

# "Politi-Fact-Only": A Political Domain Benchmark Dataset

Anonymous ACL submission

## Abstract

The rapid proliferation of online information has made it increasingly challenging to differentiate factual content from misinformation. Traditional fact-checking methods, which require extensive manual effort, are not scalable given the volume of misinformation spreading online. Automated fact-checking has emerged as a promising solution, leveraging machine learning models trained on datasets derived from fact-checking websites (Wang, 2017; Augenstein et al., 2019; Gupta and Srikumar, 2021a). However, many of these datasets include post-analysis commentary from annotators, which may introduce bias and provide implicit cues that aid model performance. To address this limitation, we introduce Politi-Fact-Only, a benchmark dataset comprising 1,482 instances curated from PolitiFact.com, where we remove post-analysis and retain only factual evidence Fig 1. This ensures that models must rely solely on factual reasoning rather than verdict-related information. Our experiments demonstrate that state-of-the-art fact-checking models, including large language models (LLMs), struggle to accurately classify claims when deprived of post-claim analysis, highlighting their reliance on implicit cues rather than pure factual reasoning.

## 1 Introduction

The proliferation of online information has accelerated the spread of both factual and misleading content, making it increasingly difficult for the public to discern truth from falsehood. In response to this growing challenge, fact-checking platforms like *PolitiFact*<sup>1</sup>, which shows verdict for the claims into varying degrees of accuracy, from **True** to **Pants on fire** and intermediate stages like **Mostly True**, **Half True**, **Mostly False**, and **False**. These labels

<sup>1</sup><https://www.politifact.com/>

---

**claim:** A photo shows a crash in Eminence, Indiana.  
**evidence:** ~~With wintry weather striking many regions of the United States, photos are emerging from across the country showing snow plows, icy roads, and big drifts of snow. But one old picture is being mischaracterized as showing the aftermath of a crash in Eminence, Indiana.~~ This is currently at state road 42 and state road 142 in downtown Eminence, Indiana, a Feb. 2 post says, alongside a photo of a treacherous-looking pileup of cars and trucks. Black ice! Drive safe folks! Right by the Citizens Bank and Dairyland. Prayers to all involved. But this photo was taken last year, and about 900 miles southwest of Eminence. ~~This post was flagged as part of Facebook's efforts to combat false news and misinformation on its News Feed. (Read more about our partnership with Facebook.)~~ Photographer Lawrence Jenkins took it on Feb. 11, 2021, in Fort Worth, Texas, where 133 vehicles crashed after freezing rain coated the roads there, sending dozens of people to the hospital and leaving at least six dead, the Dallas Morning News reported at the time. North Texas and Central Indiana are both experiencing wintry weather, ~~but this photo doesn't show it.~~  
**label:** false

---

**source:** social\_media  
**speaker:** Facebook posts  
**claim\_data:** 19/11/2018  
**factchecker:** Jill Terrell Ramos  
**fact\_check\_data:** 7/12/2018  
**factcheck\_analysis\_link:** <https://www.politifact.com/factchecks/>....

---

Figure 1: An example from the dataset where sentences with strike-through lines represent information added only after the claim was verified. These lines were manually deleted to ensure the evidence contained only factual details sufficient for fact-checking the claim. Also, we show the meta-data available for an instance.

reflect the complexity of misinformation, where claims often contain elements of truth mixed with misleading or omitted details, whereas half-truths present unique challenges. They selectively expose the truth, exploiting human cognitive biases to manipulate perceptions (Estornell et al., 2020). Unlike outright falsehoods, which are often easier to detect, half-truths thrive on ambiguity. This makes them highly effective in shaping public opinion, particularly in areas like politics, advertisement, and finance, where they are strategically employed to influence decision-making.

Fact-checking is a laborious process that requires significant time and effort. Journalists need to sift through multiple sources to verify claims, assess the credibility of those sources, and draw meaningful comparisons. This process, which can take several hours or even days for professional fact-

checkers (Hassan et al., 2015), is often further strained by tight deadlines, especially for internal fact-check procedures (Godler and Reich, 2017). Research indicates that less than half of the published articles undergo verification (Lewis et al., 2008). With the rapid pace of information generation and dissemination, manual fact-checking alone is not scalable, highlighting the need for automation (Guo et al., 2022).

Several studies have explored automated fact-checking, including works by (Wang, 2017), (Augenstein et al., 2019), and (Gupta and Srikumar, 2021a), which have contributed valuable datasets. While these datasets contain real-world claims, they are primarily derived from fact-checking websites. The articles on such platforms often present a post-analysis of claims, incorporating assessments from annotators based on factual evidence. However, this does not fully reflect real-world fact-checking scenarios. To address this limitation, some researchers, such as (Yang et al., 2022) and (Khan et al., 2022), have focused on utilizing premise articles or sources that were published before the claim itself. This approach brings the fact-checking process closer to real-world settings. However, while relevant information is extracted from these documents based on the claim, there is no guarantee that the retrieved content is sufficient for verification. To bridge this gap, we propose a benchmark political domain test set Politifact-Facts-Only Section 4, a subset of Misra (2022). In this dataset, we manually remove the post-analysis provided by annotators and retain only the factual information. This ensures a more realistic evaluation of automated fact-checking models.

Figure 1 illustrates an example from the dataset, where annotators manually reviewed the instances and removed sentences containing post-claim analysis. In real-world scenarios, fact-checking relies solely on factual information, requiring reasoning based on these facts without the aid of post-publication commentary. While previous approaches (Khan et al., 2022) have attempted to address this by using review or premise articles to avoid post-analysis content, this raises a critical concern about the effectiveness of an abstract summary of evidence extracted for accurate fact-checking of claims.

Our contributions are:

1. We introduce *Politi-Fact-Only*, a benchmark dataset comprising 1,482 instances curated

from PolitiFact.com<sup>2</sup>. As detailed in Table 1, the dataset has been manually filtered to remove post-claim analyses and verdict-related information originally present in the articles (Section 4). This ensures that the evidence consists solely of factual content, eliminating potential annotator bias introduced by verdict cues and improving the reliability of fact-checking models. Initially, we collected 1,500 instances. However, to maintain accuracy and credibility, we removed 18 instances due to reasons outlined in Section 5.

2. Through experiments we show the performance of *Politi-Fact-Only* along with other various datasets, Table 3 and 2. We observe that on our test set models are struggling to reason about the facts to support or refute the claim. Large language models (LLMs) struggle to reason effectively when limited to fact-only evidence from the *Politi-Fact-Only* dataset. In contrast, LLMs perform comparatively better on other datasets in zero-shot settings, highlighting their reliance on implicit cues and verdict-related information rather than pure factual reasoning.

## 2 Problem Statement

**Input:** A claim  $C = \{c_1, c_2, \dots, c_n\}$  and its corresponding evidence  $E = \{e_1, e_2, \dots, e_m\}$ , where  $c_i$  and  $e_j$  represent individual tokens in the claim and evidence, respectively.

