



# ELKP: Knowledge-Powered Rumor Detection

Enhancing social media fact-checking with external knowledge graphs. A minimalist implementation exploring how "world knowledge" bridges the gap between language understanding and factual verification.

# Why Language Models Get Fooled

## The Core Problem

Standard models focus on grammar and style, not facts. A fake tweet can sound perfectly official.

## The Limitation

BERT analyzes linguistic patterns—tone, formality, structure—but lacks access to factual databases.

## The Risk

"The Eiffel Tower has been sold to a private investor" sounds credible, yet is completely fabricated.

## The Gap

Without external verification, AI cannot distinguish well-written fiction from truth, leading to false positives.

# The ELKP Framework



## What is ELKP?

Entity-Link Knowledge-Powered framework—a system that extracts entities from claims and cross-references them against trusted sources.



## The Strategy

Instead of immediate classification, we first "fact-check" claims against Wikipedia's knowledge base before making decisions.



## The Paradigm Shift

Moving from "Does this sound like a rumor?" to "Does this contradict documented facts?"—a fundamental rethinking of verification.

# Bridging the Knowledge Gap



01

## Entity Identification

Extract core entities from the claim using NER

02

## Knowledge Retrieval

Query trusted sources for contextual information

03

## Context Fusion

Combine retrieved knowledge with original content

04

## Classification

Analyze augmented text for contradictions

The research paper identifies that rumors frequently contain verifiable entities—people, places, events—that can be cross-referenced against established knowledge sources.

# The ELKP Pipeline



The system transforms a raw tweet like "Michael Jackson seen on an island" by detecting entities, retrieving contradictory facts ("Died in 2009"), and feeding the augmented context to BERT for final classification.

# Modular & Clean Architecture



## prepare\_data.py

### The ETL Pipeline

Parses complex nested JSON from PHEME dataset into clean CSV format for analysis.



## preprocess.py

### The Knowledge Engine

KnowledgeEngine class handles Spacy NER and Wikipedia API to inject contextual knowledge.



## main.py

### The Trainer

Manages data splitting to prevent leakage, fine-tunes BERT, calculates performance metrics.



## demo\_notebook.ipynb

### Interactive Testing

Real-time interface for testing the system with custom inputs and exploring predictions.

# Performance Metrics

85.14%

Accuracy

Overall classification accuracy

0.81

Rumor F1

F1-score for rumor detection

0.88

Non-Rumor F1

F1-score for verified content

## Model Configuration

Architecture: BERT-Base-Uncased fine-tuned for 3 epochs on augmented PHEME dataset.

**Key Observation:** Significant improvement over baseline raw BERT by leveraging injected knowledge context. The model demonstrates stronger performance on non-rumor classification, suggesting effective knowledge integration.

## ERROR ANALYSIS

# When the Model Fails

"The PM's office releases a statement about #sydneyseige."

**Prediction:** Real (Incorrect) | **Actual:** Rumor

## NER Failure

Failed to recognize "#sydneyseige" as a key event entity—hashtag format confused the detector.

## Bad Retrieval

Retrieved Wikipedia info about the word "About" instead of the Sydney siege event.

## Garbage In, Garbage Out

Without proper context, model defaulted to analyzing formal tone, missing the rumor entirely.

# Where Do We Go From Here?



## Real-World Applications

- Automated news filtering systems
- Crisis management platforms
- Social media content moderation

## System Strengths

High accuracy on verifiable claims with robust, modular architecture enabling easy experimentation and deployment.

## Future Improvements

- Upgrade to Transformer-based NER (RoBERTa)
- Integrate real-time search APIs (Google) instead of static Wikipedia
- Multi-lingual knowledge retrieval for global coverage

# Thank You

## Questions & Discussion

Exploring the intersection of natural language processing, knowledge graphs, and misinformation detection.

