

# Contrastive Pre-Training Strategies and Cross-Lingual Transfer in CodeT5 for Rumor Detection

Assignee Research

June 9, 2026

## Abstract

This report synthesises findings from 9 peer-reviewed papers addressing the following research question: What is the impact of contrastive pre-training strategies on the cross-lingual transfer performance of CodeT5 for rumor detection when evaluated on Weibo and Twitter benchmarks. 16 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.8/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: A Unified Contrastive Transfer Framework with Propagation Structure for Boosting Low-Resource Rumor Detection. Research question: What is the impact of contrastive pre-training strategies on the cross-lingual transfer performance of CodeT5 for rumor detection when evaluated on Weibo and Twitter benchmarks?.

## 2 Methodology

Systematic literature search across multiple databases yielded 9 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 3.8/10.

## 3 Results

9 papers retrieved. 16 claims extracted; 0 independently verified. Quality review score: 3.8/10.

## 4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.



## 5 Extracted Claims

Claim	Verified	Confidence
There are no public benchmarks available for detecting low-resource rumors with propagation tree structure in tweets.	×	0.12
The study considers the COVID-19 breaking event as a low-resource domain.	×	0.10
Data was collected from Twitter in English, Cantonese, and Arabic, and from Sina Weibo in Chinese.	×	0.04
The Twitter datasets (English-COVID19, Cantonese-COVID19, Arabic-COVID19) were extended by collecting propagation thread	×	0.03
The Chinese-COVID19 dataset from Sina Weibo gathered rumorous claims from the Sina community management center.	×	0.01
Non-rumorous claims in the Chinese-COVID19 dataset were obtained by randomly filtering out posts not reported as rumors.	×	0.03
All datasets in the study contain two binary labels: Rumor and Non-rumor.	×	0.05
The EngCovid dataset contains 1,154 events and 60,409 tree nodes.	×	0.02
The ChiCovid dataset contains 4,649 events and 1,956,449 tree nodes.	×	0.02
The CanCovid dataset contains 400 events and 406,185 tree nodes.	×	0.02
The AraCovid dataset contains 399 events and 26,687 tree nodes.	×	0.02
The average depth per tree in the Cantonese-COVID19 (CanCovid) dataset is 143.03.	×	0.01
The proposed framework transforms microblog posts into language-independent vectors by semantically aligning source and	×	0.06
The framework represents conversation propagation threads as an undirected topology to allow full-duplex interactions be	×	0.06
The model utilizes a multi-scale Graph Convolutional mechanism to aggregate features from claim semantics and event stru	×	0.14
The study proposes a domain-adaptive contrastive learning paradigm to minimize domain discrepancy.	×	0.12

## References

- <http://arxiv.org/abs/2002.12612v2>
- <http://arxiv.org/abs/2204.08143v2>
- <http://arxiv.org/abs/2304.01492v5>