# Face the Facts!
# Evaluating RAG-based Fact-checking Pipelines in Realistic Settings

**Daniel Russo**[1,2], **Stefano Menini**[1], **Jacopo Staiano**[2], **Marco Guerini**[1]

[1]Fondazione Bruno Kessler, Italy
[2]University of Trento, Italy

`{drusso, menini, guerini}@fbk.eu, jacopo.staiano@unitn.it`

## Abstract

Natural Language Processing and Generation systems have recently shown the potential to complement and streamline the costly and time-consuming job of professional fact-checkers. In this work, we lift several constraints of current state-of-the-art pipelines for automated fact-checking based on the Retrieval-Augmented Generation (RAG) paradigm. Our goal is to benchmark, under more realistic scenarios, RAG-based methods for the generation of verdicts - i.e., short texts discussing the veracity of a claim - evaluating them on stylistically complex claims and heterogeneous, yet reliable, knowledge bases. Our findings show a complex landscape, where, for example, LLM-based retrievers outperform other retrieval techniques, though they still struggle with heterogeneous knowledge bases; larger models excel in verdict faithfulness, while smaller models provide better context adherence, with human evaluations favouring zero-shot and one-shot approaches for informativeness, and fine-tuned models for emotional alignment.

## 1 Introduction

Despite the efforts to validate the accuracy of online content, professional fact-checkers are increasingly struggling to keep up with the rapid spread of misinformation (Lewis et al., 2008; Adair et al., 2017; Godler and Reich, 2017; Wang et al., 2018). Therefore, Natural Language Processing (NLP) has been proposed as a viable solution to partially automate the costly process of verifying misleading claims online (Vlachos and Riedel, 2014).

Within this context, the task of *verdict production*, i.e. explaining why a claim is true or false, stands as one of the most challenging (Kotonya and Toni, 2020a; Guo et al., 2022).

Framing verdict production as a summarization task over fact-checking articles is a suitable solution due to the possibility of generating highly readable verdicts even for non-expert users (Atanasova, 2024; Kotonya and Toni, 2020b; Russo et al., 2023b). Despite their promising results, summarization-based approaches suffer from two main limitations: *(i)* they rely on the assumption that a fact-checking article always exists for a given claim; and *(ii)*, they further assume that claims are already paired with a fact-checking article, which is typical in fact-checking websites but not on social media platforms, where most of the misinformation spreads (Lazer et al., 2018).

Grounding textual generation on retrieved evidence, an approach named Retrieval-Augmented Generation (RAG), has been shown effective for knowledge-intensive tasks like fact verification (Lewis et al., 2020; Liao et al., 2023; Xu et al., 2023); it also allows to address the limitations of previous summarization approaches, e.g. the assumption that the claims are already paired with gold fact-checking articles. Moreover, RAG-based approaches have proven useful in reducing potential factual inconsistencies, often referred to as *"hallucinations"* (Zellers et al., 2019; Solaiman et al., 2019), during text generation (Lewis et al., 2020), making them attractive for fact-checking tasks. Thus, researchers have increasingly adopted RAG to enhance the accuracy of the generated verdicts (Zeng and Gao, 2024; Yao et al., 2023).

Current studies depend on fact-checking websites, resulting in verdicts that are marked by formal and dry language. This style contrasts sharply with the language used on Social Media Platforms (SMPs), which is typically more complex and includes *noise* such as additional text, personal commentary, or emotional content that surrounds the core fact. The mismatch between a social claim and the journalistic response might be highly ineffective when countering misinformation online (Colliander, 2019).

This work challenges and refines common assumptions in previous studies on verdict generation
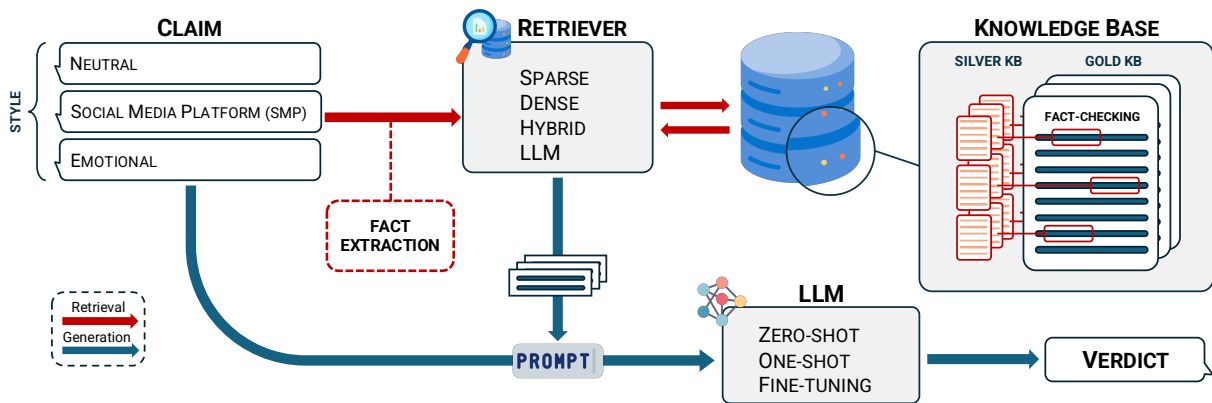
Figure 1: Visual representation of our RAG-based experimental design (the steps for retrieval and generation are indicated by the red and blue lines, respectively). We explored various configurations to tackle increasingly realistic scenarios across different claim styles (neutral, SMP, emotional) and Knowledge Bases (Gold vs. Silver), as well as varying computational demands through multiple retriever architectures (sparse, dense, hybrid, and LLM-based) and five distinct LLMs generation setups (zero-shot, one-shot, fine-tuning).

by testing various RAG-based pipelines in increasingly realistic scenarios and with different levels of computational complexity.

To this end, we present a thorough evaluation of verdict production along the several dimensions that characterize a RAG pipeline. For each dimension, we tested the main available options (see Figure 1). In particular, we tested:

- Three styles of claims that differ in their level of realism: *neutral* for journalistic claims, *SMP* for claims that resemble social media platform styles, and *emotional* whose SMP style is enriched with emotional elements;

- Four different retrieval approaches with varying computational requirements (*sparse*, *dense*, *hybrid*, *LLM-based*);

- Two claim pre-processing configurations to better address noise, namely with and without fact extraction;

- Five LLMs of varying sizes and training characteristics;

- Three LLM setups: zero-shot, few-shot, and fine-tuned models;

- Two types of knowledge bases (KB): a gold KB and a silver KB, with or without gold fact-checking articles included, respectively;

- Two document storage strategies in the KB, requiring the retriever to retrieve either entire articles or smaller segments (chunks).

Our findings show that LLM-based retrievers consistently outperform other methods, though they face challenges with silver knowledge bases. Dense retrievers manage stylistic variations of the claim effectively but fall short compared to LLMs, whereas sparse retrievers exhibit high sensitivity to noise present in emotional and SMP claims. Hybrid approaches and query pre-processing improved performance. Turning to generation, larger models excelled in faithfulness and alignment with gold verdicts, while smaller ones were more consistent in context adherence. Fine-tuning boosted verdict similarity but reduced contextual accuracy, with human evaluations favouring verdicts generated under zero/one -shot strategies for informativeness and fine-tuned models for emotional alignment.[1]

## 2 Related Work

Early approaches to verdict generation leveraged either attention modules to highlight salient tokens from the evidence text (Popat et al., 2018; Shu et al., 2019; Lu and Li, 2020; Wu et al., 2020; Yang et al., 2019), or Horn Rules to reason upon structured knowledge bases (Gad-Elrab et al., 2019; Ahmadi et al., 2019). However, both lack readability, being hard to interpret by common users (Guo et al., 2022). To overcome this issue, researchers started casting verdict production as a summarization task over fact-checking articles, either through extractive (Atanasova et al., 2020) or abstractive (Kotonya and Toni, 2020a; Stammbach and Ash, 2020) summarization, or using a combination of the

---

[1] Code and data are publicly available on GitHub: https://github.com/drusso98/face-the-facts.

two approaches in a hybrid summarization pipeline (Russo et al., 2023b). More recently, He et al. (2023) introduced a reinforcement learning-based framework for generating counter-misinformation responses to social media content.

Ad-hoc data collection strategies for verdict generation rely on synthetic data generation, like e-FEVER (Stammbach and Ash, 2020), and journalistic sources, such as LIARPLUS (Alhindi et al., 2018), PUBHEALTH (Kotonya and Toni, 2020b), LIAR++, and FullFact (Russo et al., 2023b). For more realistic, SMP-style claims, He et al. (2023) developed MisinfoCorrect, and Russo et al. (2023a) extended FullFact with VerMouth, incorporating emotional claims and verdicts grounded in trustworthy fact-checking articles.

While the latest approaches provide readable verdicts, they may lack faithfulness due to language models generating factual inaccuracies. Additionally, summarization approaches assume that trustworthy evidence is always available for the claim under inspection. Therefore, RAG-based approaches have been employed to guide the generation of reliable verdicts upon trustworthy evidence previously retrieved from a knowledge base (KB). To this end, Zeng and Gao (2024) proposed JustiLM, a few-shot RAG-based approach for the generation of verdicts for real-world claims, by leveraging both fact-checking articles and auxiliary evidence during model training. Yao et al. (2023) developed an end-to-end RAG-based system to jointly perform verdict prediction and production in a multimodal setting. Nevertheless, both studies concentrate on journalistic data and writing styles, without considering the communication style employed on SMP.

Building on the work of Zeng and Gao (2024), we present an extensive evaluation of RAG pipelines for verdict generation, testing various combinations of retrieval and generation strategies. Progressing toward increasingly realistic scenarios, we address the challenge of handling claims and sources that vary in style and complexity. Specifically, we assess the impact of claim writing style on retriever performance, highlighting where differences arise as the style shifts toward that used by SMP users. Finally, we explore the extreme scenario where no fact-checking article exists, relying solely on reliable supporting evidence.

