

Fact-checking AI-generated news reports: Can LLMs catch their own lies?

Anonymous ACL submission

Abstract

In this paper, we evaluate the ability of Large Language Models (LLMs) to assess the veracity of claims in “news reports” generated by themselves or other LLMs. Our goal is to determine whether LLMs can effectively fact-check their own content, using methods similar to those used to verify claims made by humans. Our findings indicate that LLMs are more effective at assessing claims in national or international news stories than in local news stories, better at evaluating static information than dynamic information, and better at verifying true claims compared to false ones. We hypothesize that this disparity arises because the former types of claims are better represented in the training data. Additionally, we find that incorporating retrieved results from a search engine in a Retrieval-Augmented Generation (RAG) setting significantly reduces the number of claims an LLM cannot assess. However, this approach also increases the occurrence of incorrect assessments, partly due to irrelevant or low-quality search results. This diagnostic study highlights the need for future research on fact-checking machine-generated reports to prioritize improving the precision and relevance of retrieved information to better support fact-checking efforts. Furthermore, claims about dynamic events and local news may require human-in-the-loop fact-checking systems to ensure accuracy and reliability.

1 Introduction

Large Language Models (LLMs) have revolutionized the field of Natural Language Processing (NLP), effortlessly performing tasks that were traditionally considered highly challenging. Their performance is particularly impressive in generating natural language text. Models like GPT-4 can generate coherent, fluent summaries, accurately translate text between languages (especially those with a strong online presence and ample training

data), and refine human writing to enhance fluency and appropriateness in tone and style for specific purposes. This technology has the potential to significantly increase productivity across many industries, offering endless applications. However, with this potential also come risks if they are not used properly. One of the main risks is that they can be easily used to generate convincing and yet factually incorrect text, either intentionally or unintentionally. For example, with a simple prompt like “Generate a news report about volcano eruption in Massachusetts, USA”, GPT-4 can generate a news report starting with this first paragraph:

“Massachusetts, USA – May 29, 2024 – In an unprecedented and shocking event, a volcanic eruption has occurred in the state of Massachusetts, an area not typically associated with volcanic activity. The eruption took place early this morning in the central part of the state, near the town of Worcester, sending residents and scientists alike into a state of disbelief and concern.”

Although there has never been a volcanic eruption in reality, the news report is coherent and fluent. Coupled with modern media platforms, such LLM-generated content can quickly spread and reach a large audience. An example is the emergence of AI “news” farms that produce news reports with LLMs to generate advertising revenue with little concern for their impact on society (Puccetti et al., 2024). The machine-generated reports can cause confusion and chaos and disrupt the proper functioning of the society. In fact, studies show that false news tends to spread “farther, faster, deeper” than true news, as it often contains novel content that people are more likely to share (Vosoughi et al., 2018).

In this study, we present experimental results to answer the questions of whether LLMs are capable of telling if the news stories they generate

are truthful and how well they can catch factually incorrect claims in those news stories. We generated 92 news stories with a simple prompt such as “Write a story about Kobe Bryant rejoining the Lakers” with two LLMs, GPT-4o (OpenAI et al., 2024) and GLM (Du et al., 2022), all stories with some incorrect claims. These false stories vary by how untruthful they are. Some stories report events that are simply impossible, such as a story about Kobe Bryant rejoining the Lakers, as the former Lakers star has passed away. Some stories report events that are not out of the realm of possibilities but are highly unlikely, such as a volcano eruption in Massachusetts, as the area is not known for volcanic activities. Other stories are about events that have actually happened or are scheduled to happen, but with the wrong time, location, or participants.

We performed our experiments in two settings. In the first setting, we simply provided the full story to GPT-4o and GLM as input and asked if they are truthful. In the second setting, we manually decompose each story into individual checkable *atomic* claims. A checkable claim can be either an event with specific participants, location, or time, or they can be a state (e.g., Massachusetts borders New Hampshire) or recurring event that holds for an extended period of time such that the exact time is irrelevant. We perform manual “decontextualization” (Choi et al., 2021) on these checkable claims so that they can be verified outside of the context of their document. For this setting, we also experimented with using these checkable claims as queries, providing the results retrieved via the Google Search Serper Api¹ to GPT-4 to assist in evaluating the veracity of the claims within a Retrieval-Augmented Generation (RAG) framework (Lewis et al., 2020).

The results of our experiments show that GPT-4o and GLM are very good at detecting stories that contain incorrect claims (and all of them do) when they involve well-known entities (e.g., Kobe Bryant rejoining the Lakers), but they are quite uncertain about recent events that are unlikely. At the level of atomic claims, a significant proportion of them are incorrectly assessed: either a factually correct claim is judged to be wrong or a factually wrong claim is identified as being correct. For an even larger proportion of atomic claims, the LLMs simply cannot decide. When provided with results retrieved via the Google Search Serper Api, the number of non-assessments decreases significantly,

accompanied by an increase in both correct and incorrect assessments. Interestingly, even when the Google Search Serper Api returns no results for a claim, GPT-4 still attempts to provide an assessment instead of declining to answer. It appears that simply knowing no results were retrieved is enough to prompt GPT-4 to make a guess. Even with RAG, there is still a significant proportion of claims that the LLM cannot provide an assessment. This means that any solution to fact-checking machine-generated news reports needs to include functionalities on checking claims about new event occurrences that are not checkable against existing knowledge sources. While there has been recent research that shows the promise of using external resources or tools to improve the factuality of LLMs (Gou et al., 2023), such an approach is not applicable to fact-checking machine-generated news stories and novel human-in-the-loop methods may need to be developed to check such claims.