**Output:** A verdict label  $L$ , where  $L$  belongs to {True, Mostly True, Half True, Mostly False, False}

The goal is to classify the claim  $C$  based on the factual content of  $E$ , determining its degree of truthfulness.

## 3 Related Work

Existing fact-checking datasets can be broadly categorized into meta-based and text-based datasets. Meta-based datasets, such as LIAR (Wang, 2017) and (Rashkin et al., 2017), primarily include claims with metadata like speaker identity and historical records but lack supporting textual evidence, limiting their utility for verification. Similarly, Vlachos and Riedel (2014) compiled a small dataset of political claims, but without explicit evidence, restricting its effectiveness in real-world fact-checking.

Text-based datasets offer stronger evidence-grounded verification, with FEVER (Thorne et al., 2018), HOVER (Jiang et al., 2020), relying solely on Wikipedia as a knowledge base. While valuable, these datasets fail to capture misinformation from diverse sources beyond Wikipedia. Some datasets, such as Multifc (Augenstein et al., 2019) and X-fact (Gupta and Srikumar, 2021b), incorporate evidence from broader domains. LIAR-PLUS (Alhindi et al., 2018) attempted to provide evidence by extracting the last five sentences from source articles, but this often resulted in incomplete or irrelevant context. L++ (Russo et al., 2023) and ru22fact (Zeng et al., 2024) make use of fact-checking website for the dataset creation. Khan et al. (2022) and (Yang et al., 2022) worked on more real-world situations by fetching content from the article that exists before the claim was published. To address these limitations, Politi-Fact-Only builds upon PolitiFact (Misra, 2022) by filtering out post-claim analysis and retaining only factual evidence. This ensures a more realistic test set for evaluating models on factual reasoning without relying on annotator cues.

Label	Count	Token <sub><math>\mu</math></sub>	Sent <sub><math>\mu</math></sub>	BPE <sub><math>\mu</math></sub>
True	296	596.05	26.05	755.94
Mostly True	298	682.90	30.88	865.19
Half True	293	756.69	33.61	954.85
Mostly False	300	780.14	34.88	978.39
False	295	559.09	27.26	705.01
<b>Total</b>	1482	675.18	30.55	852.13

Table 1: Statistics for *Politi-fact-only* dataset. Token <sub>$\mu$</sub> , Sent <sub>$\mu$</sub> , and BPE <sub>$\mu$</sub>  represent the average number of standard tokens, sentences, and BPE tokens per evidence, respectively.

#### 4 Politi-fact-only: A Fact Only Benchmark Dataset

Our dataset, *Politi-Fact-Only*, ensures that each claim is accompanied by evidence containing facts related to the claim, supporting its veracity. We randomly selected 1,500 instances from the PolitiFact dataset (Misra, 2022), sourced from Politifact.com. Section 5 details our manual filtration and annotation process. We removed instances where the predictions did not match, resulting in a final dataset of 1,482 instances, as shown in Table 1. Upon comparing Table 1 with Table 4 in the Appendix, we observe a decrease of approximately 15% in the average after filtration. This suggests that around 15% of the content in the article consisted of com-

mentary from the annotators.

Each record in *Politi-Fact-Only* contains nine attributes. We retain the following key attributes: *label*, *claim*, *evidence*, *speaker*, *factcheck\_analysis\_link*, *factcheck\_date*, *factchecker*, *claim\_date*, and *claim\_source*. Additionally, the *false* and *pants-fire* labels were merged into a single *false* label, as both categories represent completely false information. Examples from each class can be found in the Appendix, Figures 2 - 6.

We retain the claim’s publish date and the fact-check date for relevant use cases, such as extracting evidence articles using the Google API. By using the claim verbatim, as done by (Yang et al., 2022), we can filter out articles published after the claim’s publish date and fact-check date, depending on the research needs. The speaker of a claim is crucial, as it provides insights into their reliability. If a speaker frequently makes false statements, it indicates the speaker’s lack of credibility. The same applies to the source, whose history of publications can reflect its credibility.

An example is shown in Figure 1, where strikethrough lines indicate content present solely due to annotator commentary. Our goal is to make this dataset as representative as possible of real-world scenarios, where only factual information is available. To achieve this, we must extract factual content from the web. However, extracting relevant information up to a predefined limit (e.g.,  $k$ ) does not guarantee that the summary will be sufficient to assess the veracity of the claim.

#### 5 Annotation and Filtration Process

We hired three annotators, all proficient in English, to clean 500 instances each. The annotators were instructed to annotate the instances while also applying a filtration process. The evidence for each instance was obtained by scraping fact-checking websites. These articles included post-claim analyses written by fact-checkers from sources such as PolitiFact, who provide ratings ranging from “True” to “Pants on Fire.”

As shown in Appendix B, we established specific guidelines for the filtration process and provided some filtered instances as reference. By adding two new fields for each instance: **leaked** and **annotator\_prediction** we asked annotators to fill this accordingly. We removed sentences from the evidence if they were directly related to the verdict,

Models	MultiFC	Liar Plus	RU22fact	L++	Politi-Fact-Only	Unfiltered
Mistral	0.1419/ <b>0.2952</b>	0.2060/ <b>0.2822</b>	0.3239/ <b>0.6554</b>	0.2874/ <b>0.3607</b>	0.2667/ <b>0.3475</b>	0.3725/ <b>0.4575</b>
LLaMA	0.1579/ <b>0.2532</b>	0.1853/ <b>0.2580</b>	0.2933/ <b>0.6462</b>	0.2851/ <b>0.3553</b>	0.2115/ <b>0.2861</b>	0.4811/ <b>0.5142</b>
Gemma	0.1586/ <b>0.2540</b>	0.0708/ <b>0.1723</b>	0.2710/ <b>0.6427</b>	0.1845/ <b>0.3107</b>	0.2057/ <b>0.2821</b>	0.5904/ <b>0.6005</b>

Table 2: Performance comparison of models ranging from 7B to 9B parameters using Zero-Shot prompting (Kojima et al., 2024) across various fact-checking datasets. The results are reported in macro F1/micro F1-score. The "Unfiltered" dataset represents the unfiltered version of *PolitiFact-Fact-Only*. We use models like Meta’s LLaMA (Dubey et al., 2024) (*meta-llama/Meta-Llama-3.1-8B*), Mistral’s version 3 (Jiang et al., 2023) (*mistralai/Mistral-7B-v0.3*) and, model from Google, such as the GEMMA series (Team et al., 2024) (*google/gemma-2-9b*) from huggingface.

Dataset - Model	Respective Dataset	PolitiFact-Fact-Only
LIAR-PLUS - SVM (Alhindi et al., 2018)	0.25	0.27
LIAR - CNN (Wang, 2017)	0.27	0.28
LIAR-PLUS - LR (Alhindi et al., 2018)	0.37	0.27
AVERTeC - BERT-large (Schlichtkrull et al., 2023)	0.49	0.29

Table 3: Comparison of *PolitiFact-Only* with other fact-checking datasets. All results are reported in Macro F1-score. The values for the respective datasets are sourced from the original authors’ reported results.

such as statements like “This post was flagged as part of...”. Similarly, we eliminated annotator commentary, such as “But this photo doesn’t show it”. These logical or inference-based cues can make reasoning easier; however, in real-world scenarios, only factual information is available, making it significantly harder to assess the claim’s veracity based solely on the provided evidence.