## 3 Experimental Design

In this section, we provide details on the experimental design: from the datasets used, through the retrieval methods adopted, to the configurations of the LLMs employed for verdict generation.

### 3.1 Dataset

To study the impact of different styles on a RAG-based verdict production task, we used FullFact (Russo et al., 2023b) and VerMouth (Russo et al., 2023a) datasets. The two datasets comprise eight different versions of the same claims and verdicts. FullFact provides data written in a journalistic style scraped from fullfact.org while VerMouth proposes the same data rewritten in a SMP style, and also enriched with the six emotional components defined by Ekman (1992).

In both datasets, each claim-verdict pair is linked to a human-written fact-checking article, thus compounding to 8 different versions of the same claim: journalistic style (*neutral* hereafter), *SMP* style, anger, surprise, disgust, joy, fear, and sadness. In VerMouth, verdicts were also rewritten to reflect the various styles and emotions present in the claims. Throughout the paper, we will refer to emotion-styled subsets as *emotional* data.[2]

### 3.2 Retrieval Module

This comprises three elements: a *query* (a claim in our case), a *knowledge base* (KB), and a *retriever*.

**Claim** We used claims from FullFact and VerMouth datasets as queries. The two datasets offer three aligned variations of a claim: neutral, SMP, and emotional. Due to *noise* in SMP and emotional data, i.e., irrelevant information surrounding the main facts, directly using claims as queries can negatively impact the retrievers' performance. *Query rewriting*, which transforms context-dependent user queries into self-contained ones, has proven to be an effective approach for enhancing retriever performance (Elgohary et al., 2019; Ye et al., 2023). For this reason, we implemented a **fact extraction module** to simplify claims and remove noise around the main fact we need to retrieve evidence for. In particular, we employed `Llama-2-13b-chat-hf`, in a one-shot learning setup,[3] to extract the main facts from all SMP and emotional claims. An example of an emotional claim and its related fact is provided below:

---

[2]More details on the datasets are provided in Appendix A.
[3]The full instruction prompt is reported in Appendix B.

*Unbelievable! Just heard that 53 people have lost their lives in Gibraltar within 10 days of receiving Pfizer's Covid-19 vaccine. This is beyond alarming and I am absolutely furious. How can we trust these vaccines when they're causing more harm than good?! #PfizerVaccine #COVID19*

*53 people have lost their lives in Gibraltar within 10 days of receiving Pfizer's Covid-19 vaccine.*

**Retriever** We evaluated several retrieval strategies, with varying computational demands to accommodate the potential computational constraints of the target users: *(i)* sparse: BM25 and BM25+ (Robertson et al., 1995), a popular and effective extension of tf-idf; *(ii)* dense: Dragon+ (Lin et al., 2023) and Contriever (Izacard et al., 2021); *(iii)* hybrid, combining BM25+ and Dragon+ retrievers, using BAAI/bge-reranker-large as a reranker (Xiao et al., 2023); and, *(iv)* an instruction-tuned LLM for text embedding, e5-mistral-7b-instruct (Wang et al., 2023),[4] LLM-Retriever hereafter.

**Knowledge Base** To build the KB, we employed articles from the FullFact dataset, aligned with Ver-Mouth data. We named this KB as **Gold KB**. We experimented with two approaches: *(i)* indexing entire articles (Gold_KB$_{art}$); *(ii)* indexing small portions of text, *chunks*,[5] of each article as separate documents (Gold_KB$_{chunks}$).

In a realistic scenario, an up-to-date KB of fact-checking articles may not be available, or a fact-checking article might not exist (yet) for a given claim. To approximate this scenario, we leveraged knowledge from reliable sources to build a **Silver KB**. Specifically, we discarded gold fact-checking articles and extracted the evidence used to write and fact-check claims from FullFact's articles. This was collected by following the URLs present in the articles and getting their textual content. FullFact articles also typically link to the sources of the claims; however, the reliability of these sources is questionable, and thus we filtered them out. Also, we ignored all links to social networks (Twitter, Facebook, Instagram, TikTok, and Reddit).

Finally, from the remaining URLs, we extracted the text using the Newspaper3k[6] Python library.[7]

### 3.3 Verdict Generation

For the generation of the verdicts, we tested five LLMs, selected based on differences in sizes or the presence of guardrails: Mistral, in its v1.0 and v2.0 versions (Jiang et al., 2023); Llama-2 (Touvron et al., 2023), in its 7B and 13B chat versions;[8] and Llama-3-8B-Instruct.[9] We combined the claim and the retrieved evidence to prompt the LLM (see Appendix E.1), and tested generation under different setups, namely zero-shot, one-shot, and fine-tuning. For fine-tuning, we employed Llama-2-13b, the best-performing model in zero-shot and one-shot settings.

## 4 Retrieval Experiments

We tested the retrievers on FullFact and VerMouth test sets with an increasing number of retrieved documents ($k = 1, \ldots, 10$). For each claim used as a query, we considered as relevant documents the fact-checking article, or its chunks, linked to the claim. For space reasons, results on the emotional datasets will be presented in aggregated form, referred to as the *'emotional'* set. Retrieval experiments were carried out integrating into the LlamaIndex (Liu, 2022) framework either the Rank-BM25 library (Brown, 2020) or HugginFace's models (Wolf et al., 2020); retrieval performance was assessed with ranx (Bassani, 2022).

Table 1 presents retrieval results for each retrieval approach (sparse, dense, hybrid, LLM-Retriever) across all claim's styles (neutral, SMP, emotional) and fact-extraction pre-processings (SMP$_{facts}$, emotional$_{facts}$) using both KB configurations (Gold_KB$_{art}$ and Gold_KB$_{chunks}$). For Gold_KB$_{art}$, we report hit_rate@1 and MRR@1 (Mean Reciprocal Rank), as each claim had only one gold related article. For Gold_KB$_{chunks}$, we report hit_rate@10 and map@10 (Mean Average Precision) to assess whether the retrievers could consistently include at least one gold chunk among the top 10 and the precision of the retrieval system across different recall levels.[10]

---

[4]https://hf.co//intfloat/e5-mistral-7b-instruct

[5]We used the LlamaIndex (Liu, 2022) sentence splitter, which minimizes text fragmentation by keeping sentences integrity, with a maximum chunk token size of 100.

[6]https://github.com/codelucas/newspaper/

[7]Statistics related to the extra evidence collection are presented in Appendix C.

[8]https://hf.co/meta-llama/Llama-2-7b-chat-hf; https://hf.co/meta-llama/Llama-2-13b-chat-hf

[9]https://hf.co/meta-llama/Meta-Llama-3-8B-Instruct

[10]More details in Appendix D.

|  |  | sparse | dense | hybrid | LLM-Retriever | sparse | dense | hybrid | LLM-Retriever |
|---|---|---|---|---|---|---|---|---|---|
|  |  | hit_rate@1 | | | | mrr@10 | | | |
| **Articles** | Neutral | 0.903 | 0.905 | **0.966** | <u>0.960</u> | 0.931 | 0.938 | <u>0.972</u> | **0.978** |
|  | SMP | 0.770 | <u>0.799</u> | **0.937** | **0.937** | 0.817 | 0.866 | **0.963** | <u>0.962</u> |
|  | Emotional | 0.778 | 0.839 | <u>0.905</u> | **0.938** | 0.838 | 0.866 | <u>0.933</u> | **0.964** |
|  | $\text{SMP}_{facts}$ | 0.778 | 0.801 | **0.937** | <u>0.914</u> | 0.837 | 0.891 | **0.963** | <u>0.947</u> |
|  | $\text{Emotional}_{facts}$ | 0.835 | 0.846 | <u>0.905</u> | **0.932** | 0.883 | 0.897 | <u>0.933</u> | **0.958** |
|  |  | hit_rate@10 | | | | map@10 | | | |
| **Chunks** | Neutral | 0.963 | <u>0.992</u> | **1.000** | **1.000** | 0.392 | 0.552 | <u>0.573</u> | **0.655** |
|  | SMP | 0.856 | 0.974 | <u>0.994</u> | **1.000** | 0.275 | 0.484 | <u>0.536</u> | **0.619** |
|  | Emotional | 0.904 | 0.977 | <u>0.994</u> | **1.000** | 0.273 | 0.482 | <u>0.545</u> | **0.599** |
|  | $\text{SMP}_{facts}$ | 0.905 | 0.972 | <u>0.994</u> | **1.000** | 0.304 | 0.505 | <u>0.526</u> | **0.601** |
|  | $\text{Emotional}_{facts}$ | 0.939 | 0.978 | <u>0.994</u> | **0.999** | 0.345 | 0.518 | <u>0.552</u> | **0.615** |

Table 1: Results for the retrieval experiments. We report hit_rate, mrr, and map for retrieval over the $\text{Gold\_KB}_{art}$ (Articles) and the $\text{Gold\_KB}_{chunks}$ (Chunks) KBs. $\text{SMP}_{facts}$ and $\text{Emotional}_{facts}$ indicate input preprocessing with the fact extraction module. The first and second best results are in bold and underlined, respectively.

For article retrieval, all four retrievers achieved high accuracy on neutral claims, with a hit_rate@1 above 90%. When the correct article was not immediately retrieved, they still ranked it highly, as shown by strong MRR@10 scores. Performance declined with noisier claims, especially for sparse retrievers, while dense models and the LLM-Retriever showed greater robustness. Hybrid retrieval (combining sparse and dense) performed comparably to the LLM-Retriever.

In chunk retrieval, the LLM-Reranker excelled, achieving a hit_rate@10 that always included a gold chunk and reaching an average MAP@10 of 70%. Sparse retrievers showed low map scores (∼40% for neutral claims), particularly for SMP and emotional claims (<30%), while dense and hybrid approaches followed trends similar to article retrieval. Thus, the low MAP scores indicate that sparse retrievers must retrieve more chunks from the knowledge base to select the relevant content. However, this comes at the cost of retrieving more non-relevant content, which could potentially compromise the subsequent generation phase.