The rest of the paper is organized as follows. In Section 2, we discuss related work. In Section 3, we present our method for generating news stories, extract “atomic” claims, using LLMs to assess the veracity of these stories and claims, and manually verifying the assessments performed by LLMs themselves. We present experimental results in Section 4, and discuss these results in Section 5. We conclude in Section 6.

2 Related Work

Fact-checking human or machine-generated content. There is an active NLP research community focused on developing automatic methods to fact-check false claims, such as those made by politicians (Nakov et al., 2021; Deng et al., 2024; Yuan and Vlachos, 2024; Schlichtkrull et al., 2024). There is also more recent work on fact-checking machine-generated content (Min et al., 2023; Wang et al., 2024; Fadeeva et al., 2024). Previous work on fact-checking false claims made by either humans or machines typically assume there is an information source, usually a published source on the Internet, against which the claims can be checked. However, events reported in machine-generated news stories that we are interested in, such as the volcano eruption example, are often assumed to be new occurrences that cannot be cross-verified against any existing public sources, although they may still contain claims about the real world that can be fact-checked. This poses novel challenges

¹<https://serper.dev/>

that are not present in people biographies used in previous studies (Min et al., 2023; Fadeeva et al., 2024).

3 Method

Our experiment on fact-checking LLM-generated news stories consists of four steps. First, we use two LLMs to generate a set of news stories with varying levels of factual inaccuracy. Next, from these stories, we manually extract verifiable atomic claims and decontextualize them, creating standalone claims that can be verified independently of the original story. In the third step, we prompt each LLM to evaluate the veracity of news stories generated by itself or the other LLM, as well as to assess the individual atomic claims. Finally, we conduct a human evaluation to determine the accuracy of the LLMs’ veracity assessments.

3.1 News Report Generation with LLMs

To evaluate the claim verification capabilities of GPT and GLM, we first prompt both models to generate a set of 92 news articles, including 47 news articles generated by GPT-4o and 44 articles generated by GLM. Each prompt is designed around scenario-based inputs that intentionally contain factual inconsistencies. The following is an example prompt that contains a time error, as the time of 2024 Australian Open women’s final is January 27, not January 20:

"Generate a news report about Aryna Sabalenka winning the 2024 Australian Open Women’s final, held at Rod Laver Arena on January 20, as Aryna Sabalenka beat Zheng Qinwen (6-3, 6-2)."

All these inconsistencies are designed around four critical aspects of a scenario: the event itself, along with its time, location, and participants. To rigorously test the models’ understanding of both nationally recognized and locally relevant information, we control the scope of the generated content by introducing both local and national news categories. The distinction between these categories serves as a critical factor in our evaluation, allowing us to evaluate how effectively each model handles claims involving specific local information versus those based on widely known national knowledge. This is motivated by prior research suggesting that LLMs may have greater exposure to widely discussed national or international events, given the

nature of the large, diverse datasets they are trained on (Kandpal et al., 2023). When generating the news stories, we ensure that the same general template is used for all prompts, varying only the scenarios for each different story. By using consistent prompts, we ensure that differences in model performance can be attributed to the model’s capabilities rather than variability in the inputs. This approach allows us to build a diverse and representative dataset that rigorously tests each LLM’s ability to identify and evaluate issues across different aspects of the generated content.

3.2 Manual claim Extraction

After generating the news reports, we manually extracted all checkable claims from the GPT-generated content. Each claim is a clear, verifiable statement with specific details such as time, location, participants, or events. We adhered to criteria that required each checkable claim to contain precise, unambiguous information—such as exact dates, locations, or identifiable participants. Vague or generic statements, like “Sabalenka had a great match” were excluded, as they lack objective, verifiable details. This approach ensured that only claims containing concrete, factual information were selected for manual extraction. We manually decontextualize claims by resolving pronominal and other anaphoric expressions, and by supplementing events with time, location, and participant details when they are clear from the context, ensuring that each claim is independently verifiable.

The following are example claims illustrating various types of factual inaccuracies:

- **Time error:** “Aryna Sabalenka triumphed over Zheng Qinwen to win the 2024 Australian Open Women’s final at Rod Laver Arena on January 20, 2024.”
- **Location error:** “Aryna Sabalenka played against Zheng Qinwen in the 2024 Australian Open Women’s final at Margaret Court Arena on January 27.”
- **Event error:** “In the third set of the 2024 Australian Open Women’s Final at Rod Laver Arena on January 27, Zheng Qinwen broke Aryna Sabalenka’s serve at 5-5 and won the set 7-5 to clinch the championship.”
- **Participant and location error:** “Naomi Osaka and Iga Swiatek are battling for the prestigious Grand Slam title at

280
281
282

283
284
285
286
287
288
289
290
291
292
293

294
295
296
297
298
299
300
301

302
303
304
305
306
307
308
309
310
311

312
313
314
315
316
317
318

319
320

321
322
323
324
325

the 2024 Australian Open Women’s Final at the Margaret Court Arena on January 27, 2024.”