Consequently, we removed 15 instances where the annotator’s prediction did not match the expected outcome. Additionally, we excluded three more instances because, after the cleaning process, the reduced context no longer provided enough information to support or refute the claim. By publishing this dataset, we aim to provide a resource that can be used for both explanation generation and fact-checking classification tasks.

## 6 Experimental and Results Analysis

We conducted two experiments: one utilizing existing fact-checking datasets sourced from PolitiFact.com and another leveraging large language models, including *meta-llama/Meta-Llama-3.1-8B* (Dubey et al., 2024), *mistralai/Mistral-7B-v0.3* (Jiang et al., 2023), and *google/gemma-2-9b* (Team et al., 2024). As shown in Table 3, our dataset outperformed LIAR-PLUS when evaluated using an SVM classifier. Alhindi et al. (2018) primarily extracted the "Our Ruling" or "Our Rating" section as evidence when available; otherwise, they relied on the last five lines of the article. Yang et al. (2022) further refined this approach by retaining only instances containing the "Our Ruling" or "Our

Rating" section, using it as gold-standard evidence for comparison with their generated explanations. We also tried different prompts for different models and we analyze that a single keyword can effect LLM’s output Appendix A.

## Conclusion and Future Work

We introduced *PolitiFact-Only*, a benchmark dataset for evaluating fact-checking models using only factual evidence, without post-claim analysis. Our experiments show that models struggle significantly when deprived of annotator cues, resulting in a notable performance drop compared to unfiltered datasets. This highlights their reliance on implicit signals rather than pure factual reasoning. While large language models (LLMs) perform well on traditional datasets, their difficulty in classifying claims accurately in our test set suggests dependence on verdict-related information. This underscores the challenge of building robust fact-checking systems that operate without pre-annotated guidance.

For future work, we propose using PolitiFact-Only as a test set to evaluate retrieved or summarized information from web articles and documents, assessing whether automated methods preserve enough factual content for verification. Additionally, we aim to enhance fact-checking models by incorporating reasoning-driven approaches that rely solely on factual evidence. Addressing these challenges will contribute to the development of more transparent, and effective fact-checking systems for real-world misinformation detection.



## Limitation

This dataset is collected from a fact-checking website. While we have attempted to remove most annotator cues, some sentences could not be eliminated without compromising the context necessary to support or refute the claim.

## References

- 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.
- Isabelle Augenstein, Christina Lioma, Dongsheng Wang, Lucas Chaves Lima, Casper Hansen, Christian Hansen, and Jakob Grue Simonsen. 2019. [MultiFC: A real-world multi-domain dataset for evidence-based fact checking of claims](#). 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 4685–4697, Hong Kong, China. Association for Computational Linguistics.
- Abhimanyu Dubey, Abhinav Jauhari, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Alonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz

- Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Rapparth, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhennde, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gouget, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkan Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán,

426	Frank Kanayet, Frank Seide, Gabriela Medina Flo-	489
427	rez, Gabriella Schwarz, Gada Badeer, Georgia Swee,	490
428	Gil Halpern, Govind Thattai, Grant Herman, Grigory	491
429	Sizov, Guangyi, Zhang, Guna Lakshminarayanan,	492
430	Hamid Shojanazeri, Han Zou, Hannah Wang, Han-	493
431	wen Zha, Haroun Habeeb, Harrison Rudolph, He-	494
432	len Suk, Henry Aspegren, Hunter Goldman, Ibrahim	
433	Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena	495
434	Veliche, Itai Gat, Jake Weissman, James Geboski,	496
435	James Kohli, Japhet Asher, Jean-Baptiste Gaya,	497
436	Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen,	498
437	Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong,	
438	Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill,	499
439	Jon Shepard, Jonathan McPhie, Jonathan Torres,	500
440	Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou	501
441	U, Karan Saxena, Karthik Prasad, Kartikay Khan-	
442	delwal, Katayoun Zand, Kathy Matosich, Kaushik	502
443	Veeraraghavan, Kelly Michelena, Keqian Li, Kun	503
444	Huang, Kunal Chawla, Kushal Lakhota, Kyle Huang,	504
445	Lailin Chen, Lakshya Garg, Lavender A, Leandro	505
446	Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng	
447	Yu, Liron Moshkovich, Luca Wehrstedt, Madian	506
448	Khabsa, Manav Avalani, Manish Bhatt, Maria Tsim-	507
449	poukelli, Martynas Mankus, Matan Hasson, Matthew	508
450	Lennie, Matthias Reso, Maxim Groshev, Maxim	509
451	Naumov, Maya Lathi, Meghan Keneally, Michael L.	510
452	Seltzer, Michal Valko, Michelle Restrepo, Mihir	511
453	Patel, Mik Vyatskov, Mikayel Samvelyan, Mike	512
454	Clark, Mike Macey, Mike Wang, Miquel Jubert Her-	513
455	moso, Mo Metanat, Mohammad Rastegari, Mun-	
456	ish Bansal, Nandhini Santhanam, Natascha Parks,	514
457	Natasha White, Navyata Bawa, Nayan Singhal, Nick	515
458	Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev,	516
459	Ning Dong, Ning Zhang, Norman Cheng, Oleg	517
460	Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem	518
461	Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pa-	
462	van Balaji, Pedro Rittner, Philip Bontrager, Pierre	519
463	Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratan-	520
464	chandani, Pritish Yuvraj, Qian Liang, Rachad Alao,	521
465	Rachel Rodriguez, Rafi Ayub, Raghotham Murthy,	522
466	Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah	523
467	Hogan, Robin Battey, Rocky Wang, Rohan Mah-	524
468	eswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu,	525
469	Samyak Datta, Sara Chugh, Sara Hunt, Sargun	
470	Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma,	526
471	Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lind-	527
472	say, Shaun Lindsay, Sheng Feng, Shenghao Lin,	528
473	Shengxin Cindy Zha, Shiva Shankar, Shuqiang	529
474	Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agar-	530
475	wal, Soji Sajuyigbe, Soumith Chintala, Stephanie	531
476	Max, Stephen Chen, Steve Kehoe, Steve Satterfield,	532
477	Sudarshan Govindaprasad, Sumit Gupta, Sungmin	
478	Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury,	533
479	Sydney Goldman, Tal Remez, Tamar Glaser, Tamara	534
480	Best, Thilo Kohler, Thomas Robinson, Tianhe Li,	535
481	Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook	536
482	Shaked, Varun Vontimitta, Victoria Ajayi, Victoria	537
483	Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal	538
484	Mangla, Vitor Albiero, Vlad Ionescu, Vlad Poenaru,	539
485	Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li,	
486	Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will	540
487	Constable, Xiaocheng Tang, Xiaofang Wang, Xiao-	541
488	jian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo	542
		543
		544
	Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li,	
	Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam,	
	Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach	
	Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen,	
	Zhenyu Yang, and Zhiwei Zhao. 2024. <a href="#">The llama 3</a>	
	<a href="#">herd of models</a> .	
	Andrew Estornell, Sanmay Das, and Yevgeniy Vorob-	
	eychik. 2020. <a href="#">Deception through half-truths</a> . In	
	<i>Proceedings of the AAAI Conference on Artificial</i>	
	<i>Intelligence</i> , volume 34, pages 10110–10117.	
	Yigal Godler and Zvi Reich. 2017. Journalistic evi-	
	dence: Cross-verification as a constituent of mediated	
	knowledge. <i>Journalism</i> , 18(5):558–574.	
	Zhijiang Guo, Michael Schlichtkrull, and Andreas Vla-	
	chos. 2022. <a href="#">A survey on automated fact-checking</a> .	
	<i>Transactions of the Association for Computational</i>	
	<i>Linguistics</i> , 10:178–206.	
	Ashim Gupta and Vivek Srikumar. 2021a. <a href="#">X-fact: A</a>	
	<a href="#">new benchmark dataset for multilingual fact check-</a>	
	<a href="#">ing</a> . In <i>Proceedings of the 59th Annual Meeting of the</i>	
	<i>Association for Computational Linguistics and the</i>	
	<i>11th International Joint Conference on Natural Lan-</i>	
	<i>guage Processing (Volume 2: Short Papers)</i> , pages	
	675–682, Online. Association for Computational Lin-	
	guistics.	
	Ashim Gupta and Vivek Srikumar. 2021b. X-FACT: A	
	New Benchmark Dataset for Multilingual Fact Check-	
	ing. In <i>Proceedings of the 59th Annual Meeting of the</i>	
	<i>Association for Computational Linguistics</i> , Online.	
	Association for Computational Linguistics.	
	Naeemul Hassan, Chengkai Li, and Mark Tremayne.	
	2015. <a href="#">Detecting check-worthy factual claims in pres-</a>	
	<a href="#">idential debates</a> . In <i>Proceedings of the 24th ACM In-</i>	
	<i>ternational on Conference on Information and Knowl-</i>	
	<i>edge Management, CIKM '15</i> , page 1835–1838, New	
	York, NY, USA. Association for Computing Machin-	
	ery.	
	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-	
	sch, Chris Bamford, Devendra Singh Chaplot, Diego	
	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	
	laume Lample, Lucile Saulnier, L�lio Renard Lavaud,	
	Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,	
	Thibaut Lavril, Thomas Wang, Timoth�e Lacroix,	
	and William El Sayed. 2023. <a href="#">Mistral 7b</a> .	
	Yichen Jiang, Shikha Bordia, Zheng Zhong, Charles	
	Dognin, Maneesh Singh, and Mohit Bansal. 2020.	
	<a href="#">HoVer: A dataset for many-hop fact extraction and</a>	
	<a href="#">claim verification</a> . In <i>Findings of the Association</i>	
	<i>for Computational Linguistics: EMNLP 2020</i> , pages	
	3441–3460, Online. Association for Computational	
	Linguistics.	
	Kashif Khan, Ruizhe Wang, and Pascal Poupart. 2022.	
	<a href="#">WatClaimCheck: A new dataset for claim entailment</a>	
	<a href="#">and inference</a> . In <i>Proceedings of the 60th Annual</i>	
	<i>Meeting of the Association for Computational Lin-</i>	
	<i>guistics (Volume 1: Long Papers)</i> , pages 1293–1304,	