Overall, the LLM-Retriever consistently outperforms other approaches. Notably, it remains stable even when exposed to different input claim styles, exhibiting minimal performance degradation when noise is introduced. Similarly, dense retrievers exhibit lower performance but demonstrate robustness to stylistic variations. In contrast, sparse retrievers are significantly affected by data noise, resulting in performance drops across all three datasets, while the fact extraction module we included yields consistent performance gains, particularly when using sparse and dense retrievers.

**Retrieval from Silver KB** Turning to the Silver KB setup, we tested the optimal retriever methodology from the previous experiments, specifically the LLM-Retriever. As outlined in Section 3.2, the Silver KB consists of reliable sources that have been extracted from the initial fact-checking articles. The evidence compounded 9983 chunks, each corresponding to a fact-checking article considered a gold standard during evaluation. The results, illustrated in Table 2, indicate that modifying the knowledge base strongly impacted retrieval performance across all three datasets. Unlike the previous setting with the Gold KB, the LLM's performance in this context is markedly influenced by the stylistic nature of the claims: neutral claims yielded better results.

Prepending a fact extraction module significantly improves retriever results, implying that even robust retrievers can benefit from preprocessing when dealing with heterogeneous KB (that do not contain gold fact-checking articles) and realistic claims, such as those from SMPs.

| Neutral | SMP | Emotional | $\text{SMP}_{facts}$ | $\text{Emotional}_{facts}$ |
|---|---|---|---|---|
| 0.683 | 0.652 | 0.637 | 0.689 | 0.671 |

Table 2: Hit_rate@10 scores for retrieval with LLM-Retrieval over the Silver KB.

# 5 Generation Experiments

For verdict generation, the claim and its corresponding retrieved evidence were combined into a prompt fed to the five different LLMs (see Section 3.3).

| | | Articles | | | | Chunks | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | ROUGE-LSum | BARTScore | SummaC | GoldSim | ROUGE-LSum | BARTScore | SummaC | GoldSim |
| **zero shot** | Neutral | **0,16** | **-2,13** | **0,35** | -3,03 | **0,25** | **-2,61** | **0,25** | -3,06 |
| | SMP | 0,15 | -2,31 | 0,32 | -3,01 | 0,24 | -2,69 | 0,24 | -3,03 |
| | emotional | 0,14 | -2,48 | 0,31 | **-2,94** | 0,22 | -2,79 | 0,23 | **-2,98** |
| **one shot** | Neutral | **0,16** | **-2,13** | **0,33** | -3,03 | **0,26** | **-2,53** | 0,24 | -2,99 |
| | SMP | 0,15 | -2,42 | 0,32 | -3,01 | 0,24 | -2,73 | **0,23** | -2,95 |
| | emotional | 0,14 | -2,60 | 0,32 | **-2,95** | 0,22 | -2,85 | **0,23** | **-2,91** |
| **fine tuning** | Neutral | **0,10** | **-1,45** | **0,53** | -2,71 | 0,08 | **-1,42** | **0,52** | -2,75 |
| | SMP | **0,10** | -2,30 | 0,33 | -2,58 | **0,10** | -2,45 | 0,32 | **-2,63** |
| | emotional | **0,10** | -2,43 | 0,31 | **-2,48** | **0,10** | -2,76 | 0,31 | -2,68 |

Table 3: Generation results per dataset, *averaged across the LLMs*. Retrieved articles or chunks were employed in the generation. The best results for each generation configuration are in bold.

| | | Articles | | | | Chunks | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | ROUGE-LSum | BARTScore | SummaC | GoldSim | ROUGE-LSum | BARTScore | SummaC | GoldSim |
| **zero-shot** | mistral-v0.1 | **0,19** | -2,16 | **0,35** | -3,02 | **0,19** | -2,20 | **0,35** | -3,05 |
| | mistral-v0.2 | 0,13 | -2,55 | 0,32 | -3,12 | 0,15 | -2,50 | 0,33 | -3,11 |
| | llama3-8b | 0,05 | -3,21 | 0,30 | -3,54 | 0,04 | -3,37 | 0,32 | -3,65 |
| | llama2-7b | 0,16 | -2,12 | 0,32 | -2,84 | **0,19** | -2,06 | 0,33 | -2,87 |
| | llama2-13b | 0,17 | **-1,89** | 0,34 | **-2,73** | **0,19** | **-1,89** | **0,35** | **-2,77** |
| **one-shot** | mistral-v0.1 | **0,16** | -2,38 | **0,33** | -3,02 | 0,17 | -2,47 | 0,32 | -3,07 |
| | mistral-v0.2 | 0,12 | -2,61 | 0,31 | -3,16 | 0,13 | -2,63 | 0,36 | -3,16 |
| | llama3-8b | 0,13 | **-2,25** | **0,33** | -2,98 | 0,14 | -2,26 | **0,33** | -2,97 |
| | llama2-7b | **0,16** | -2,35 | 0,32 | -2,93 | **0,19** | -2,34 | 0,31 | -2,89 |
| | llama2-13b | **0,16** | -2,31 | 0,32 | **-2,90** | 0,17 | **-2,22** | 0,31 | **-2,79** |
| **ft** | llama2-13b | 0,10 | -2,08 | 0,39 | -2,59 | 0,10 | -2,21 | 0,38 | -2,69 |

Table 4: Generation results per LLM, *averaged across the three datasets* (neutral, SMP, emotional). Retrieved articles or chunks were employed in the generation. The best results for each generation configuration are in bold.

For evidence retrieval, we employed the best-performing retriever (i.e. LLM-Retriever). We tested the five LLMs under three setups: zero-shot, one-shot, and fine-tuned.[11] For fine-tuning, we employed the best-performing model in both the zero-shot and one-shot configurations, namely Llama-2-13b.

When using the Gold_KB$_{art}$, we included the top-1 article (i.e. the most relevant) in the prompt, a choice justified by the LLM retriever's remarkable hit_rate@1 results. Conversely, with retrieval from Gold_KB$_{chunsk}$, we fed the model with 10 retrieved chunks, based on its map@10 performance (see Table 1 and Figure D).[12]

**Automatic Metrics** We use ROUGE-LSum (Lin, 2004), BARTScore Yuan et al. (2021), and SummaC (Laban et al., 2022) to evaluate lexical adherence and faithfulness of the generated text to the context provided to the LLM. Further, we used BARTScore, which is unaffected by differences in text length, to compute the semantic similarity (*GoldSim*) between the generated and gold verdicts from FullFact and VerMouth.

For average performances per dataset and per model see Tables 3 and 4, respectively.[13]

**Zero-shot and One-shot** Generations with neutral claims yielded better results compared to the SMP and emotional data in both zero-shot and one-shot experiments, indicating that the complexity of claims affects not only the retrieval phase but also the generation phase. Interestingly, when comparing the similarity between the generated and the gold verdicts, the generations with emotional data produced the best results (Table 3). This can be due to recurrent patterns in the SMP and emotional data, such as expressions of empathy and politeness generated by ChatGPT ("I understand your frustration", "It is important to note that"). These patterns were also replicated by the models used in this study. Notably, zero-shot experiments generally outperformed one-shot experiments when results were averaged across all the data (Table 3) and LLMs (Table 4).

Turning to individual model performances (Table 4), larger models (Llama-2-13b) demonstrated higher faithfulness to the context and similarity to gold standard verdicts, whereas smaller

---

[11]The instructions used are provided in Appendix E.1.
[12]Examples of generations are provided in Appendix E.5.

[13]Complete results can be found in Appendix E.4, Table 9.

ones (`mistral-7b-v0.1`, `llama3-8b`) showed better contextual adherence in terms of overlaps (ROUGE-LSum) and consistency (SummaC).

To sum up, generations with neutral claims outperformed those with emotional data in both zero-shot and one-shot experiments but produced more accurate results when paired with emotional data. Larger models generally performed better than smaller ones in terms of faithfulness to context and consistency, while still showing promise in contextual adherence metrics.

**LLM Fine-Tuning** We fine-tuned the `Llama-2-13b` model, namely the best performing model in the previous in-context learning experiments, disjointly on the three claim styles. To this end, the model was fed with claims and gold verdicts from the FullFact and VerMouth datasets, complemented by positive and negative contextual information.[14] We randomly sampled 200 entries from both neutral and SMP training datasets; similarly, to obtain a comparably sized dataset for emotional data, we randomly sampled 35 examples for each of the 6 emotional dimensions available.[15]

Fine-tuning results are reported in Tables 3 and 4, showing that `LLaMA-2-13b` models improve in faithfulness and similarity with respect to gold verdicts across all three datasets and for both full articles and text chunks. However, this comes at the cost of lower ROUGE-L scores: without fine-tuning, the models tend to extract and replicate not only the necessary information but also the exact wording from the context. Thus, fine-tuning models abstract better from the context, as would be expected, and also show better performance in selecting relevant and reliable information. Further, after fine-tuning, the emotional models achieved higher similarity scores to the original claims. Akin to the zero and one-shot configurations, analysis of the generated verdicts reveals that these models produce empathetic expressions similar to those found in the VerMouth dataset.

**Generation with Silver KB** Finally, we tackled the more realistic scenario wherein the information useful to generate a verdict is spread across several documents: this allows us to lift the typical constraint of existing datasets wherein a claim is paired to a single article.

|  |  | ROUGE-LSum | BARTScore | SummaC | GoldSim |
|---|---|---|---|---|---|
| zero-shot | Neutral | 0,23 | -2,17 | 0,28 | -3,04 |
|  | SMP | 0,21 | -2,30 | 0,28 | -2,84 |
|  | Emotional | 0,19 | -2,50 | 0,26 | -2,73 |
|  | Mean | 0,21 | -2,32 | 0,28 | -2,87 |
| one-shot | Neutral | 0,24 | -2,39 | 0,24 | -3,02 |
|  | SMP | 0,21 | -2,57 | 0,25 | -2,90 |
|  | Emotional | 0,18 | -2,68 | 0,26 | -2,76 |
|  | Mean | 0,21 | -2,55 | 0,25 | -2,89 |
| fine tuned | Neutral | 0,12 | -2,18 | 0,39 | -3,00 |
|  | SMP | 0,13 | -2,64 | 0,26 | -2,59 |
|  | Emotional | 0,13 | -2,97 | 0,25 | -2,70 |
|  | Mean | 0,13 | -2,60 | 0,30 | -2,77 |

Table 5: Generation results with extra-evidence.

