

# 概要介绍

## 前言

在人工智能与信息传播技术高速发展的背景下，伪造新闻的泛滥已成为威胁社会信任的重要问题。虚假信息不仅误导公众，还可能引发社会恐慌，甚至影响政策决策。针对这一挑战，本项目“面向新闻场景的伪造检测平台”应运而生。该平台创新性地结合深度学习语言模型与因果推理、检索增强生成等技术，致力于解决传统伪造新闻检测方法在准确性、实时性和可解释性上的不足，为新闻读者、从业者及内容平台提供高效、可信的假新闻检测工具。

## 创意描述

本项目的核心创新点在于将高性能预训练语言模型（Chinese-RoBERTa-wwm-ext）与因果分析（Causal Analysis）、检索增强生成（RAG）及思维链（CoT）推理技术深度融合，使检测过程不仅“知其然”，更能“知其所以然”。传统方法通常依赖表面特征或黑箱模型，导致结果难以解释且易受对抗攻击。相比之下，我们的方案通过多维度分析（文本特征提取、外部知识检索、逻辑推理和因果分析），显著提升了检测精度与可解释性。

此外，我们引入 DGM4 框架（Detecting and Grounding Multi-Modal Media Manipulation and Beyond），采用 HAMMER 模型，通过融合图像和文本信息，增强对复杂伪造手段的识别能力。该模型不仅能检测伪造内容，还能定位伪造区域，帮助用户直观理解新闻的真实性。

## 功能简介

### 1. 高精度新闻文本真伪识别

基于微调后的 Chinese-RoBERTa-wwm-ext 模型，平台可快速、准确地分类新闻真伪，目标准确率>90%，处理速度达 1000 字/2 秒。

### 2. 可解释性推理过程生成

结合 RAG 检索的外部知识与 CoT 生成的推理步骤，平台提供详细的判断依据，例如展示模型如何对比新闻内容与权威事实库，使用户清晰理解检测逻辑。

### 3. 因果关系与特征分析

通过因果分析算法，探究文本特征（如特定词汇、句式）对判断的影响，帮助用户识别虚假新闻的常见模式。

### 4. 多模态伪造检测

DGM4 采用双流架构分析新闻中的图像与文本一致性，识别篡改图片、误导性配图等，并标注伪造区域，提升复杂场景下的检测能力。

### 5. 浏览器插件一键检测

提供 Chrome、Firefox 等浏览器的插件支持，用户点击即可实时核查新闻，操作便捷，无缝融入日常信息获取流程。

这些功能相互协作，共同实现了构建一个准确、高效且高度透明的新闻伪造检测工具的整体目标。

## 特色综述

### 1. 领先的模型与性能

采用表现优异的 **Chinese-RoBERTa-wwm-ext** 中文预训练模型，结合精心设计的训练策略，确保高检测准确率（目标 90% 以上）和处理效率（目标 1000 字/2 秒内）。引入 **HAMMER** 模型，结合图像和文本数据进行联合推理，增强了新闻内容中图文不符、视频篡改等复杂伪造检测的能力。由此不仅可以精准识别伪造新闻，还能对伪造区域进行定位，帮助用户直观地理解新闻伪造的具体情况。

### 2. 深度可解释性

结合 CoT 和因果分析，提供分步骤推理链条，打破“黑箱”模式。

可视化伪造区域定位，提升结果可信度。

### 3. 动态知识增强

通过 RAG 技术实时检索外部知识库（如事实核查数据库），适应新型谣言。

### 4. 用户友好设计

浏览器插件即点即用，支持 API 接入新闻平台或社交媒体审核系统。

## 开发工具与技术

软件：Python（模型开发）、Hugging Face Transformers（RoBERTa）、JavaScript（插件）、FastAPI（后端）。

硬件：推荐 NVIDIA GPU 加速训练与推理。

关键技术：NLP、深度学习、因果推断、RAG、DGM4。

兼容性：支持 Windows 10+ 及主流浏览器（Chrome/Firefox）。

## 应用对象

普通用户：快速识别虚假新闻，提升媒介素养。

新闻从业者：辅助内容真实性核查。

内容平台：集成 API 自动化审核信息流。

研究人员：提供数据分析工具。

## 应用场景

社交媒体：实时拦截虚假信息。

新闻网站：嵌入检测插件或 API。

学术研究：虚假信息传播模式分析。

## 结语

本平台通过创新技术融合与用户友好设计，为伪造新闻检测提供了高效、透明的解决方案。未来将持续优化检测能力，拓展应用场景，助力构建更健康的信息生态。我们相信，通过本项目，能够为提升新闻信息环境的健康度、增强公众媒介素养、辅助专业人员进行内容核查带来显著的社会价值和应用价值，并在中国大学生服务外包创新创业大赛中取得优异成绩。