545	Dublin, Ireland. Association for Computational Lin-	Lipschultz, Josh Newlan, Ju yeong Ji, Kareem Mo-	602
546	guistics.	hamed, Kartikeya Badola, Kat Black, Katie Mil-	603
547	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yu-	lican, Keelin McDonell, Kelvin Nguyen, Kiranbir	604
548	taka Matsuo, and Yusuke Iwasawa. 2024. Large	Sodhia, Kish Greene, Lars Lowe Sjoesund, Lau-	605
549	language models are zero-shot reasoners. In <i>Pro-</i>	ren Usui, Laurent Sifre, Lena Heuermann, Leti-	606
550	<i>ceedings of the 36th International Conference on</i>	cia Lago, Lilly McNealus, Livio Baldini Soares,	607
551	<i>Neural Information Processing Systems, NIPS '22,</i>	Logan Kilpatrick, Lucas Dixon, Luciano Martins,	608
552	Red Hook, NY, USA. Curran Associates Inc.	Machel Reid, Manvinder Singh, Mark Iverson, Mar-	609
553	Justin Matthew Wren Lewis, Andy Williams,	tin Görner, Mat Velloso, Mateo Wirth, Matt Davi-	610
554	Robert Arthur Franklin, James Thomas, and	dow, Matt Miller, Matthew Rahtz, Matthew Watson,	611
555	Nicholas Alexander Mosdell. 2008. The quality and	Meg Risdal, Mehran Kazemi, Michael Moynihan,	612
556	independence of british journalism.	Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi	613
557	Rishabh Misra. 2022. <a href="#">Politifact fact check dataset</a> .	Rahman, Mohit Khatwani, Natalie Dao, Nenshad	614
558	Hannah Rashkin, Eunsol Choi, Jin Yea Jang, Svitlana	Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay	615
559	Volkova, and Yejin Choi. 2017. <a href="#">Truth of varying</a>	Chauhan, Oscar Wahltinez, Pankil Botarda, Parker	616
560	<a href="#">shades: Analyzing language in fake news and po-</a>	Barnes, Paul Barham, Paul Michel, Pengchong	617
561	<a href="#">litical fact-checking</a> . In <i>Proceedings of the 2017</i>	Jin, Petko Georgiev, Phil Culliton, Pradeep Kup-	618
562	<i>Conference on Empirical Methods in Natural Lan-</i>	pala, Ramona Comanescu, Ramona Merhej, Reena	619
563	<i>guage Processing</i> , pages 2931–2937, Copenhagen,	Jana, Reza Ardeshtir Rokni, Rishabh Agarwal, Ryan	620
564	Denmark. Association for Computational Linguis-	Mullins, Samaneh Saadat, Sara Mc Carthy, Sarah	621
565	tics.	Cogan, Sarah Perrin, Sébastien M. R. Arnold, Se-	622
566	Daniel Russo, Serra Sinem Tekiroğlu, and Marco	bastian Krause, Shengyang Dai, Shruti Garg, Shruti	623
567	Guerini. 2023. <a href="#">Benchmarking the Generation of Fact</a>	Sheth, Sue Ronstrom, Susan Chan, Timothy Jordan,	624
568	<a href="#">Checking Explanations</a> . <i>Transactions of the Associa-</i>	Ting Yu, Tom Eccles, Tom Hennigan, Tomas Ko-	625
569	<i>tion for Computational Linguistics</i> , 11:1250–1264.	ciskiy, Tulsee Doshi, Vihan Jain, Vikas Yadav, Vilobh	626
570	Michael Sejr Schlichtkrull, Zhijiang Guo, and Andreas	Meshram, Vishal Dharmadhikari, Warren Barkley,	627
571	Vlachos. 2023. <a href="#">AVeriTeC: A dataset for real-world</a>	Wei Wei, Wenming Ye, Woohyun Han, Woosuk	628
572	<a href="#">claim verification with evidence from the web</a> . In	Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan	629
573	<i>Thirty-seventh Conference on Neural Information</i>	Wei, Victor Cotruta, Phoebe Kirk, Anand Rao, Minh	630
574	<i>Processing Systems Datasets and Benchmarks Track</i> .	Giang, Ludovic Peran, Tris Warkentin, Eli Collins,	631
575	Gemma Team, Morgane Riviere, Shreya Pathak,	Joelle Barral, Zoubin Ghahramani, Raia Hadsell,	632
576	Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupati-	D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov,	633
577	raju, Léonard Hussenot, Thomas Mesnard, Bobak	Oriol Vinyals, Jeff Dean, Demis Hassabis, Koray	634
578	Shahriari, Alexandre Ramé, Johan Ferret, Peter	Kavukcuoglu, Clement Farabet, Elena Buchatskaya,	635
579	Liu, Pouya Tafti, Abe Friesen, Michelle Casbon,	Sebastian Borgeaud, Noah Fiedel, Armand Joulin,	636
580	Sabela Ramos, Ravin Kumar, Charline Le Lan,	Kathleen Kenealy, Robert Dadashi, and Alek An-	637
581	Sammy Jerome, Anton Tsitsulin, Nino Vieillard,	dreev. 2024. <a href="#">Gemma 2: Improving open language</a>	638
582	Piotr Stanczyk, Sertan Girgin, Nikola Momchev,	<a href="#">models at a practical size</a> .	639
583	Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill,	James Thorne, Andreas Vlachos, Christos	640
584	Behnam Neyshabur, Olivier Bachem, Alanna Wal-	Christodoulopoulos, and Arpit Mittal. 2018.	641
585	ton, Aliaksei Severyn, Alicia Parrish, Aliya Ah-	<a href="#">FEVER: a large-scale dataset for fact extraction</a>	642
586	mad, Allen Hutchison, Alvin Abdagic, Amanda	<a href="#">and VERification</a> . In <i>Proceedings of the 2018</i>	643
587	Carl, Amy Shen, Andy Brock, Andy Coenen, An-	<i>Conference of the North American Chapter of</i>	644
588	thony Laforge, Antonia Paterson, Ben Bastian, Bilal	<i>the Association for Computational Linguistics:</i>	645
589	Piot, Bo Wu, Brandon Royal, Charlie Chen, Chintu	<i>Human Language Technologies, Volume 1 (Long</i>	646
590	Kumar, Chris Perry, Chris Welty, Christopher A.	<i>Papers)</i> , pages 809–819, New Orleans, Louisiana.	647
591	Choquette-Choo, Danila Sinopalnikov, David Wein-	Association for Computational Linguistics.	648
592	berger, Dimple Vijaykumar, Dominika Rogozińska,	Andreas Vlachos and Sebastian Riedel. 2014. <a href="#">Fact</a>	649
593	Dustin Herbison, Elisa Bandy, Emma Wang, Eric	<a href="#">checking: Task definition and dataset construction</a> .	650
594	Noland, Erica Moreira, Evan Senter, Evgenii Elty-	In <i>Proceedings of the ACL 2014 Workshop on Lan-</i>	651
595	shev, Francesco Visin, Gabriel Rasskin, Gary Wei,	<i>guage Technologies and Computational Social Sci-</i>	652
596	Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna	<i>ence</i> , pages 18–22, Baltimore, MD, USA. Associa-	653
597	Klimczak-Plucińska, Harleen Batra, Harsh Dhand,	tion for Computational Linguistics.	654
598	Ivan Nardini, Jacinda Mein, Jack Zhou, James Svens-	William Yang Wang. 2017. <a href="#">“liar, liar pants on fire”:</a>	655
599	son, Jeff Stanway, Jetha Chan, Jin Peng Zhou, Joana	<a href="#">A new benchmark dataset for fake news detection</a> .	656
600	Carrasqueira, Joana Iljazi, Jocelyn Becker, Joe Fer-	In <i>Proceedings of the 55th Annual Meeting of the</i>	657
601	nandez, Joost van Amersfoort, Josh Gordon, Josh	<i>Association for Computational Linguistics (Volume 2:</i>	658
		<i>Short Papers)</i> , pages 422–426, Vancouver, Canada.	659
		Association for Computational Linguistics.	660
		Zhiwei Yang, Jing Ma, Hechang Chen, Hongzhan Lin,	661
		Ziyang Luo, and Chang Yi. 2022. <a href="#">A coarse-to-fine</a>	662