We used all the documents from the Silver KB (see Section 3.2) and tested with the LLM-Retriever and `Llama-2-13b` (i.e. the best-performing setup from the above experiments) for RAG-based verdict generation.

We focused solely on the chunk-based configuration as the information required to build a verdict is *(i)* often distributed across multiple extra documents,[16] and *(ii)* it is more likely to be located in specific sections of these extra evidence articles. The results are reported in Table 5. Except for ROUGE-LSum, `Llama-2-13b`'s performance is consistently slightly worse when compared with results from the generation using the Gold KB.[17] Still, a qualitative analysis showed that using the Silver KB resulted in verdicts that in most cases were consistent with the claim, faithful to the context, and informative.

The lower results can be explained by the substantial difference between the Gold and the Silver KBs: indeed, a gold fact-checking article refers to a single claim and contains all the information needed to generate the verdict. Therefore, when using the Gold KB, out of a total of ten chunks, a robust retriever – such as LLM-Retriever – can identify a larger number of informative chunks. Conversely, in the case of the Silver KB, for each claim, on average there exist four related articles (see Appendix C) that are most likely to provide partial information about the verdict. Hence, in a realistic scenario, the retriever must search in a collection of documents that is larger, highly sparse in terms of information, and potentially repetitive: among 10 retrieved chunks, it cannot be assumed that the gold chunks retrieved hold enough information to generate a good-quality verdict.

---

[14] Details on the creation of the training data are provided in Appendix E.2.

[15] Fine-tuning details are reported in Appendix E.3.

[16] On average, it is found in four documents, see Appendix C, Table 7.

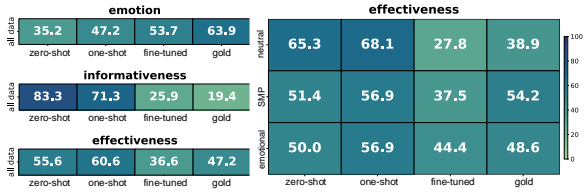[17] Compare with Tables 4, 3, and Appendix E.4 - Table 9.

Figure 2: Human evaluation results: Percentages of preference for the four generation setups across the datasets.

## 5.1 Human Evaluation

Automatic metrics for NLG evaluation are known to correlate poorly with human judgements. Several works showed how optimizing for such metrics (e.g. ROUGE) is largely suboptimal (Paulus et al., 2018; Scialom et al., 2019), and they suffer from weak interpretability and failure to capture nuances (Sai et al., 2022). Thus, we provide a comprehensive human evaluation of the generated verdicts.

We enlisted three expert evaluators,[18] and we focused on the data generated using the Gold KB, as this configuration yielded the best results in both the retrieval and generation phases (see Sections 4, 5). The evaluators were provided with pairs of verdicts (either gold or generated using zero-shot, one-shot, or fine-tuned models) and their corresponding claims. They were instructed to assess the best verdict based on three aspects: effectiveness, informativeness, and emotional/empathetic coverage. We sampled 72 claims equally distributed among the test sets, and provided the evaluators with six combinations of verdict pairs, compounding to 432 samples to be evaluated on the three evaluation dimensions; thus, the reported human evaluation is based on a total of 1296 data points. In Figure 2 (left) we report how many times, in percentage, the human annotators preferred each of the four verdict generation setups (gold, zero-shot, one-shot, fine-tuning) across the three claim styles (neutral, SMP, emotional).

Results show that zero-shot and one-shot approaches were largely preferred in terms of informativeness. Gold and fine-tuned configurations were considered the best in terms of emotional matching between the claim and the related verdicts, which can be explained by the fact that fine-tuned models learned to mimic the emotional style of the gold training data.

This is supported by the effectiveness evaluation in Figure 2 (right), where the delta between the preferences for zero/one-shot configurations and gold/fine-tuned ones is higher for neutral claims and it decreases as the emotional component increases. Nevertheless, on average, zero and one-shot configurations remain the most preferred. This finding can also be explained by the fact that the verdicts provided by FullFact are not designed to be used in isolation for social correction. They are intended to be read alongside the associated article. Therefore, they are not self-contained in terms of grounding, as the full context is available on the same page of the article.[19]

## 6 Conclusions

Fighting misinformation is a demanding endeavour aggravated over the latest years by the reach and scale provided by popular social networking platforms. With this work, we provided a comprehensive overview of how RAG-based approaches can effectively be employed for the automatic generation of verdicts when addressing realistic scenarios. We thus presented an extensive experimental analysis of the different approaches that can be adopted at each stage of such pipelines, designing scenarios which also account for the varying computational constraints that researchers, practitioners, and professional fact-checking organizations might have.

Our results show that LLM-based retrievers consistently demonstrate superior adaptability and performance in retrieval tasks, outperforming other methods. However, they face challenges when dealing with heterogeneous knowledge bases, where performance declines. Hybrid retrieval approaches offer a cost-effective alternative, while dense retrievers remain robust to stylistic variations. Notably, incorporating fact extraction modules enhances retrieval effectiveness across setups; thus query preprocessing is particularly beneficial for addressing realistic claims. Although claim complexity reduces generation accuracy, larger models achieve greater faithfulness to context and alignment with gold-standard verdicts. Human evaluations also favour zero-shot and one-shot approaches, particularly for their informativeness and emotional alignment.

---

[18]A senior researcher and two MSc graduates; all volunteer evaluators are proficient in English, experts in NLP, and knowledgeable about social media platforms' communication styles and dynamics, particularly in the context of misinformation.

[19]Further details are provided in Appendix F.

# 7 Limitations

The data employed are limited to the English language only. However, we believe that the proposed RAG-based pipeline can still be adapted to a multilingual setting by adopting appropriate/multilingual retrievers and LLMs.

Moreover, the datasets employed in the study are closely associated with FullFact topic domain (and style) and are confined to the SMP/emotional communication style adopted in VerMouth dataset. However, to the best of our knowledge, FullFact and VerMouth are the sole available aligned datasets that allow a thorough examination of the impact of the claim's communication style on a RAG-based verdict production task. Nonetheless, we are confident that the single modules employed in our work can be easily adapted to novel domains or SMP communication styles.

In this work, we focus on utilizing the extra evidence provided by fact-checking articles as a proxy for using the web as a knowledge base. This approach allows us to bypass the inherently complex challenges of determining source relevance and reliability, issues beyond the scope of this paper. Consequently, we strike a balance between leveraging external information and maintaining manageable complexity. Furthermore, our experiments reveal that even this constrained setting presents significant challenges. Therefore, we argue that mastering the use of extra evidence in this context is a crucial first step before moving towards more sophisticated and resource-intensive methods.

# 8 Ethical Statement

Our work is motivated by the potential to improve the accuracy and efficiency of automated fact-checking systems. However, we acknowledge that the development of such technologies can potentially, as any human artefact, be exploited by malicious actors. In this case, the technological building blocks (e.g. the LLMs) can be tuned to accomplish goals opposite to ours (e.g. generate persuasive fake news). We argue that, while malicious actors would keep pursuing their goals regardless of the community efforts, our work provides a contribution to keep up to their pace and fostering advancements by relying exclusively on publicly available data and models, and by publicly releasing novel artefacts (e.g. the fine-tuned `Llama2-13b` checkpoints).

# References

Bill Adair, Chengkai Li, Jun Yang, and Cong Yu. 2017. Progress toward "the holy grail": The continued quest to automate fact-checking. In *Computation+ Journalism Symposium,(September)*.

Naser Ahmadi, Joohyung Lee, Paolo Papotti, and Mohammed Saeed. 2019. Explainable fact checking with probabilistic answer set programming. In *Conference on Truth and Trust Online*.

Tariq Alhindi, Savvas Petridis, and Smaranda Muresan. 2018. Where is your evidence: Improving fact-checking by justification modeling. In *Proceedings of the First Workshop on Fact Extraction and VERification (FEVER)*, pages 85–90, Brussels, Belgium. Association for Computational Linguistics.

Pepa Atanasova. 2024. Generating fact checking explanations. In *Accountable and Explainable Methods for Complex Reasoning over Text*, pages 83–103. Springer.

Pepa Atanasova, Jakob Grue Simonsen, Christina Lioma, and Isabelle Augenstein. 2020. Generating fact checking explanations. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7352–7364, Online. Association for Computational Linguistics.

Elias Bassani. 2022. ranx: A blazing-fast python library for ranking evaluation and comparison. In *ECIR (2)*, volume 13186 of *Lecture Notes in Computer Science*, pages 259–264. Springer.

Dorian Brown. 2020. Rank-BM25: A Collection of BM25 Algorithms in Python.

Jonas Colliander. 2019. "this is fake news": Investigating the role of conformity to other users' views when commenting on and spreading disinformation in social media. *Computers in Human Behavior*, 97:202–215.

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2023. Qlora: Efficient finetuning of quantized llms. In *Advances in Neural Information Processing Systems*, volume 36, pages 10088–10115. Curran Associates, Inc.

Paul Ekman. 1992. An argument for basic emotions. *Cognition & Emotion*, 6:169–200.

Ahmed Elgohary, Denis Peskov, and Jordan Boyd-Graber. 2019. Can you unpack that? learning to rewrite questions-in-context. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5918–5924, Hong Kong, China. Association for Computational Linguistics.