Each article typically yields between 10-20 checkable decontextualized claims, depending on its length and complexity. This process ensures that the claims include all the necessary contextual information required for verification, maintaining the integrity and relevance of the claims within the broader context of the news reports. From the 92 articles we have extracted 1,337 total atomic claims, including 697 claims from the 47 news reports generated by GPT-4o, and 640 claims from the 44 reports generated by GLM.

3.3 Claim verification with LLMs

Both GPT-4o and GLM models are tasked with verifying the veracity of each entire article as well as each atomic claim. To assess claim veracity, we prompted GPT-4o and GLM to evaluate the accuracy of all 92 news articles and their corresponding atomic claims. The following are the prompts we use for the evaluation:

- **Article-level prompt:** “Today is August 1st, 2024. You are a helpful assistant that performs the below tasks: verify if the following news is accurate or false. Respond as concisely as possible.”
- **Claim-level prompt:** “Today is August 1st, 2024. You are a helpful assistant that performs the below tasks: verify if the following claim extracted from a news report is accurate or false. Respond as concisely as possible.”

The models are first prompted to assess the veracity of each entire article and provide a rationale for their evaluations. They are then prompted to evaluate the veracity of each atomic claim extracted from the articles, along with a rationale for each assessment. Three different prompting approaches are used in this pipeline.

3.3.1 Deterministic Prompting (Temperature 0.0)

We prompt the models to provide a singular, deterministic evaluation for each article or claim. Setting temperature to 0 minimizes randomness and allows us to observe the models’ baseline claim verification performance under controlled conditions.

3.3.2 Self-consistency Prompting (Temperature 1.0)

We use a higher temperature setting (1.0) to introduce variability in the responses of the models. Models are prompted multiple times (5 times per article / claim in our experiment), and a majority voting mechanism is used to determine the final assessment. This setting simulates the potential variability in model reasoning and robustness across multiple prompts.

In each instance, the model outputs a determination (correct or false) along with a rationale for its assessment. These rationales are crucial for error analysis, offering insights into whether the model’s reasoning aligns with the factual basis of the claim.

3.3.3 RAG Prompting

. We queried the Google Search Serper Api with manually extracted atomic claims and incorporated the retrieved results into the prompt for GPT-4 when evaluating the veracity of claims within a Retrieval-Augmented Generation (RAG) framework. The goal of this experiment was to assess whether providing search results improves the evaluation accuracy of LLMs. Due to cost constraints and the length limitation of the search engine, we did not perform this experiment with the entire article. Instead, we focused on atomic claims extracted from news reports generated by GPT-4 itself, assuming the results would generalize to other settings.

3.4 Comparing model verification with human judgments

To validate the models’ evaluations, we manually verify each claim by conducting targeted web searches and cross-referencing the findings with our existing information. We use independent online sources, including reputable news databases, fact-checking websites, and government records. The human judgments serve as the gold standard for evaluating model assessments, enabling us to quantify both false positives and false negatives in the models’ evaluations. Additionally, we performed error analysis to understand whether the type of news (local vs. national) had a measurable impact on the model’s performance. Special attention was paid to cases where the models provided no assessment, incorrect reasoning, or inaccurate evaluations.

326
327

328
329
330
331
332
333
334
335
336
337
338
339
340

341
342
343
344
345
346
347
348
349
350
351
352
353
354
355

356
357

358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373

4 Experiments

We conduct a comprehensive set of experiments to evaluate the performance of GPT-4o and GLM models in verifying claims within generated news articles. Both models are assessed in the contexts of local and national news generation, with claim verification performed across all relevant dimensions. For the claim verification task, we classify the assessment results into five possible categories, as outlined below:

- **Correct Assessment (CA):** The model correctly identifies the veracity of the claim without providing a rationale.
- **Correct Assessment and Correct Reasoning (CA/CR):** The model correctly identifies the veracity of the claim and provides a correct justification for its assessment.
- **Correct Assessment and Wrong Reasoning (CA/WR):** The model correctly classifies the claim but with flawed reasoning.
- **Wrong Assessment (WA):** The model incorrectly classifies the veracity of the claim.
- **No Assessment (NA):** The model fails to provide any assessment.

| Generator Evaluator | GPT-4o | | | GLM | | |
|------------------------|--------|-------------|-------|--------|-------------|-------|
| | GPT-4o | GPT-4-turbo | GLM-4 | GPT-4o | GPT-4-turbo | GLM-4 |
| CA | 1 | 1 | 0 | 1 | 0 | 0 |
| CA/CR | 30 | 26 | 37 | 31 | 31 | 31 |
| CA/WR | 5 | 2 | 1 | 5 | 6 | 2 |
| WA | 9 | 8 | 8 | 2 | 2 | 10 |
| NA | 2 | 10 | 1 | 6 | 6 | 2 |
| Total | 47 | 47 | 47 | 45 | 45 | 45 |

Table 1: Count of LLM-generated articles for each assessment category

4.1 Entire news articles

Table 1 presents the performance data of GPT-4 (gpt-4o-20240806 and gpt-4-turbo-20240409) and GLM-4 (GLM-4-0520) in evaluating entire articles. Both models were prompted to generate news reports, followed by self-evaluation and cross-evaluation of the generated articles.