cascaded evidence-distillation neural network for explainable fake news detection. In *Proceedings of the 29th International Conference on Computational Linguistics (COLING)*, pages 2608–2621.

Yirong Zeng, Xiao Ding, Yi Zhao, Xiangyu Li, Jie Zhang, Chao Yao, Ting Liu, and Bing Qin. 2024. RU22Fact: Optimizing evidence for multilingual explainable fact-checking on Russia-Ukraine conflict. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 14215–14226, Torino, Italia. ELRA and ICCL.

## A Prompt Selection

In this section, we present the various prompts explored to identify the most effective one for the 5-class fact-checking task. We also report the weighted F1 scores in table 5 for each prompt evaluated on the validation set, providing insight into the performance differences across the prompt variations.

### A.1 Zero Shot Prompts

#### Base Model Prompts

In this section, we provide the seven prompts used for the base model in the zero-shot setting for the 5-class fact-checking task.

P1 Given claim and evidence, predict if the claim is true, mostly-true, half-true, mostly-false, or false.  
claim: {{claim}}  
evidence: {{evidence}}  
label:

P2 Given the evidence, decide if the given claim is true, mostly-true, half-true, mostly-false, or false.  
claim: {{claim}}  
evidence: {{evidence}}  
label:

P3 Given claim and evidence, find if the claim is true, mostly-true, half-true, mostly-false, or false.  
claim: {{claim}}  
evidence: {{evidence}}  
label:

P4 Identify if the claim is true, mostly-true, half-true, mostly-false, or false based on the evidence.  
claim: {{claim}}  
evidence: {{evidence}}  
label:

P5 Given claim and evidence, classify if the claim is true, mostly-true, half-true, mostly-false, or false.  
claim: {{claim}}  
evidence: {{evidence}}  
label:

P6 You need to determine the accuracy of a claim based on the evidence. Use one of following 5 labels for the claim: true, mostly-true, half-true, mostly-false, or false. Examine the evidence and choose the most likely label based on the claim’s accuracy without explaining your reasoning.  
claim: {{claim}}  
evidence: {{evidence}}  
label:

P7 Given claim and evidence, you are tasked with evaluating the truthfulness of claims based on the provided evidence. Each claim can be categorized into one of 5 labels: true, mostly-true, half-true, mostly-false, false. Assess the claim given the evidence and classify it appropriately without providing an explanation.  
claim: {{claim}}  
evidence: {{evidence}}  
label:

#### Mistral Instruct Models Prompts

In this section, we provide the seven prompts used for the Mistral instruct model in the zero-shot setting for the 5-class fact-checking task.