Mohamed H. Gad-Elrab, Daria Stepanova, Jacopo Urbani, and Gerhard Weikum. 2019. Exfakt: A framework for explaining facts over knowledge graphs and text. In *Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining*, WSDM '19, page 87–95, New York, NY, USA. Association for Computing Machinery.

Yigal Godler and Zvi Reich. 2017. Journalistic evidence: Cross-verification as a constituent of mediated knowledge. *Journalism*, 18(5):558–574.

Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. 2022. A survey on automated fact-checking. *Transactions of the Association for Computational Linguistics*, 10:178–206.

Bing He, Mustaque Ahamad, and Srijan Kumar. 2023. Reinforcement learning-based counter-misinformation response generation: A case study of covid-19 vaccine misinformation. In *Proceedings of the ACM Web Conference 2023*, WWW '23, page 2698–2709, New York, NY, USA. Association for Computing Machinery.

Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. Unsupervised dense information retrieval with contrastive learning. *arXiv preprint arXiv:2112.09118*.

Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. Mistral 7b. *arXiv preprint arXiv:2310.06825*.

Neema Kotonya and Francesca Toni. 2020a. Explainable automated fact-checking: A survey. *arXiv preprint arXiv:2011.03870*.

Neema Kotonya and Francesca Toni. 2020b. Explainable automated fact-checking for public health claims. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7740–7754, Online. Association for Computational Linguistics.

Philippe Laban, Tobias Schnabel, Paul N. Bennett, and Marti A. Hearst. 2022. SummaC: Re-visiting NLI-based models for inconsistency detection in summarization. *Transactions of the Association for Computational Linguistics*, 10:163–177.

David M. J. Lazer, Matthew A. Baum, Yochai Benkler, Adam J. Berinsky, Kelly M. Greenhill, Filippo Menczer, Miriam J. Metzger, Brendan Nyhan, Gordon Pennycook, David Rothschild, Michael Schudson, Steven A. Sloman, Cass R. Sunstein, Emily A. Thorson, Duncan J. Watts, and Jonathan L. Zittrain. 2018. The science of fake news. *Science*, 359(6380):1094–1096.

Justin Matthew Wren Lewis, Andy Williams, Robert Arthur Franklin, James Thomas, and Nicholas Alexander Mosdell. 2008. The quality and independence of british journalism.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474. Curran Associates, Inc.

Hao Liao, Jiahao Peng, Zhanyi Huang, Wei Zhang, Guanghua Li, Kai Shu, and Xing Xie. 2023. Muser: A multi-step evidence retrieval enhancement framework for fake news detection. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, KDD '23, page 4461–4472, New York, NY, USA. Association for Computing Machinery.

Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Sheng-Chieh Lin, Akari Asai, Minghan Li, Barlas Oguz, Jimmy Lin, Yashar Mehdad, Wen-tau Yih, and Xilun Chen. 2023. How to train your dragon: Diverse augmentation towards generalizable dense retrieval. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 6385–6400, Singapore. Association for Computational Linguistics.

Jerry Liu. 2022. LlamaIndex.

Yi-Ju Lu and Cheng-Te Li. 2020. GCAN: Graph-aware co-attention networks for explainable fake news detection on social media. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 505–514, Online. Association for Computational Linguistics.

Romain Paulus, Caiming Xiong, and Richard Socher. 2018. A deep reinforced model for abstractive summarization. In *International Conference on Learning Representations*.

Kashyap Popat, Subhabrata Mukherjee, Andrew Yates, and Gerhard Weikum. 2018. Declare: Debunking fake news and false claims using evidence-aware deep learning. *arXiv preprint arXiv:1809.06416*.

Stephen E Robertson, Steve Walker, Susan Jones, Micheline M Hancock-Beaulieu, Mike Gatford, et al. 1995. Okapi at trec-3. *Nist Special Publication Sp*, 109:109.

Daniel Russo, Shane Kaszefski-Yaschuk, Jacopo Staiano, and Marco Guerini. 2023a. Countering misinformation via emotional response generation. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 11476–11492, Singapore. Association for Computational Linguistics.

Daniel Russo, Serra Sinem Tekiroğlu, and Marco Guerini. 2023b. Benchmarking the generation of fact checking explanations. *Transactions of the Association for Computational Linguistics*, 11:1250–1264.

Ananya B. Sai, Akash Kumar Mohankumar, and Mitesh M. Khapra. 2022. A survey of evaluation metrics used for nlg systems. *ACM Comput. Surv.*, 55(2).

Thomas Scialom, Sylvain Lamprier, Benjamin Piwowarski, and Jacopo Staiano. 2019. Answers unite! unsupervised metrics for reinforced summarization models. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3246–3256.

Kai Shu, Limeng Cui, Suhang Wang, Dongwon Lee, and Huan Liu. 2019. defend: Explainable fake news detection. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, KDD '19, page 395–405, New York, NY, USA. Association for Computing Machinery.

Irene Solaiman, Miles Brundage, Jack Clark, Amanda Askell, Ariel Herbert-Voss, Jeff Wu, Alec Radford, and Jasmine Wang. 2019. Release strategies and the social impacts of language models.

Dominik Stammbach and Elliott Ash. 2020. e-fever: Explanations and summaries forautomated fact checking. In *Conference for Truth and Trust Online*.

Serra Sinem Tekiroğlu, Yi-Ling Chung, and Marco Guerini. 2020. Generating counter narratives against online hate speech: Data and strategies. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1177–1190, Online. Association for Computational Linguistics.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Andreas Vlachos and Sebastian Riedel. 2014. Fact checking: Task definition and dataset construction. In *Proceedings of the ACL 2014 workshop on language technologies and computational social science*, pages 18–22.

Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2023. Improving text embeddings with large language models. *arXiv preprint arXiv:2401.00368*.

Patrick Wang, Rafael Angarita, and Ilaria Renna. 2018. Is this the era of misinformation yet: Combining social bots and fake news to deceive the masses. In *Companion Proceedings of the The Web Conference 2018*, WWW '18, page 1557–1561, Republic and Canton of Geneva, CHE. International World Wide Web Conferences Steering Committee.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. Huggingface's transformers: State-of-the-art natural language processing.

Lianwei Wu, Yuan Rao, Yongqiang Zhao, Hao Liang, and Ambreen Nazir. 2020. DTCA: Decision tree-based co-attention networks for explainable claim verification. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1024–1035, Online. Association for Computational Linguistics.

Shitao Xiao, Zheng Liu, Peitian Zhang, and Niklas Muennighoff. 2023. C-pack: Packaged resources to advance general chinese embedding.

Peng Xu, Wei Ping, Xianchao Wu, Lawrence McAfee, Chen Zhu, Zihan Liu, Sandeep Subramanian, Evelina Bakhturina, Mohammad Shoeybi, and Bryan Catanzaro. 2023. Retrieval meets long context large language models. *arXiv preprint arXiv:2310.03025*.

Fan Yang, Shiva K. Pentyala, Sina Mohseni, Mengnan Du, Hao Yuan, Rhema Linder, Eric D. Ragan, Shuiwang Ji, and Xia (Ben) Hu. 2019. Xfake: Explainable fake news detector with visualizations. In *The World Wide Web Conference*, WWW '19, page 3600–3604, New York, NY, USA. Association for Computing Machinery.

Barry Menglong Yao, Aditya Shah, Lichao Sun, Jin-Hee Cho, and Lifu Huang. 2023. End-to-end multimodal fact-checking and explanation generation: A challenging dataset and models. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '23, page 2733–2743, New York, NY, USA. Association for Computing Machinery.

Fanghua Ye, Meng Fang, Shenghui Li, and Emine Yilmaz. 2023. Enhancing conversational search: Large language model-aided informative query rewriting. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5985–6006, Singapore. Association for Computational Linguistics.

Weizhe Yuan, Graham Neubig, and Pengfei Liu. 2021. Bartscore: Evaluating generated text as text generation. In *Advances in Neural Information Processing Systems*, volume 34, pages 27263–27277. Curran Associates, Inc.

Rowan Zellers, Ari Holtzman, Hannah Rashkin, Yonatan Bisk, Ali Farhadi, Franziska Roesner, and Yejin Choi. 2019. Defending against neural fake news. *Advances in neural information processing systems*, 32.

Fengzhu Zeng and Wei Gao. 2024. JustiLM: Few-shot Justification Generation for Explainable Fact-Checking of Real-world Claims. *Transactions of the Association for Computational Linguistics*, 12:334–354.

## A   Dataset Details

In this work, we employed data from the FullFact (Russo et al., 2023b) and VerMouth (Russo et al., 2023a) datasets. The FullFact dataset consists of `claim-article-verdict` triplets extracted from the FullFact website[20]. The VerMouth dataset extends the FullFact dataset. Starting from FullFact's triplets, Russo et al. (2023a) leveraged an *author-reviewer pipeline* (Tekiroğlu et al., 2020) to rewrite the data in according to social media platform style, either in a general manner or with an embedded emotional component. For more comprehensive details on the datasets, we encourage readers to refer to the original papers (Russo et al., 2023b,a). In Table 6, we detail the distribution of entries across the training, evaluation, and test sets for each dataset, namely FullFact (Russo et al., 2023b) and VerMouth (Russo et al., 2023a).

|  |  | Train | Eval | Test |
|---|---|---|---|---|
| **FullFact** | | 1470 | 184 | 174 |
| **SMP** | | 1470 | 184 | 174 |
| | anger | 1265 | 158 | 158 |
| | disgust | 1339 | 164 | 163 |
| VerMouth / emotions | fear | 1440 | 179 | 171 |
| | happiness | 1200 | 165 | 149 |
| | sadness | 1404 | 173 | 171 |
| | surprise | 1433 | 181 | 170 |

Table 6: FullFact and VerMouth data distribution.