GPT-4 and GLM-4 demonstrate comparable performance in the number of correct and incorrect assessments they produce. In contrast, GPT-4 Turbo is more likely to refrain from making assessments, reflecting a more cautious approach compared to GPT-4o and GLM-4. This suggests that GPT-4-turbo prioritizes minimizing errors, even if it results in fewer overall judgments.

4.2 Individual atomic claims

In evaluating LLMs in verifying atomic claims, we conducted experiments with GPT-4o and GLM-4 to ensure our findings are generalizable across LLMs. The performance of GPT and GLM models was assessed across different temperature settings to better assess their strengths and limitations in claim verification tasks. Both models were tasked with verifying the veracity of claims extracted from LLM-generated news articles, with their evaluations measured using the identical 5-dimensional protocol we use for entire articles.

The assessment results are presented in Table 2 and we can make several key observations. First, GPT-4o consistently provides more correct assessments (including those with and without correct reasoning) than GLM, regardless of whether it is evaluating claims from articles it generated or those generated by GLM. This trend holds across all temperature settings. Interestingly, both GPT-4o and GLM produce more incorrect assessments (WA) when evaluating claims from articles they generated themselves. The most notable finding is the high number of cases with no assessment (NA), with GLM showing a significantly higher number (about 20%) of no assessments than GPT-4.

4.2.1 Claims in National vs Local news stories

We also attempted to evaluate the ability of LLMs to assess claims in national and local news stories. The following are example claims from national and local news stories we generated with LLMs:

- **Claims in local news:** The free rave hosted by Watertown, MA on July 15, 2024 will be held at Arsenal Park.
- **Claims in national or international news:** The 2024 Paris Olympics opening ceremony is set to take place on July 26.

Table 3 presents a comparative error analysis of GPT and GLM models when evaluating claims from national and local news sources, across different temperature settings. Errors in assessments include cases where the model provides the correct assessment with wrong reasoning (CA/WR), wrong assessment (WA), or no assessment (NA). As we can see from the table, while GPT slightly outperforms GLM as indicated by the generally lower number of errors, the error rate is relatively consistent across temperatures.

The most notable finding is the substantial difference in error rates between the models’ assess-

| Generator | GPT-4o | | | | GLM | | | |
|-----------|------------|------------|------------|------------|-------------|------------|------------|------------|
| Evaluator | GPT/0 | GPT/1 | GLM/0 | GLM/1 | GPT/0 | GPT/1 | GLM/0 | GLM/1 |
| CA (%) | 38(5.45) | 44(6.31) | 1(0.14) | 0(0.00) | 15(2.34) | 14(2.19) | 3(0.47) | 5(0.78) |
| CA/CR(%) | 306(43.90) | 291(41.75) | 271(38.88) | 276(39.60) | 349(54.53) | 353(55.16) | 312(48.75) | 306(47.81) |
| CA/WR(%) | 5(0.72) | 10(1.43) | 13(1.87) | 15(2.15) | 42(6.56) | 31(4.84) | 24(3.75) | 24(3.75) |
| WA(%) | 33(4.73) | 42(6.03) | 12(1.72) | 14(2.00) | 9(1.40) | 12(1.88) | 15(2.34) | 29(4.53) |
| NA(%) | 315(45.19) | 310(44.48) | 400(57.39) | 392(56.24) | 225 (35.16) | 230(35.94) | 286(44.69) | 276(43.13) |
| Total | 697 | 697 | 697 | 697 | 640 | 640 | 640 | 640 |

Table 2: Count and percentage of individual atomic claims for each assessment category across models at different temperature settings. GPT/0 and GPT/1 indicate GPT at temperature 0 and 1 respectively. Similarly, GLM/0 and GLM/1 indicate GLM at temperature 0 and 1.

| Generator | GPT-4o | | | | | GLM | | | | |
|-------------|--------|------------|------------|------------|------------|-------|------------|------------|------------|------------|
| Evaluator | Subt. | GPT/0 | GPT/1 | GLM/0 | GLM/1 | Subt. | GPT/0 | GPT/1 | GLM/0 | GLM/1 |
| National(%) | 496 | 193(38.91) | 208(41.94) | 252(50.81) | 247(49.80) | 462 | 143(30.95) | 141(30.52) | 194(41.99) | 197(42.64) |
| Local(%) | 201 | 160(79.60) | 154(76.62) | 173(86.07) | 174(86.57) | 178 | 133(74.72) | 132(74.16) | 131(73.60) | 132(74.16) |
| Total | 697 | 353(50.65) | 362(51.94) | 425(60.98) | 421(60.40) | 640 | 276(43.13) | 273(42.66) | 325(50.78) | 329(51.41) |

