



OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

When retrieval outperforms generation: Dense evidence retrieval for scalable fake news detection

Alamgir Qazi

Supervisors: Dr. Jamal Abdul Nasir, Dr. John McCrae



University
ofGalway.ie

Agenda

1. Problem Statement & Background
2. Current Approaches & Limitations
3. DeReC Architecture & Technical Implementation
4. Experimental Setup
5. Results & Analysis
6. Future Work





OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

1. Problem Statement & Background

Challenges of Misinformation

- Fake news and misinformation are major issues, impacting society by spreading false information that can influence public opinion and decisions.
- Automated fact verification systems are crucial to combat this, but current methods using explanation generating autoregressive Large Language Models (LLMs) face challenges.
- These include high computational costs, making real-time checking difficult, and risks of hallucinations—where LLMs generate incorrect or misleading information.



Research Objective

Research Question: Can dense retrieval-based classification systems achieve comparable or superior verification performance compared to explanation-generating LLM approaches for fact verification?

- Develop a computationally efficient, evidence-based fact verification system.
- Address limitations of current state-of-the-art approaches and provide evidence-grounded results without expensive LLM generation.
- Demonstrate that specialized models can outperform general-purpose LLMs in targeted tasks





OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

2. Current Approaches & Limitations

2. Current Approaches & Limitations

- State-of-the-art systems rely on autoregressive LLMs (7B+ parameters) to generate explanations.
- LLM-based methods, such as L-Defense and FactLLaMA, use models to generate explanations, aiming for detailed reasoning.
- 7B+ parameters (open-weight models like Llama3.1:8B, gemma3:27B, deepseek-r1:8B)
- 1.76T parameters (GPT-4) [rumoured]
- 5T parameters (GPT-4.5) [rumoured]



Key Limitations

- **Computational Overhead:** LLM inference is too slow for real-time fact checking.
- **Hallucination Risk:** Generated explanations frequently contain factual inconsistencies.
- **Resource Requirements:** 7B+ parameter models require significant GPU memory.





OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

3. DeReC Architecture & Technical Implementation

3. DeReC Three-Stage Architecture

DeReC: a light-weight Dense-Retrieval-Classification framework

1. Evidence Extraction

2. Evidence Retrieval

3. Veracity Prediction

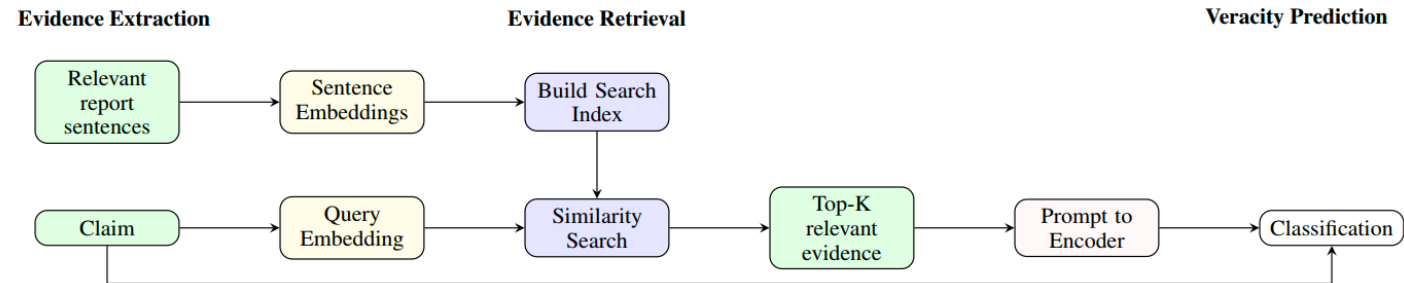


Figure 1: DeReC: Three-Stage Pipeline for Evidence-Based Fact Verification.



3. DeReC Three-Stage Architecture

1. Evidence Extraction (DeReC-qwen / DeReC-nomic)

- Dense embedding generation using gte-Qwen2-1.5B-instruct (1.5B parameters) or nomic-embed-text-v1.5 (137M parameters)
- Semantic representation learning



3. DeReC Three-Stage Architecture

2. Evidence Retrieval:

FAISS Configuration: IndexFlatIP for exact search

- Optimized inner product index for cosine similarity

Search Optimization:

- Top-K retrieval ($k=10$)



3. DeReC Three-Stage Architecture

3. Veracity Classification:

Model: DeBERTa-v3-large

- 304M parameters
- Evidence-enhanced input sequences



Datasets

Metric	RAWFC			LIAR-RAW		
	Train	Val	Test	Train	Val	Test
Number of Claims	1,612	200	200	10,065	1,274	1,251
Number of Reports	33,862	4,127	4,278	114,721	18,243	21,408
Total Sentences	248,343	31,191	31,453	626,573	102,147	118,449
Avg Sentences/Claim	154.06	155.96	157.26	62.25	80.18	94.68

Table 1: Analysis of dataset splits across LIAR-RAW and RAWFC datasets.

Baselines

- **Traditional Models:** dEFEND, SentHAN, SBERT-FC, CofCED, GenFE
- **LLM-based Models:** FactLLaMA, L-Defense (ChatGPT & LLaMA2)



Dataset Distribution

Veracity Label	RAWFC	LIAR-RAW
pants-fire	-	1,013
false	646	2,466
barely-true	-	2,057
half-true	671	2,594
mostly-true	-	2,439
true	695	2,021
Total Claims	2,012	12,590
Veracity Labels	3	6

Table 2: Distribution of veracity labels across RAWFC and LIAR-RAW datasets.





OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

Results & Analysis

DeReC: Results & Analysis

DeReC-qwen & DeReC-nomic achieve state-of-the-art results across both RAWFC and LIAR-RAW datasets.

	RAWFC			LIAR-RAW		
	P	R	F1	P	R	F1
<i>Traditional approach</i>						
dFEND (Shu et al., 2019)	44.90	43.20	44.00	23.00	18.50	20.50
SentHAN (Ma et al., 2019)	45.70	45.50	45.60	22.60	20.00	21.20
SBERT-FC (Kotonya and Toni, 2020a,b)	51.10	46.00	48.40	24.10	22.10	23.10
CofCED (Yang et al., 2022)	53.00	51.00	52.00	29.50	29.60	29.50
GenFE (Atanasova, 2024)	44.29	44.74	44.43	28.01	26.16	26.49
GenFE-MT (Atanasova, 2024)	45.64	45.27	45.08	18.55	19.90	15.15
<i>LLM-based approach</i>						
FactLLaMA (Cheung and Lam, 2023)	53.76	54.00	53.76	29.98	31.57	32.32
FactLLaMA _{know} (Cheung and Lam, 2023)	55.65	55.50	56.11	30.44	32.05	32.46
L-Defense _{ChatGPT} (Wang et al., 2024a)	61.72	61.91	61.20	30.55	<u>32.20</u>	30.53
L-Defense _{LLaMA2} (Wang et al., 2024a)	60.95	60.00	60.12	31.63	31.71	31.40
<i>Ours</i>						
DeReC-qwen	65.58	<u>64.56</u>	<u>64.60</u>	35.94	32.24	33.13
DeReC-nomic	<u>64.48</u>	65.57	64.61	<u>33.19</u>	31.50	<u>31.79</u>

Table 3: Performance comparison across RAWFC and LIAR-RAW datasets using Deberta-v3-large as classifier. Best scores are in **bold** and second-best scores are underlined for each metric.



Computational Efficiency

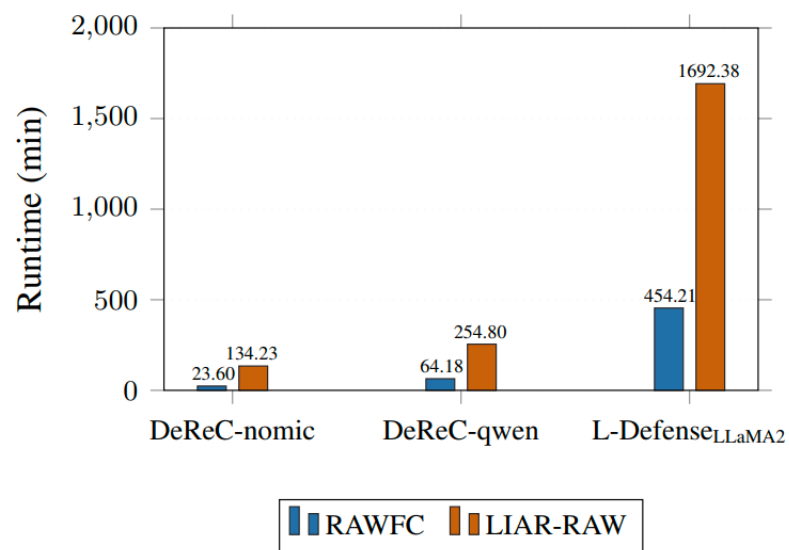


Figure 2: Complete pipeline runtime comparison (in minutes) on RAWFC and LIAR-RAW datasets.



DeReC: Results & Analysis

Dataset	Step	DeReC-nomic	DeReC-qwen	L-Defense _{LLaMA2}
RAWFC	Evidence Extraction	3m 50s	35m 15s	61m 39s
	Evidence Retrieval	2m 2s	7m 26s	-
	LLM-generated Explanations	-	-	381m 31s
	Veracity Prediction	17m 44s	21m 30s	11m 2s
	Total Runtime	23m 36s	64m 11s	454m 12s
LIAR-RAW	Evidence Extraction	9m 17s	89m 21s	185m 59s
	Evidence Retrieval	30m 12s	45m 13s	-
	LLM-generated Explanations	-	-	1466m 8s
	Veracity Prediction	94m 45s	89m 53s	40m 16s
	Total Runtime	134m 14s	254m 48s	1692m 23s

Table 4: Step-wise runtime breakdown (in minutes and seconds) for different models.



DeReC: Results & Analysis

Memory Requirements:

DeReC-nomic: 0.5GB (137M parameters)

DeReC-qwen: 6GB (1.5B parameters)

L-Defense: 28GB+ (7B+ parameters)





OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

Future Work

Future Work

- **Lightweight Explanation Generation:** Adding interpretability without sacrificing efficiency.
- **State Space Models:** Work on State Space Models and Mamba and try to include them in my pipeline to improve efficiency.





OLLSCOIL NA GAILLIMHE
UNIVERSITY OF GALWAY

Questions ?

University
ofGalway.ie