## B   Fact Extraction Module

In order to remove noise from VerMouth's claims (Russo et al., 2023a), we prepend a fact extraction module before passing the claim to the retriever. To this end, we prompted `Llama-2-13b` and we provided an example of the expected output (one-shot). Hereafter we report the prompt employed:

```
Extract from the following text the main fact.
    Remove possible opinions or emotional
    statement.
Report results in the following format: FACT:[
    main fact]

Here there is an example:
TEXT:"I just heard about the Covid-19 vaccines &
    sadly they don't seem to be very effective
    in preventing the virus. Really
    disappointing! #vaccineineffective #
    covid19vaccine"
FACT:"The Covid-19 vaccines offer very little
    protection against the disease."
```

---

[20]https://fullfact.org

```
Now extract the main fact from the following
    text:
TEXT:{claim}
```

## C   Extra Evidence Extraction Details

In Table 7, we present detailed information about
the additional evidence extracted from FullFact
fact-checking articles, used to approximate the real-
istic scenario in which a gold fact-checking article
is not available or does not exist (yet). From the
original FullFact fact-checking articles we removed
links to social networks and the source URL of the
claim. Indeed, the claims fact-checked by FullFact
vary in nature, often originating from social media
posts, images, videos, and sometimes misleading
headlines. Consequently, the source of the claim
might not always provide additional information
beyond the claim itself that can be used for veri-
fication. Furthermore, even if the claim's source
contains extra text, the information can potentially
be misleading. Therefore, following our "reliability
requirement" we filtered out the claim sources.

|      | extra art | extra | words | sent | chunks |
|------|-----------|-------|-------|------|--------|
| all  | 4093      | 4     | 970   | 38   | 69412  |
| test | 672       | 4     | 868   | 35   | 9983   |

Table 7: Statistics for all additional evidence extracted
from FullFact fact-checking articles and the test set used
in our experiments. We report the total number of extra
evidence documents (*extra art*); the average number of
extra documents per fact-checking article (*extra*); the
average number of words (*words*) and sentences (*sent*);
and the total number of chunks (*chunks*).

## D   Retrieval Experiments Details

In Figure 3 we report hit_rate, Mean Reciprocal
Rank (MRR), and Mean Average Precision (MAP)
for the retrieval experiments.

## E   Generation Experiments Details

### E.1   Model's instruction

Hereafter the instruction employed for the zero-
shot setting. A similar instruction was modified,
by adding an example from the training sets in the
one-shot configuration.

```
Based on the provided context, respond to the
    claim, ensuring a thorough explanation. Use
    only the given context. Reply in no more
    than three sentences. Avoid mentioning the
    context in the reply. Match the
    communication style of the claim and address
```

```
    the possible emotional component present in
    it, if needed. If the context is
    insufficient, state that you don't know.
    Format your response as follows:

Reply: [your_reply]

The context information is provided below (in
    between xml tags).

<context>

{context_str}

</context>

Claim: "{query_str}"
```

### E.2   Training Set Creation

In order to fine-tune the LLM for the RAG-based
verdict production task, three main elements are
needed: a claim, a gold answer, and the context
comprising the knowledge needed to reply. In Full-
Fact and VerMouth, the knowledge is present in the
form of a fact-checking article. This comes in use-
ful when the entire article is used as a context, but
when working with chunks a proper selection of the
most informative chunks must be performed. To
this end, we started from the gold verdicts present
in the two aforementioned datasets, and we ranked
each article's chunks given the verdict information
using a cross-encoder reranking model, i.e. the
BAAI/bge-reranker-large[21]. Both for the arti-
cles and the chunk configurations we add to the
context of each training entry some negative exam-
ples, as in testing time the retrieved content might
comprise articles or chunks that are not gold. To
do so, we employed BM25 for retrieving 10 arti-
cles/chunks for each gold verdict and selected the
non-gold retrieved context. An example of training
input is provided in Table 8.

### E.3   Fine-Tuning Details

We fine-tuned the Llama-2-13b[22] chat model on
different subsamples of training data from the Full-
Fact and VerMouth datasets. From the training
dataset, created following the procedure explained
in Section E.2, we randomly extracted 200 entries
each from the FullFact and SMP datasets. For the
emotional datasets, we sampled 35 entries for each
emotion, totalling 210 training entries. An example
of input is shown in Table 8.
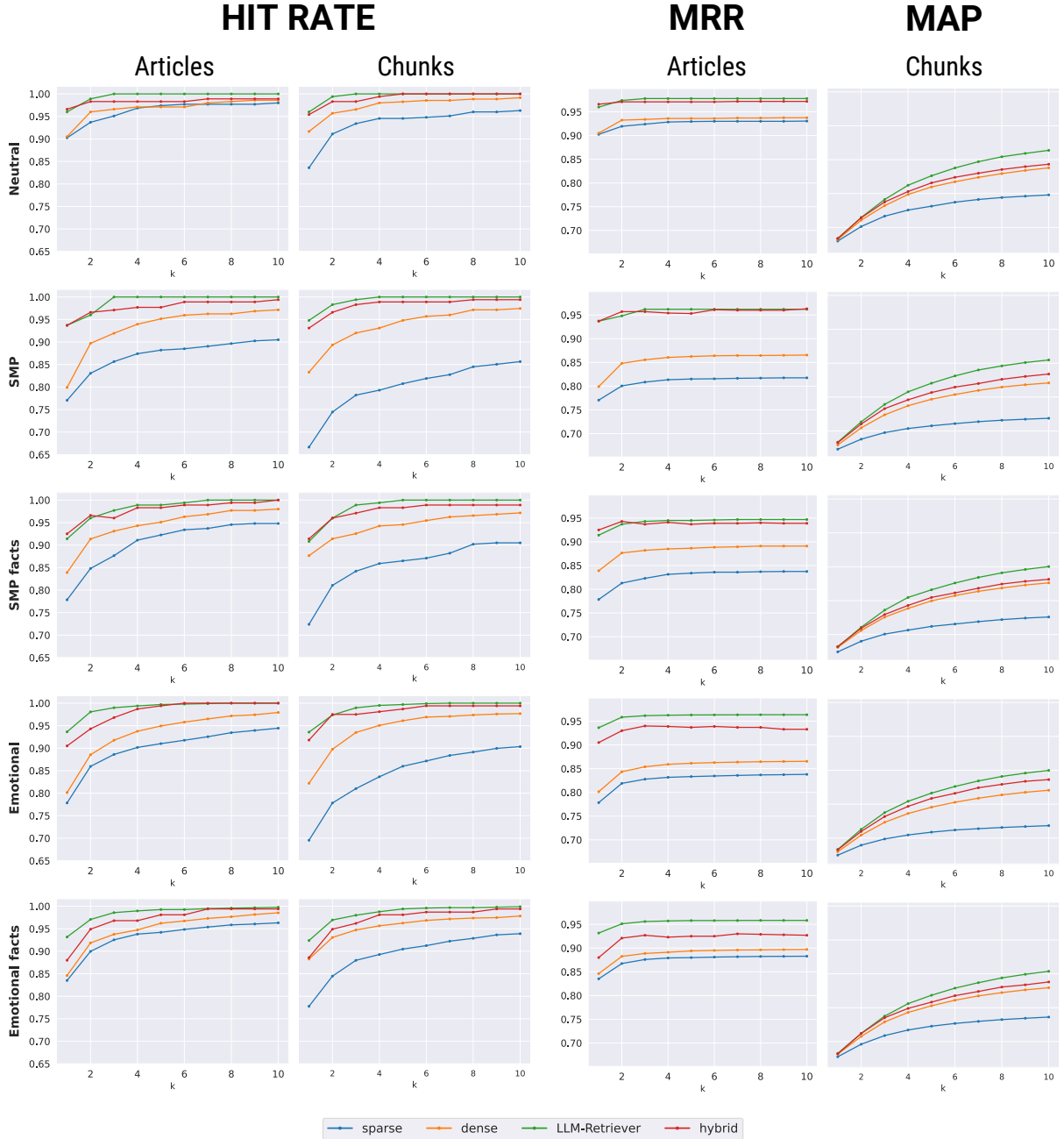
  All models were trained on a single Ampere

---

[21]https://hf.co/BAAI/bge-reranker-large
[22]https://hf.co/meta-llama/Llama-2-13b-chat-hf

Figure 3: Retrieval results for each type of retriever (sparse, dense, LLM, hybrid) across Gold_KB$_{art}$ and Gold_KB$_{chunks}$ are presented for all claim styles, both with (SMP Facts, Emotional Facts) and without (neutral, SMP, emotional) claim pre-processing. The metrics reported include hit_rate and MRR for retrieval over Gold_KB$_{art}$, and hit_rate and MAP for Gold_KB$_{chunks}$, for increasing values of retrieved documents/chunks ($k = 1, ..., 10$).

A40 with 48GB memory using the QLoRA strategy (Dettmers et al., 2023), with a low-rank approximation set to 64, a low-rank adaptation set to 16, and a dropout rate of 0.1. Evaluation steps were set at 25, and the batch size was 4. All models were trained for 3 epochs with a learning rate of $10^{-4}$.

### E.4 Generation Results

In Table 9 we report the complete results for zero-shot and one-shot experiments using chunks and articles as information context.

### E.5 Generation Examples

In Table 10 and 11 we show examples of generations with claims from both FullFact and VerMouth
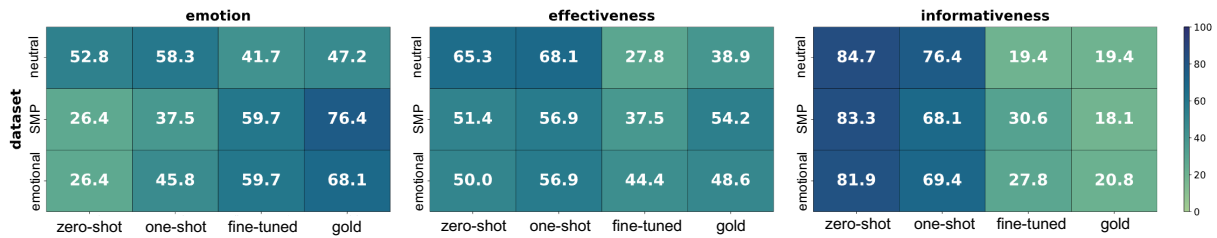
Figure 4: Complete results for the human evaluation. Each matrix refers to the results obtained for each verdict evaluation aspect. The matrices report how many times, in percentage, the human annotators preferred each of the four generation setups (gold, zero-shot, one-shot, fine-tuning).