Table 3: Errors from evaluating claims in national or local news. Each cell represents the percentage of claims that are incorrectly assessed for that category (national vs local), with the last row representing the number of errors / the total claims for that generator.

ments of claims from national and local news, with significantly higher error rates for local news. One possible explanation is that claims in national news often pertain to major events or widely recognized topics that are well-documented across diverse on-line sources, making these claims more likely to appear in the models’ training data and thus easier to assess. In contrast, claims in local news may involve niche, region-specific issues that receive limited attention and documentation, leaving the models less prepared to verify such claims accurately. This discrepancy highlights how the scope and distribution of training data can impact the models’ performance in evaluating claims with different degrees of specificity and familiarity.

4.2.2 Assessment of true claims vs false claims

Table 4 evaluates the accuracy of LLMs in assessing both factually correct and wrong claims. We analyze whether the LLMs make accurate or inaccurate assessments when presented with claims that are either true or false. Correct Assessment includes cases where (i) the claim is factually true, and the LLM assesses it as true. (ii) The claim is factually false, and the LLM assesses it as false. And wrong assessment includes cases where (i) the claim is factually false, but the LLM assesses it as true and (ii) the claim is factually true, but the LLM assesses it as false. We aim to investigate whether there is a difference in the accuracy

with which LLMs assess factually true versus false claims. Our hypothesis is that factually true claims are more likely to be represented in the training data than factually false ones, making it more probable that factually false claims will be incorrectly assessed. Our hypothesis is born out, as results in Table 4 show that both the GPT and GLM generally have a higher rate of correct assessments when the claim was factually correct while both models struggle with factually wrong claims and made wrong assessments. Among all the cases where the model made correct assessments but provided incorrect reasoning, a considerable portion of them is from claims that are factually wrong. This suggests that while the model can arrive at the correct conclusion, its internal logic or justifications may still be flawed, which happens mostly when the claims are factually incorrect.

4.2.3 State and event claims

We also experimented with asking LLMs to assess claims that are linguistic states and those that are not. Here, a state refers to a specific condition or phase in the existence of something, characterized by stability and consistency over time, whereas a non-state claim typically involves an event, signifying a significant occurrence that brings about change. A non-state claim is typically associated with a time, location, and participants. The following shows example claims categorized as state and

| Generator | GPT-4 | | | | GLM | | | |
|-----------|----------|----------|----------|----------|---------|----------|---------|----------|
| | GPT/0 | | GLM/0 | | GPT/0 | | GLM/0 | |
| Evaluator | FC(%) | FW(%) | FC(%) | FW(%) | FC(%) | FW(%) | FC(%) | FW(%) |
| Veracity | | | | | | | | |
| CA (CR) | 143 (87) | 201 (38) | 135 (82) | 137 (26) | 92 (89) | 272 (51) | 86 (84) | 229 (43) |
| CAWR | 0 (0.0) | 5 (0.9) | 1 (0.6) | 12 (2.3) | 2 (1.9) | 40 (7.4) | 0 (0.0) | 24 (4.5) |
| WA | 4 (2.4) | 29 (5.5) | 4 (2.4) | 8 (1.5) | 1 (1.0) | 8 (7.8) | 5 (4.9) | 10 (1.9) |
| NA | 18 (11) | 297 (56) | 25 (15) | 375 (71) | 8 (7.8) | 217 (40) | 12 (12) | 274 (51) |
| Total | 165 | 532 | 165 | 532 | 103 | 537 | 103 | 537 |

Table 4: Comparison of LLM assessment accuracy for factually correct (FC) and factually incorrect (FW) claims with GPT and GLM as evaluators at 0 temperature.

non-state:

- **State claim:** Aryna Sabalenka is Belarusian.
- **Non-state claim:** The 2024 Australian Open Women’s final was held at Margaret Court Arena on January 27.

We hypothesize that LLMs perform better on state claims because states are more stable and likely to be documented in training data, whereas events are often new and undocumented. Consequently, LLMs are more prone to errors, including wrong assessments (WA) and no assessments (NA), when evaluating non-state claims, as supported by the higher error rates observed for these claims. This hypothesis is largely born out by the higher error rate for non-states than states. We also observed a significant temperature effect and found that higher temperatures yield better results for state claims, potentially due to improved pattern recognition from broad, consistent data, while for non-state claims, the same high temperatures lead to worse outcomes as they inhibit the verification of event-specific details, causing increased uncertainty and wrong assessments. More information about this can be found in Appendix A.3.