P1 <s>[INST] You are a helpful AI assistant, and you are tasked with evaluating the truthfulness of claims based on the provided evidence. Each claim can be categorized into one of 5 labels: "true", "mostly-true", "half-true", "mostly-false", "false". Assess the claim given the evidence and classify it appropriately without providing an explanation. [/INST] I am excited to work on this classification problem. Can you please provide me with the label description for all 5 labels? [/INST][Label Descriptions]



761	true: The claim is accurate and	mostly-true, half-true, mostly-false,	812
762	includes all relevant information.	false.	813
763	There are no omissions or distortions	</s> [INST]	814
764	that could mislead the audience.	claim: {{claim}}	815
765	mostly-true: The claim is accurate,	evidence: {{evidence}}	816
766	but it might benefit from additional	label: [/INST]	817
767	context to provide a complete picture.		
768	However, the absence of this context	P3 <s>[INST] You need to judge the truth	818
769	does not alter the claim's accuracy.	of a claim based on the evidence	819
770	half-true: The claim is true	given. Use one of these 5 labels	820
771	in a limited context. However,	for each claim: true, mostly-true,	821
772	it omits crucial information	half-true, mostly-false, or false.	822
773	that could significantly alter	Review the evidence and classify	823
774	its interpretation, leading to	the claim without explaining your	824
775	potential misunderstanding or	reasoning. [/INST]	825
776	misinterpretation.	Now, can you please provide me with	826
777	mostly-false: The claim contains	a claim and evidence so that based	827
778	some elements of truth but distorts	on the evidence I can classify the	828
779	or misrepresents critical facts.	claim into one of the 5 labels: true,	829
780	Important information is omitted,	mostly-true, half-true, mostly-false,	830
781	which could lead to a misleading	false.	831
782	impression despite some truthful	</s> [INST]	832
783	elements.	claim: {{claim}}	833
784	false: The claim is inaccurate and	evidence: {{evidence}}	834
785	contradicts established facts. The	label: [/INST]	835
786	claim has no truth, and it is likely		
787	to mislead those who encounter it.	P4 <s> Given claim and evidence, you	836
788	[/INST] Now, can you please provide	are tasked with evaluating the	837
789	me with a claim and evidence so that	truthfulness of claims based on	838
790	based on the evidence I can classify	the provided evidence. Each claim	839
791	the claim into one of the 5 labels:	can be categorized into one of 5	840
792	"true", "mostly-true", "half-true",	labels: true, mostly-true, half-true,	841
793	"mostly-false", "false".	mostly-false, false. Assess the claim	842
794	[/INST]	given the evidence and classify it	843
795	claim: {{claim}}	appropriately without providing an	844
796	evidence: {{evidence}}	explanation.	845
797	label: [/INST]	claim: {{claim}}	846
		evidence: {{evidence}}	847
		label:	848
798	P2 <s>[INST] Given claim and evidence,		
799	you are tasked with evaluating the	P5 <s> Given a claim and evidence, you	849
800	truthfulness of claims based on	need to decide how accurate a claim is	850
801	the provided evidence. Each claim	based on the evidence given. Select	851
802	can be categorized into one of 5	one of the five labels to classify the	852
803	labels: true, mostly-true, half-true,	claim: true, mostly-true, half-true,	853
804	mostly-false, false. Assess the claim	mostly-false, or false. Review the	854
805	given the evidence and classify it	evidence, decide how well it supports	855
806	appropriately without providing an	the claim, and then pick the best	856
807	explanation. [/INST]	label for the truthfulness of the	857
808	Now, can you please provide me with	claim.	858
809	a claim and evidence so that based	claim: {{claim}}	859
810	on the evidence I can classify the	evidence: {{evidence}}	860
811	claim into one of the 5 labels: true,	label:	861

862	P6 <s> You need to determine the accuracy	Thoroughly review the evidence and	910
863	of a claim based on the evidence. Use	accurately categorize the claim	911
864	one of the following 5 labels for the	without explaining your decision.	912
865	claim: true, mostly-true, half-true,	claim: {{claim}}	913
866	mostly-false, or false. Examine the	evidence: {{evidence}}	914
867	evidence and choose the most likely	label:	915
868	label based on the claim's accuracy		
869	without explaining your reasoning.	P4 You need to determine the accuracy of	916
870	claim: {{claim}}	a claim based on the evidence.	917
871	evidence: {{evidence}}	Use one of the following 5 labels	918
872	label:	for each claim: true, mostly-true,	919
		half-true, mostly-false, or false.	920
873	P7 <s> Given claim and evidence, find	Examine the evidence and pick the most	921
874	if the claim is true, mostly-true,	probable label for the claim without	922
875	half-true, mostly-false, or false.	explaining your reasoning.	923
876	claim: {{claim}}	claim: {{claim}}	924
877	evidence: {{evidence}}	evidence: {{evidence}}	925
878	label:	label:	926
879	<b>Llama/Gemma Instruct Models Prompts</b>		
880	In this section, we provide the seven prompts used	P5 You need to determine the accuracy of	927
881	for the LLaMA/Gemma instruct model in the zero-	a claim based on the evidence.	928
882	shot setting for the 5-class fact-checking task.	Use one of the following 5 labels	929
		for each claim: true, mostly-true,	930
883	P1 You need to judge the truth of a claim	half-true, mostly-false, or false.	931
884	based on the evidence given.	Examine the evidence and pick the	932
885	Use one of these 5 labels for each	most probable label according to the	933
886	claim: true, mostly-true, half-true,	truthfulness of the claim without	934
887	mostly-false, or false.	explaining your reasoning.	935
888	Review the evidence and classify	claim: {{claim}}	936
889	the claim without explaining your	evidence: {{evidence}}	937
890	reasoning.	label:	938
891	claim: {{claim}}		
892	evidence: {{evidence}}	P6 You need to determine the accuracy of	939
893	label:	a claim based on the evidence.	940
		Use one of the following 5 labels	941
894	P2 You need to decide how accurate a	for the claim: true, mostly-true,	942
895	claim is based on the evidence given.	half-true, mostly-false, or false.	943
896	Use one of these 5 labels to classify	Examine the evidence and choose the	944
897	each claim: true, mostly-true,	most likely label based on the	945
898	half-true, mostly-false, or false.	claim's accuracy without explaining	946
899	Read the evidence, decide how well it	your reasoning.	947
900	supports the claim, and then pick the	claim: {{claim}}	948
901	best label.	evidence: {{evidence}}	949
902	claim: {{claim}}	label:	950
903	evidence: {{evidence}}		
904	label:	P7 Given claim and evidence, you	951
		are tasked with evaluating the	952
905	P3 Determine the validity of a claim	truthfulness of claims based on the	953
906	using the provided evidence.	provided evidence.	954
907	Select one of the following 5	Each claim can be categorized into	955
908	labels: true, mostly-true, half-true,	one of 5 labels: true, mostly-true,	956
909	mostly-false, or false.	half-true, mostly-false, false.	957
		Assess the claim given the evidence	958

and classify it appropriately without providing an explanation.  
 claim: {{claim}}  
 evidence: {{evidence}}  
 label:

## A.2 Few Shot Prompts

### Base/Instruct Models Prompts

In this section, we provide the seven prompts used for Base/Instruct models in the few-shot setting for the 5-class fact-checking task.