(anger emotion) datasets. Each table comprises the following information: the claim; the gold verdict; the generations with `Llama-2-13b-chat` model in zero-shot, one-shot, and fine-tuning settings; the relevant evidence retrieved (either chunks, Table 10, or articles, Table 11).

## F   Human Evaluation Details

For the human evaluation of the generated verdict, we enrolled three volunteer evaluators. They were provided with pairs of verdicts (either gold or generated using zero-shot, one-shot, or fine-tuned models) and their corresponding claims. They were instructed to assess the best verdict based on three aspects: effectiveness, informativeness, and emotional/empathetic coverage. Hereafter, we list the tasks/questions that evaluators were required to follow when judging the verdicts pair.

### Informativeness

Tell which of the two verdicts contains more information supporting its stance.

### Emotional Coverage

Some of the claims can express a variety of emotions. Tell which of the two verdicts better takes into consideration the emotion of the claim by responding with empathy.

### Effectiveness

Which of the two is an overall better verdict (with respect to the claim) that could be used to answer the claim?

Given that the claims presented may cover sensitive subjects, we have incorporated a cautionary note in the task description: *"This task may contain text that some readers find offensive."*. Additionally, we briefed the evaluators on the study's objectives and assured them that all collected data would be anonymized and solely utilized for research purposes.

In Figure 4 we report the full outcome of the human evaluation, showing the details of the selected verdicts over the three datasets according to emotional coverage, informativeness and effectiveness.

<s>[INST] <<SYS>> Based on the provided context, respond to the claim, ensuring a thorough explanation. Use only the given context. Reply in no more than three sentences. Avoid mentioning the context in the reply. Match the communication style of the claim and address the possible emotional component present in it, if needed. If the context is insufficient, state that you don't know. Format your response as follows:

Reply: [your_reply]

<</SYS>>

The context information is provided below (in between xml tags).

**<context>**

A meme shared on Facebook features actor John Krasinski in The Office with a whiteboard with edited text, which says: "3 countries refused the covid vaccine", followed by: "Now all 3 of their presidents have died unexpectedly". Beneath the image are the names of the former presidents of Haiti (Jovenel Moïse), Tanzania (John Magufuli) and Zambia (Kenneth Kaunda).

The president did not refuse the Covid-19 vaccines for Zambia. In fact, in March 2021, the Zambian health minister announced plans to vaccinate all over 18s in the country. Similar claims have been fact checked before.

This survey covers households in England and Wales and so does not cover groups (such as those living in student halls of residence), who have "potentially high proportions of drug use", meaning the true figure could be higher. Comparing England & Wales to other countries in Europe is difficult because not all countries have up to date data.

It's correct that cocaine use among 16 to 24 years olds in England and Wales is at its highest level for around a decade. In 2017/18 6% said they had used at least once in the previous year. The claim referred to Britain, but used data covering only England & Wales. We're focusing on England & Wales as data for Scotland and Wales are not available for the most recent year.

There is no evidence to suggest that the death was related to Mr Magufuli's stance on the Covid-19 vaccines. There has been some speculation from Tanzanian opposition leaders, and on social media, that Mr Magufuli's death may have been caused by Covid-19, however this has been discredited. President Kaunda died of pneumonia at a military hospital in Lusaka in June 2021, age 97.

President Magufuli reportedly said that home treatments such as steam inhalation were preferable to "dangerous foreign vaccines", and in February 2021 the country's health minister said that Tanzania had no plans to accept Covid-19 vaccines. Mr Magufuli's successor, president Samia Suluhu Hassan announced that the president's death in March 2021 was due to heart disease.

Arrests have been made but there are still many unknown details about the assasination. There is no evidence to suggest that there is a link to the lack of progress made regarding Haiti's vaccine roll out.

President Moïse did not explicitly refuse all of the Covid-19 vaccines, but the country did initially refuse the AstraZeneca vaccine due to safety concerns. There have also been administrative and supply difficulties around the arrival of vaccinations to Haiti, and concerns about the logistics of organising the vaccine roll out. Mr Moïse was fatally shot at home on 7 July 2021.

The posts suggest that the death of these three men is linked to a refusal of the Covid-19 vaccines. There is no evidence to support this claim. Haiti is eligible for provision of Covid-19 vaccines through the Covax programme, but roll out has been slow, with Haiti yet to give any doses to residents, and only receiving its first vaccines this week.

The Royal College of Obstetricians & Gynaecologists (RCOG) says that there is "no evidence" to suggest that the Covid-19 vaccines affect fertility. Similar statements have been issued by the British Fertility Association, the NHS and Public Health England. The RCOG also states that there is no biologically plausible mechanism by which the vaccine would cause fertility problems.

**</context>**

**Claim:** "The presidents of Haiti, Tanzania, & Zambia all passed away suddenly & unexpectedly. Reports suggest it may be connected to their refusal to accept the Covid-19 vaccines for their countries. This is a tragic loss for the world & for the people of those countries. #vaccinesaveslives #RIP" [/INST]

**Reply:** "I'm sorry to hear about the recent passing of the presidents of Haiti, Tanzania, and Zambia. However, it's important to clarify that the president of Zambia did not refuse vaccines and passed away from pneumonia at a ripe old age, while the president of Tanzania passed away from a heart condition. The assassination of the president of Haiti is a tragic event, but there is no evidence linking it to vaccine refusal. Let's honor their legacies and focus on ways to prevent further loss of life."

Table 8: Training example for `Llama-2-13b` model. The positive passages are highlighted in green, while negative in red.

| | | Articles | | | | Chunks | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | **Zero-shot** | | | | | | | |
| **model** | **dataset** | **ROUGE-LSum** | **BARTScore** | **SummaC** | **GoldSim** | **ROUGE-LSum** | **BARTScore** | **SummaC** | **GoldSim** |
| **mistral-7b-v0.1** | fullfact | 0,216 | -1,948 | 0,381 | -3,101 | 0,202 | -1,992 | 0,361 | -3,112 |
| | SMP | 0,185 | -2,153 | 0,333 | -2,990 | 0,190 | -2,207 | 0,351 | -3,060 |
| | emotional | 0,174 | -2,389 | 0,322 | -2,953 | 0,188 | -2,411 | 0,327 | -2,990 |
| **mistral-7b-v0.2** | fullfact | 0,130 | -2,315 | 0,331 | -3,084 | 0,144 | -2,251 | 0,334 | -3,101 |
| | SMP | 0,140 | -2,516 | 0,327 | -3,144 | 0,145 | -2,453 | 0,335 | -3,078 |
| | emotional | 0,131 | -2,803 | 0,312 | -3,134 | 0,145 | -2,789 | 0,317 | -3,143 |
| **llama3-8b** | fullfact | 0,050 | -3,143 | 0,310 | -3,516 | 0,035 | -3,439 | 0,319 | -3,648 |
| | SMP | 0,047 | -3,330 | 0,295 | -3,644 | 0,032 | -3,461 | 0,328 | -3,715 |
| | emotional | 0,054 | -3,142 | 0,282 | -3,466 | 0,041 | -3,217 | 0,305 | -3,575 |
| **llama2-7b** | fullfact | 0,168 | -1,979 | 0,330 | -2,914 | 0,181 | -1,988 | 0,330 | -2,932 |
| | SMP | 0,160 | -2,118 | 0,321 | -2,861 | 0,195 | -2,022 | 0,330 | -2,898 |
| | emotional | 0,159 | -2,253 | 0,311 | -2,755 | 0,197 | -2,171 | 0,318 | -2,781 |
| **llama2-13b** | fullfact | 0,176 | -1,714 | 0,353 | -2,787 | 0,195 | -1,718 | 0,355 | -2,811 |
| | SMP | 0,169 | -1,849 | 0,331 | -2,714 | 0,186 | -1,879 | 0,352 | -2,752 |
| | emotional | 0,156 | -2,118 | 0,323 | -2,691 | 0,183 | -2,058 | 0,338 | -2,754 |
| | | **One-shot** | | | | | | | |
| **mistral-7b-v0.1** | fullfact | 0,173 | -2,186 | 0,336 | -3,073 | 0,180 | -2,224 | 0,341 | -3,166 |
| | SMP | 0,160 | -2,369 | 0,339 | -3,043 | 0,161 | -2,523 | 0,325 | -3,018 |
| | emotional | 0,145 | -2,586 | 0,307 | -2,947 | 0,162 | -2,664 | 0,302 | -3,012 |
| **mistral-7b-v0.2** | fullfact | 0,134 | -2,347 | 0,321 | -3,177 | 0,136 | -2,393 | 0,320 | -3,206 |
| | SMP | 0,123 | -2,630 | 0,301 | -3,160 | 0,130 | -2,623 | 0,304 | -3,117 |
| | emotional | 0,116 | -2,865 | 0,318 | -3,138 | 0,125 | -2,871 | 0,321 | -3,148 |
| **llama3-8b** | fullfact | 0,133 | -1,960 | 0,351 | -3,003 | 0,135 | -1,905 | 0,344 | -2,988 |
| | SMP | 0,138 | -2,322 | 0,323 | -3,011 | 0,137 | -2,377 | 0,321 | -3,013 |
| | emotional | 0,131 | -2,469 | 0,328 | -2,927 | 0,134 | -2,485 | 0,334 | -2,902 |
| **llama2-7b** | fullfact | 0,171 | -2,091 | 0,327 | -2,977 | 0,197 | -2,016 | 0,325 | -2,889 |
| | SMP | 0,158 | -2,408 | 0,305 | -2,922 | 0,181 | -2,467 | 0,304 | -2,911 |
| | emotional | 0,147 | -2,561 | 0,317 | -2,875 | 0,191 | -2,549 | 0,306 | -2,869 |
| **llama2-13b** | fullfact | 0,185 | -2,057 | 0,322 | -2,942 | 0,199 | -1,992 | 0,309 | -2,840 |
| | SMP | 0,154 | -2,355 | 0,314 | -2,920 | 0,172 | -2,235 | 0,316 | -2,803 |
| | emotional | 0,139 | -2,523 | 0,314 | -2,839 | 0,152 | -2,441 | 0,303 | -2,737 |