4.2.4 Fact-checking with Retrieval Augmented Generation (RAG)

Retrieval-Augmented Generation (RAG) (Lewis et al., 2020) has emerged as a popular method for fact-checking (Rothermel et al., 2024; Khaliq et al., 2024; Raina and Gales, 2024; Ullrich et al., 2024; Adjali, 2024), particularly when LLMs struggle to find information relevant to a given claim. The process typically involves transforming the claim into questions that can be used to query a knowledge source, such as the entire Internet or specific repositories like Wikipedia. The retrieved results, combined with the original claim, are then used to

prompt an LLM to determine whether the claim is supported or refuted by the evidence. Additionally, the LLM can conclude that there is insufficient evidence to either support or refute the claim.

In the RAG approach, each claim is treated as a search query to retrieve relevant supporting or contradictory information from the Internet. Specifically, the claim is then fed into a Serper API to fetch relevant results from online sources. The results are then filtered to ensure relevance. For textual search results, the first $k = 5$ entries are selected, prioritizing those with detailed snippets, titles, and links. For knowledge graph data, attributes like titles, entity types, and descriptions are processed into usable snippets. The retrieved snippets and contextual data are consolidated and formatted into a coherent input prompt for GPT-4o. See Appendix A.2 for an example prompt.

The assessment results using the RAG approach are shown in Table 5. Compared to the non-RAG setting, the number of correct assessments (CACR) increases significantly by 20%, but the number of wrong assessments (WA) also rises by 8.3%, from 4.7% to 13%. Meanwhile, the number of no assessments (NA) drops dramatically, from 45% to 16%. These results suggest that when augmented with retrieval results, GPT-4o adopts a more aggressive approach in making assessments.

Interestingly, GPT often provides a "No Assessment" (NA) response even when retrieved search results (S) are available. This occurs when the LLM determines that the retrieved evidence is insufficiently relevant to support a definitive evaluation. Conversely, GPT-4o is capable of making correct assessments even when no relevant evidence is retrieved. A possible explanation may lie in the structure of the prompt given to the LLM. The sentence "Here are the related search snippets" followed by

| Generator | GPT-4 | | | | | | | |
|-----------|---------------|----------|----------|-----------|----------|----------|----------|----------|
| | GPT/0 Non-RAG | | | GPT/0 RAG | | | | |
| Evaluator | Subt. (%) | FC (%) | FW (%) | Subt. (%) | FC (%) | FW (%) | S (%) | NS (%) |
| CA (CR) | 344 (49) | 143 (87) | 201 (38) | 482 (69) | 111 (67) | 371 (70) | 424 (73) | 58 (51) |
| CAWR | 5 (0.7) | 0 (0.0) | 5 (0.9) | 11 (1.6) | 1 (0.6) | 10 (1.9) | 5 (0.9) | 6 (5.3) |
| WA | 33 (4.7) | 4 (2.4) | 29 (5.5) | 92 (13) | 24 (15) | 68 (13) | 82 (14) | 10 (0.9) |
| NA | 315 (45) | 18 (11) | 297 (56) | 112 (16) | 29 (18) | 83 (16) | 72 (12) | 40 (35) |
| Total | 697 | 165 | 532 | 697 | 165 | 532 | 583 | 114 |

Table 5: Comparison between RAG and non-RAG performance with GPT4-o at Temperature 0. “S” indicates search results are returned by the Google Serper API and “NS” means no results are returned.

an empty list might implicitly signal to the LLM that no evidence supports the claim, prompting it to guess that the claim is false. However, it is debatable whether we want the LLM to make guesses this way when acting as a fact-checking system, where credibility is paramount.

5 Discussion

In our evaluation of LLMs’ ability to assess the veracity of LLM-generated news articles and claims, we find that LLMs perform better when evaluating claims in national news compared to local news. They are also more accurate at assessing factually correct claims than factually wrong ones. Additionally, LLMs excel at evaluating claims expressed as linguistic states rather than those describing dynamic events. These seemingly distinct observations can be traced back to a common underlying factor: LLMs are more effective at processing well-documented, high-frequency information that is more likely to have been included in their training data. National news claims are typically better documented than local news claims, linguistic states are more stable and frequently recorded than rapidly evolving dynamic events, and factually accurate claims are more likely to appear in the training data than factually false ones. Using RAG significantly increases the level of correct assessments, but it also leads to a higher number of wrong assessments due to irrelevant search results (55 out of 92 cases), no search results (10 out of 92 cases), or wrong reasoning (27 out of 92 cases). There is still a significant number of no assessments (NA) even with the RAG approach, either because no search results are retrieved or the search results are noisy and irrelevant. RAG systems also have the tendency to venture guesses even in the absence of evidential support, and this is problematic even if the guess is correct. This underscores the

need for future research on fact-checking machine-generated news content to prioritize the retrieval of precise and reliable evidence. For claims the retrieval system cannot find evidence for, human-in-the-loop approaches may need to be developed to ensure accuracy and reliability.

Our study uses claims that are manually extracted and decontextualized. Fully automatic evaluation systems would require that the atomic claims are automatically extracted and decontextualized, with the goal of extracting *all and only* checkable claims from an LLM-generated text. This is especially challenging for news stories, which may contain vague and subjective language. For automatic fact-checking systems to gain the trust and confidence of users, it is critical for them to be transparent and interpretable.