P1 You need to determine the accuracy of a claim based on the evidence. Use one of the following 5 labels for each claim: true, mostly-true, half-true, mostly-false, or false. Examine the evidence and pick the most probable label according to the truthfulness of the claim without explaining your reasoning.  
 claim: {{claim}}  
 evidence: {{evidence}}  
 label:

P2 You need to judge the truth of a claim based on the evidence given. Use one of these 5 labels for each claim: true, mostly-true, half-true, mostly-false, or false. Review the evidence and classify the claim without explaining your reasoning.  
 claim: {{claim}}  
 evidence: {{evidence}}  
 label:

P3 Given claim and evidence, you are tasked with evaluating the truthfulness of claims based on the provided evidence. Each claim can be categorized into one of 5 labels: true, mostly-true, half-true, mostly-false, or false. Assess the claim given the evidence and classify it appropriately without providing an explanation.  
 claim: {{claim}}  
 evidence: {{evidence}}  
 label:

P4 Given claim and evidence, find if the claim is true, mostly-true,

Label	Count	Token <sub>μ</sub>	Sent <sub>μ</sub>	BPE <sub>μ</sub>
True	296	666.20	29.40	847.27
Mostly True	298	804.78	36.51	1022.50
Half True	293	898.12	39.95	1135.38
Mostly False	300	927.16	41.40	1165.60
False	295	679.12	33.05	862.42
<b>Total</b>	1482	795.31	36.07	1006.92

Table 4: Statistics for Unfiltered version of Politi-fact-only dataset. Token<sub>μ</sub>, Sent<sub>μ</sub>, and BPE<sub>μ</sub> represent the average number of standard tokens, sentences, and BPE tokens per entry, respectively.

half-true, mostly-false, or false.  
 claim: {{claim}}  
 evidence: {{evidence}}  
 label:

P5 Based on the provided evidence, verify the claim and classify it as true, mostly-true, half-true, mostly-false, or false.  
 claim: {{claim}}  
 evidence: {{evidence}}  
 label:

P6 Based on the provided evidence, judge whether the claim is true, mostly-true, half-true, mostly-false, or false.  
 claim: {{claim}}  
 evidence: {{evidence}}  
 label:

P7 Examine the evidence and classify the claim as true, mostly-true, half-true, mostly-false, or false.  
 claim: {{claim}}  
 evidence: {{evidence}}  
 label:

## B Guidelines

I met with the annotators regularly over a span of two months. During this time, we employed three annotators who were proficient in English and compensated by our lab. The Politifact-only dataset is a fact-checking dataset scraped from politifact.com, focusing on the political domain. It consists of 1,500 instances, each containing a political claim along with corresponding evidence. Based on the evidence, the claim’s truth value is categorized in one of the following categories: true, mostly true, half true, mostly false, false, pants on fire. I have



Zero Shot							
	P1	P2	P3	P4	P5	P6	P7
Base Models							
Mistral-7B-v0.3	0.3213	0.3213	0.3199	0.3396	0.3415	<b>0.4253</b>	0.4147
Llama-3-8B	0.29	0.4607	0.4891	0.4678	0.4468	<b>0.5202</b>	0.4781
Gemma-2-9b	0.2979	0.3180	0.3264	0.3494	0.3094	0.3473	<b>0.3769</b>
Instruct Models							
Mistral-7B-Instruct-v0.3	0.5191	0.5334	0.4060	<b>0.5428</b>	0.4832	0.5419	0.5066
Llama-3-8B-Instruct	0.6132	0.4249	0.3550	0.6239	<b>0.6276</b>	0.6240	0.4207
Gemma-2-9b-it	0.5183	0.3837	0.4281	0.4041	0.4041	0.3979	<b>0.5512</b>
Few(5) Shot							
Base Models							
Mistral-7B-v0.3	0.7690	0.7567	0.7587	<b>0.7809</b>	0.7618	0.7778	0.7785
Llama-3-8B	0.6984	0.7123	0.6883	0.7251	0.7304	0.7044	<b>0.7365</b>
Gemma-2-9b	0.6566	0.6552	0.6073	0.6914	<b>0.7127</b>	<b>0.7127</b>	0.6990
Instruct Models							
Mistral-7B-Instruct-v0.3	0.6867	0.6989	0.6856	<b>0.7360</b>	0.7215	0.7350	0.7332
Llama-3-8B-Instruct	0.4387	0.4433	0.4908	<b>0.5505</b>	0.5235	0.5235	0.5120
Gemma-2-9b-it	0.3700	0.4009	0.3774	0.3625	0.3867	<b>0.3889</b>	0.3585
2-stage CoT							
	P1	P2	P3	P4	P5	P6	
Mistral-7b-v0.3-instruct	<b>0.5317</b>	0.4129	0.4339	0.4180	0.4957	0.4604	

Table 5: Weighted F1 Scores for Different Prompts Across Various Models and Experiment Methodologies (Zero-Shot, Few-Shot, and Two-Stage CoT). The scores are reported for multiple prompt configurations for base and instruct models, demonstrating performance variations in prompt selection.

clubbed pants on fire and false into one label that is false Table 1

First, we need to clean, test and val set first so that we can make use of the dataset for the experiment, then we will move on to train the dataset.

Fields in the dataset for instance: Id, label, speaker, claim, evidence, source, speaker, claim\_data, etc given in the provided json file. I added “leaked” and “annotator prediction” for you to fill in. Label Description: True: The statement is accurate and there’s nothing significant missing. Mostly True: The statement is accurate but needs clarification or additional information Half True: The statement is partially accurate but leaves out important details or takes things out of context. Mostly False: The statement contains an element of truth but ignores critical facts that would give a different impression. False: The statement is not accurate. Pants on fire: The statement is not accurate and makes a ridiculous claim.

## **B.1 Problem with the current dataset**

The dataset contains the claim and corresponding evidence that supports or refutes the claim. We have some leakage in our dataset. Leakage means, there are evidences in the dataset which are giving away the information about the label of the corresponding claim. The evidence may contain the definition of the label, some direct intuition about the label, or the label itself.

*What needs to be done:*

Remove our ruling section if exists.

