains high standards (0.97 average credibility) with 0.61 topic relevance

# real-world results

**1. Dynamic Credibility Thresholds**
Problem: Fixed thresholds can't adapt to evolving misinformation tactics
Solution: Automatically adjust scrutiny levels during emerging threats (e.g., health crises)
Business Value: Stay ahead of misinformation without manual intervention while maintaining the 50/50 relevance/credibility balance

# recommendations

## 1. Dynamic credibility thresholds

- **Problem**: Fixed thresholds can't adapt to evolving misinformation tactics
- **Solution**: Automatically adjust scrutiny levels during emerging threats
- **Business value**: Stay ahead of misinformation without manual intervention while maintaining the 50/50 relevance/credibility balance

## 2. Personalized credibility settings

- **Problem**: One-size-fits-all approach ignores different user needs and contexts
- **Solution**: Simple slider for users to adjust relevance/credibility balance (70/30 to 30/70)
- **Business value**: Increase user satisfaction while maintaining system integrity: breaking news seekers prioritize relevance, health researchers prioritize credibility

## 3. Cross-platform credibility sharing

- **Problem**: Misinformation spreads across platforms while systems operate in isolation
- **Solution**: Privacy-preserving protocol to share anonymized credibility signals across platforms
- **Business value**: Collective defense against coordinated misinformation campaigns through shared linguistic fingerprint detection

# THANKS!

EMAIL : woodnaliesjoseph@gmail.com
LINKEDIN : www.linkedin.com/in/woodnalie-s-josephaa0011175
GITHUB : https://github.com/nah-yah