6 Conclusion and Future Work

We conducted a diagnostic study to evaluate the strengths and limitations of using LLMs and RAG systems for fact-checking claims in machine-generated “news” reports. While these systems can verify the veracity of a significant portion of claims (nearly 70%), a considerable number are either incorrectly assessed or left unassessed due to irrelevant retrievals, flawed reasoning, or insufficient evidence. This issue is particularly pronounced for rare claims with limited evidential support, which are common in news reports. Our findings underscore the need for more precise and reliable retrieval systems and the incorporation of human-in-the-loop approaches when evidence is unavailable. Future work will explore the ability of LLMs to generate verifiable claims, a crucial step toward fully automated fact-checking systems.

Limitations

In this diagnostic study, we relied on manually extracted claims, which inherently limits the size of the dataset and, consequently, the breadth of the analysis. The manual extraction process is time-consuming and labor-intensive, making it challenging to scale the dataset to include a larger number of claims. Despite this limitation, we carefully curated the dataset to ensure it is representative of the types of claims commonly found in machine-generated news reports. As a result, we are confident that the dataset is sufficiently large and diverse to support reliable and meaningful conclusions.

Ethical Statement

Machine-generated news reports can pose significant risks if they are mistaken for authentic, factual content. To mitigate these risks, when releasing the dataset for our study, we will ensure that it is clearly labeled as machine-generated and explicitly highlight that it contains false claims. This labeling is critical to prevent misuse of the dataset and to maintain transparency for researchers, developers, and the broader community. By doing so, we aim to promote ethical research practices and minimize any potential harm arising from the dissemination of this data.

In the NLP community, it is common practice to release datasets publicly by hosting them on open-source platforms like GitHub. However, in this case, it is more appropriate to store the data on a private server and provide access to fellow researchers upon request. This approach is preferable for two key reasons. First, releasing the data on an open-source platform risks it being incorporated into the training data of future LLM versions, rendering results non-comparable. Second, the dataset is primarily useful to researchers and serves little to no practical purpose for the general public.

References

Omar Adjali. 2024. [Exploring retrieval augmented generation for real-world claim verification](#). In *Proceedings of the Seventh Fact Extraction and VERification Workshop (FEVER)*, pages 113–117, Miami, Florida, USA. Association for Computational Linguistics.

Eunsol Choi, Jennimaria Palomaki, Matthew Lamm, Tom Kwiatkowski, Dipanjan Das, and Michael Collins. 2021. [Decontextualization: Making sentences stand-alone](#). *Transactions of the Association for Computational Linguistics*, 9:447–461.

Zhenyun Deng, Michael Schlichtkrull, and Andreas Vlachos. 2024. [Document-level claim extraction and de-contextualisation for fact-checking](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 11943–11954, Bangkok, Thailand. Association for Computational Linguistics.

Zhengxiao Du, Yujie Qian, Xiao Liu, Ming Ding, Jiezhong Qiu, Zhilin Yang, and Jie Tang. 2022. [GLM: General language model pretraining with autoregressive blank infilling](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 320–335, Dublin, Ireland. Association for Computational Linguistics.

Ekaterina Fadeeva, Aleksandr Rubashevskii, Artem Shelmanov, Sergey Petrakov, Haonan Li, Hamdy Mubarak, Evgenii Tsymbalov, Gleb Kuzmin, Alexander Panchenko, Timothy Baldwin, Preslav Nakov, and Maxim Panov. 2024. [Fact-checking the output of large language models via token-level uncertainty quantification](#). In *Findings of the Association for Computational Linguistics ACL 2024*, pages 9367–9385, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.

Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujia Yang, Nan Duan, and Weizhu Chen. 2023. [Critic: Large language models can self-correct with tool-interactive critiquing](#). *ArXiv*, abs/2305.11738.

Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. 2023. [Large language models struggle to learn long-tail knowledge](#). *Preprint*, arXiv:2211.08411.

Mohammed Abdul Khaliq, Paul Yu-Chun Chang, Mingyang Ma, Bernhard Pflugfelder, and Filip Milić. 2024. [RAGAR, your falsehood radar: RAG-augmented reasoning for political fact-checking using multimodal large language models](#). In *Proceedings of the Seventh Fact Extraction and VERification Workshop (FEVER)*, pages 280–296, Miami, Florida, USA. Association for Computational Linguistics.

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.

Sewon Min, Kalpesh Krishna, Xinxin Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. [FActScore: Fine-grained atomic evaluation of factual precision in long form text generation](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12076–12100, Singapore. Association for Computational Linguistics.