Table 9: Complete results for each model tested on the three datasets in zero-shot and one-shot settings. Results for both chunks and article configurations are reported.

| | CLAIM: Research says that nearly half of people who catch Covid develop depression afterwards. | |
|---|---|---|

**CLAIM:** Research says that nearly half of people who catch Covid develop depression afterwards.

**GOLD VERDICT:** This is not what the research says. It measured rates of depression among Covid survivors—but not whether they developed depression after catching the disease.

| zero-shot | one-shot | fine-tuning |
|---|---|---|
| The claim is partially accurate but misleading. The study in question found that 48% of participants experienced moderate or severe depression, but it did not establish a causal link between Covid-19 and depression. The study also did not compare depression rates among those who had Covid with those who had not, nor did it consider the normal background rate of depression in Bangladesh. Additionally, the study did not find that Covid-19 directly causes depression. Instead, it suggests that people suffering from depression may be more likely to experience long-term symptoms after recovering from Covid-19. | The claim is misleading and not entirely accurate. The research in question found that nearly half of the participants reported experiencing depression "alongside" Covid, but it does not necessarily imply that Covid caused the depression. The study did not compare the prevalence of depression among those who caught Covid to those who did not, nor did it account for the normal background rate of depression in Bangladesh. Additionally, the study did not find any causal link between Covid and depression. The researchers themselves noted that the findings may not reflect the impact of Covid-19, especially for depression, which may have developed prior to the pandemic. It is important to be cautious when interpreting the results of scientific studies and not draw conclusions that go beyond the evidence. | The study found that 48% of people who had Covid had symptoms of moderate or severe depression. But the researchers themselves say that this does not necessarily reflect the impact of Covid, and that it may have developed before the pandemic. |

**1 -** The study in question explicitly said that its findings didn't necessarily reflect the impact of Covid-19.The newspaper said: "Nearly half of people struck down with Covid-19 also develop depression, according to new research."The research in question surveyed the mental health of 1,002 people in Bangladesh who had tested positive for Covid at some point in the past.

**2 -** These people were a self-selected sample, meaning they heard about the survey and wanted to take part, so they might not be representative of people as a whole.The study found that about 48% had moderate or severe depression.However, it does not tell us that they "go on" to develop depression as a result of catching Covid, or even that they developed it after catching Covid.

**3 -** One of the study authors says it describes people suffering from depression "alongside" Covid.The researchers themselves say in the paper: "The findings may not reflect the impact of COVID-19, especially for depression, which may have developed prior to the COVID-19 pandemic."There is evidence that rates of depression have increased in Great Britain and elsewhere during the pandemic.

**4 -** The Telegraph has muddled up the meaning of a scientific study on depression among people who have had Covid-19.

**5 -** The research doesn't compare its findings with depression rates among these people before they got the disease, nor with people who hadn't been infected at all, nor with the normal background rate of depression in Bangladesh (which may have risen during the pandemic).A press release from Anglia Ruskin University, which worked on the research, does not say that Covid was found to cause depression.

**6 -** There is other evidence that people may be at higher risk of developing mental health problems after recovering from Covid.The extent of this and why it might happen are not yet fully understood, however.

**7 -** 26% had symptoms of moderate or severe levels of depression with no real change across the first six weeks, which may be what the "one in four" claim was meant to refer to. Feelings of defeat and entrapment decreased over time, while positive wellbeing increased.Although the findings of the study were not as dramatic as claimed in the video, researchers said there was cause for concern.

**8 -** So Mr Hancock was right to say that one in 20 people would have long-term symptoms after a Covid infection—if you define "long-term" as more than eight weeks.However, he was wrong to compare this to the one-in-10 figure for people under 50, because this describes those who suffered from symptoms that lasted more than four weeks—rather than eight.As the research paper explains,

**9 -** Mr Hancock was also wrong to say that there was "no correlation" between the severity of the initial illness and the length of the symptoms afterwards.The research paper says people with long Covid were "more likely to have required hospital assessment in the acute period." Overall, 13.9% of the confirmed cases in the study had visited hospital. However, among those reporting symptoms for at least four weeks this rose to 31.5%,

**10 -** and among those reporting symptoms for at least eight weeks it was 43.9%.The app data also suggested that "Individuals reporting more than 5 symptoms in the first week (the median number reported) were significantly more likely to go on to experience LC28 symptoms lasting at least four weeks."In short, more severe cases of Covid seemed more likely to result in long Covid afterwards.

Table 10: Example of generation using FullFact claim, e5-mistral as a retriever and Llama-2-13b-chat for the generation of the verdict. In green are highlighted the gold chunks retrieved.

| **CLAIM:** I'm fuming over the fact that the worst violent and sexual criminals are let off the hook halfway through their sentence. How is this justice? I feel outraged and betrayed by the justice system." #JusticeForVictims |
|---|

| **GOLD VERDICT:** I understand your frustration and anger towards the justice system. It is important to note that judges in England and Wales can already choose to keep dangerous offenders in prison for up to two thirds of their sentence. I hope this information provides some clarity. #JusticeForVictims |
|---|

| **zero-shot** | **one-shot** | **fine-tuning** |
|---|---|---|
| It's understandable that you feel outraged and betrayed by the justice system, but it's important to note that the claim that the worst violent and sexual criminals are let off the hook halfway through their sentence is not entirely accurate. While it is true that some criminals are eligible for release halfway through their sentence, this is not automatic and is subject to the discretion of the judge. Additionally, judges already have the power to keep serious criminals in prison for longer than the normal halfway point, and they consider a variety of factors when determining the appropriate sentence. It's also worth noting that the proposed policy change by the Conservative party would not necessarily lead to longer sentences for all criminals, as judges may continue to sentence in the same way as they currently do. It's important to ensure that any changes to the justice system are fair, just, and proportionate, and that the needs of victims and the community are taken into account. | I understand your frustration, but it's important to note that the current sentencing procedure already allows for judges to keep serious criminals in prison for longer periods of time. The proposed change to mandatory two-thirds sentences would not necessarily lead to more time being served, as judges would still have discretion to sentence criminals based on the circumstances of the case. Additionally, the proposal only applies to the most serious violent and sexual offenders, and it's important to remember that the justice system is designed to balance punishment with rehabilitation and public protection. Let's focus on ensuring that victims receive the support and justice they deserve, while also ensuring that the system is fair and effective for all parties involved. | Thank you for sharing your concerns. It's important to note that the justice system is designed to provide a balance between punishment and rehabilitation. While it's understandable to feel outraged when it comes to serious crimes, it's important to remember that the current sentencing system already allows for judges to keep serious criminals in prison for longer periods of time if they deem it necessary. Additionally, the government has recently announced plans to make it mandatory for the most serious violent and sexual offenders to serve two-thirds of their sentence in prison, which may address your concerns. However, it's important to remember that the justice system is complex and there are many factor |

| At the Conservative party conference last week, justice secretary Robert Buckland pledged to make it mandatory for the most serious violent and sexual offenders to serve two thirds of their sentence in prison (and not half the sentence as is normal). The Ministry of Justice told us this would apply to "the most serious sexual and violent offenders where the offence carries a maximum life sentence." However, to suggest that the most serious criminals are automatically released from prison halfway through their sentences obscures the fact that judges already have discretion to keep serious criminals in prison for longer. There are various sentences a judge can hand down and in most cases sentences are non-custodial (where no time is served in prison). By far the most common sentence for crimes in England and Wales is a fine, but what's of interest here are sentences which carry mandatory prison time. Typically in these cases a criminal will be given a standard determinate sentence. This usually requires them to spend half of their sentence in prison and the other half on license in the community, supervised by the probation service. For example, a standard two year sentence would involve one year in prison and one year on license. Being on licence means you can be recalled to prison if you breach the terms of your licence. As well as standard sentences, judges in England and Wales can hand down what are called 'extended determinate sentences' to criminals who commit any of over 100 serious offences. The judge can make this decision if: These offenders are either entitled to be released two thirds of the way through their sentence, or can apply for parole at that point. Life sentences work slightly differently. With a life sentence a criminal is required to spend a minimum time in prison and is then able to apply for parole. If they are released, they remain on license for the rest of their life. Compared to existing extended sentences, the Conservatives' proposal appears to apply to criminals who commit a slightly different group of offences (those that carry a maximum of life rather than this list of serious offences). There is also apparently no requirement for a judge to determine if a criminal poses a risk to the public when giving this new kind of sentence. While the current sentencing procedure does not dramatically change the ability to put serious criminals in prison for two thirds of their term, it would, in practice, significantly increase the number of criminals receiving two thirds sentences. That's because judges rarely hand down extended sentences. The Ministry of Justice says that in 2018 there were around 4,000 standard sentences with halfway release handed down to criminals who committed sexual or violent offences which carry the maximum penalty of life. By comparison, in 2018 judges in England and Wales handed down 398 extended sentences. There is an open question over whether the policy would in fact lead to serious criminals spending more time in prison, because it's possible that judges could change how they currently sentence. [...] |
|---|

Table 11: Example of generation using VerMouth anger claim, e5-mistral as a retriever and Llama-2-13b-chat for the generation of the verdict. The article has been cut for space reasons. The complete article's text can be found at https://fullfact.org/crime/extended-sentences/