Remove the sentence that contains a label or label definition.

Remove sentences that directly give away the information about the label.

Remove redundant conclusions by the annotator if they repeat information from the previous section or can be inferred from prior content.

Mark which evidence needed changes in the “leaked” field by writing yes/no. Give your predicted label in the “annotator prediction” field.

Examples with strike-throughs helped annotators understand what to remove. We met weekly and re-annotate if needed.

"label": "true"  
"claim": "At nearly 19 million people, the population of Florida is larger than all the earlier primary and caucus states combined."  
"evidence": "Gov. Rick Scott rallied Republican activists at Florida's presidential primary straw poll with an argument for the state's supremacy in choosing the party's presidential contender. None will have a greater impact on the selection of the nominee than our own primary in the Sunshine State, Scott told a crowd of 3,500 on Sept. 24, 2011. While other primaries or caucuses might be earlier, he said, Florida's population and diversity set it apart. At nearly 19 million people, the population of Florida is larger than all the earlier primary and caucus states combined, he said. The Republican National Committee allows just Iowa, New Hampshire, South Carolina and Nevada to vote in February 2012 without penalty. Florida has yet to choose its primary date. But state lawmakers would like to see it as early as possible, saying it better reflects the country than the four early states and should play an agenda-setting role."  
"speaker": "Rick Scott"  
"claim\_date": "9/24/2011"  
"source": "speech"  
"factchecker": "Becky Bowers"  
"factcheck\_date": "9/27/2011"  
"factcheck\_analysis\_link": "https://www.politifact.com/factchecks/2011/sep/27/rick-scott/gov-rick-scotts-primary-math-florida-has-more-peop/"

Figure 2: A true instance from the dataset.



"label": "mostly-true"  
"claim": "The failings in our civil service are encouraged by a system that makes it very difficult to fire someone even for gross misconduct."  
"evidence": "Sen. John McCain, the Arizona Republican, overstates the problem of removing federal employees for poor performance, but not by much, according to experts who examine federal work rules. It is perhaps not a surprise that a union official disputes McCain's use of the incompetent federal worker cliché. Procedures do exist to remove workers from their jobs, and many people do get fired. But it takes a long time, according to the outside experts who follow such issues closely. McCain wisely faults not an individual but a system. That puts him on pretty solid ground, where even a study by the federal government had difficulty finding supervisors who had attempted to take action against poorly performing employees."  
"speaker": "John McCain"  
"claim\_date": "3/21/2007"  
"source": "other"  
"factchecker": "Angie Drobnic Holan"  
"factcheck\_date": "9/1/2007"  
"factcheck\_analysis\_link": "https://www.politifact.com/factchecks/2007/sep/01/john-mccain/you-can-fire-federal-workers-but-its-tough/"

Figure 3: A mostly true instance from the dataset.

"label": "half-true"

"claim": "21-million Americans could have a four-year college scholarship for the money we've squandered in Iraq. 7.6-million teachers could have been hired last year if we weren't squandering this money."

"evidence": "Former U.S. Sen. Mike Gravel attacked the Iraq war during a recent debate by highlighting the increasing costs. Stop and think, he said at Howard University on June 28, 2007. When he's talking about the money we're squandering, 21-million Americans could have a four-year college scholarship for the money we've squandered in Iraq. 7.6-million teachers could have been hired last year if we weren't squandering this money. Gravel's campaign staff didn't respond to numerous requests for documentation supporting those numbers. They couldn't even say how much they think the Iraq war costs. The college board puts the average cost of tuition for a four-year public university in 2006 at \$5,836. Do the math: the sum exceeds \$490.2-billion, much higher than even the highest estimate. The U.S. Department of Education reports the average teacher salary was \$47,750 in 2005, the most recent year available. That produces a total of \$363-billion, well below the lowest estimate. The Congressional Budget Office conservatively estimates the entire bill for the Iraq war since 2001 is \$413-billion."

"speaker": "Mike Gravel"

"claim\_date": "6/28/2007"

"source": "other"

"factchecker": "John Frank"

"factcheck\_date": "9/20/2007"

"factcheck\_analysis\_link":

"https://www.politifact.com/factchecks/2007/sep/20/mike-gravel/hes-high-then-hes-low/"

Figure 4: A half true instance from the dataset.

"label": "mostly-false"

"claim": "Photo shows a semi-truck that crashed with a Chevy pickup that cut in front of it."

"evidence": "An unnerving photo of a vehicle crumpled under a semi-truck is being shared on social media with a warning: the next time you decide to cut in front of that 80,000 lb semi, remember: this was once a 4-door chevy pickup. In september 2016, wsb-tv, a news station in atlanta, aired images from a crash involving four tractor-trailers on interstate 20 in carroll county, georgia. Georgia state patrol said at the time that a tractor-trailer ran into the back of a second tractor-trailer, according to the station. The second tractor-trailer then drove over a silver pickup truck, crushing it, and ran into a third tractor-trailer. The third tractor trailer then hit a fourth one. The person driving the pickup and a passenger were killed in the crash. The atlanta journal-constitution reported the same narrative. The deadly chain reaction started when a tractor-trailer headed eastbound struck a second tractor-trailer, which then struck the silver pickup truck, killing the driver and passenger, the newspaper said. The photo suggests this is what happens to smaller vehicles that cut in front of big trucks on the highway. But this was no ordinary collision; it involved multiple vehicles that are not all pictured."

"speaker": "Viral image"

"claim\_date": "6/15/2021"

"source": "social\_media"

"factchecker": "Ciara O'Rourke"

"factcheck\_date": "6/21/2021"

"factcheck\_analysis\_link": "<https://www.politifact.com/factchecks/2021/jun/21/viral-image/crash-photo-doesnt-show-vehicle-cut-front-semi-tru/>"

Figure 5: A mostly false instance from the dataset.



"label": "false"  
"claim": "A 2022 video shows Ukrainian and Russian soldiers face to face."  
"evidence": "Footage of soldiers firing shots into the air as hundreds of unarmed people march toward an airbase in belbek, crimea, is being shared on tiktok as russia invades ukraine. Ukrainian and russian soldiers face off in big battle border, one post sharing the footage wrote. #ukrainian and #ryssland soldiers face to face, another post said. This footage was posted over 12 times on tiktok and viewed on the platform more than 20 million times as of feb. 25. Bbc news turkey shared the footage on youtube on march 4, 2014. According to the bbc article, the video depicts pro-russian troops who seized an airbase firing warning shots to prevent some 300 unarmed ukrainian soldiers from approaching. The tense standoff occurred as russia annexed crimea in 2014."  
"speaker": "TikTok posts"  
"claim\_date": "2/25/2022"  
"source": "blog"  
"factchecker": "Yacob Reyes"  
"factcheck\_date": "2/25/2022"  
"factcheck\_analysis\_link": "https://www.politifact.com/factchecks/2022/feb/25/tiktok-posts/video-standoff-between-soldiers-ukraine-2014/"

Figure 6: A false instance from the dataset.