| | | |
|-----|---|-----|
| 775 | Preslav Nakov, David Corney, Maram Hasanain, Firoj | |
| 776 | Alam, Tamer Elsayed, Alberto Barrón-Cedeño, Paolo | |
| 777 | Papotti, Shaden Shaar, and Giovanni Da San Martino. | |
| 778 | 2021. Automated fact-checking for assisting human | |
| 779 | fact-checkers . In <i>Proceedings of the Thirtieth Inter-</i> | |
| 780 | <i>national Joint Conference on Artificial Intelligence,</i> | |
| 781 | <i>IJCAI-21</i> , pages 4551–4558. International Joint Con- | |
| 782 | ferences on Artificial Intelligence Organization. Sur- | |
| 783 | vey Track. | |
| 784 | OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, | |
| 785 | Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale- | |
| 786 | man, Diogo Almeida, Janko Altmenschmidt, Sam Alt- | |
| 787 | man, Shyamal Anadkat, Red Avila, Igor Babuschkin, | |
| 788 | Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim- | |
| 789 | ing Bao, Mohammad Bavarian, Jeff Belgum, Ir- | |
| 790 | wan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, | |
| 791 | Christopher Berner, Lenny Bogdonoff, Oleg Boiko, | |
| 792 | Madeline Boyd, Anna-Luisa Brakman, Greg Brock- | |
| 793 | man, Tim Brooks, Miles Brundage, Kevin Button, | |
| 794 | Trevor Cai, Rosie Campbell, Andrew Cann, Brittany | |
| 795 | Carey, Chelsea Carlson, Rory Carmichael, Brooke | |
| 796 | Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully | |
| 797 | Chen, Ruby Chen, Jason Chen, Mark Chen, Ben | |
| 798 | Chess, Chester Cho, Casey Chu, Hyung Won Chung, | |
| 799 | Dave Cummings, Jeremiah Currier, Yunxing Dai, | |
| 800 | Cory Decareaux, Thomas Degry, Noah Deutsch, | |
| 801 | Damien Deville, Arka Dhar, David Dohan, Steve | |
| 802 | Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, | |
| 803 | Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, | |
| 804 | Simón Posada Fishman, Juston Forte, Isabella Ful- | |
| 805 | ford, Leo Gao, Elie Georges, Christian Gibson, Vik | |
| 806 | Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo- | |
| 807 | Lopes, Jonathan Gordon, Morgan Grafstein, Scott | |
| 808 | Gray, Ryan Greene, Joshua Gross, Shixiang Shane | |
| 809 | Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, | |
| 810 | Yuchen He, Mike Heaton, Johannes Heidecke, Chris | |
| 811 | Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, | |
| 812 | Brandon Houghton, Kenny Hsu, Shengli Hu, Xin | |
| 813 | Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, | |
| 814 | Joanne Jang, Angela Jiang, Roger Jiang, Haozhun | |
| 815 | Jin, Denny Jin, Shino Jomoto, Billie Jonn, Hee- | |
| 816 | woo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka- | |
| 817 | mali, Ingmar Kanitscheider, Nitish Shirish Keskar, | |
| 818 | Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, | |
| 819 | Christina Kim, Yongjik Kim, Jan Hendrik Kirch- | |
| 820 | ner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, | |
| 821 | Łukasz Kondraciuk, Andrew Kondrich, Aris Kon- | |
| 822 | stantinidis, Kyle Kosic, Gretchen Krueger, Vishal | |
| 823 | Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan | |
| 824 | Leike, Jade Leung, Daniel Levy, Chak Ming Li, | |
| 825 | Rachel Lim, Molly Lin, Stephanie Lin, Mateusz | |
| 826 | Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, | |
| 827 | Anna Makanju, Kim Malfacini, Sam Manning, Todor | |
| 828 | Markov, Yaniv Markovski, Bianca Martin, Katie | |
| 829 | Mayer, Andrew Mayne, Bob McGrew, Scott Mayer | |
| 830 | McKinney, Christine McLeavey, Paul McMillan, | |
| 831 | Jake McNeil, David Medina, Aalok Mehta, Jacob | |
| 832 | Menick, Luke Metz, Andrey Mishchenko, Pamela | |
| 833 | Mishkin, Vinnie Monaco, Evan Morikawa, Daniel | |
| 834 | Mossing, Tong Mu, Mira Murati, Oleg Murk, David | |
| 835 | Mély, Ashvin Nair, Reiichiro Nakano, Rameev Nayak, | |
| 836 | Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, | |
| 837 | Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex | |
| | Paino, Joe Palermo, Ashley Pantuliano, Giambat- | 838 |
| | tista Parascandolo, Joel Parish, Emy Parparita, Alex | 839 |
| | Passos, Mikhail Pavlov, Andrew Peng, Adam Perel- | 840 |
| | man, Filipe de Avila Belbute Peres, Michael Petrov, | 841 |
| | Henrique Ponde de Oliveira Pinto, Michael, Poko- | 842 |
| | rny, Michelle Pokrass, Vitchyr H. Pong, Tolly Pow- | 843 |
| | ell, Alethea Power, Boris Power, Elizabeth Proehl, | 844 |
| | Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, | 845 |
| | Cameron Raymond, Francis Real, Kendra Rimbach, | 846 |
| | Carl Ross, Bob Rotsted, Henri Roussez, Nick Ry- | 847 |
| | der, Mario Saltarelli, Ted Sanders, Shibani Santurkar, | 848 |
| | Girish Sastry, Heather Schmidt, David Schnurr, John | 849 |
| | Schulman, Daniel Selsam, Kyla Sheppard, Toki | 850 |
| | Sherbakov, Jessica Shieh, Sarah Shoker, Pranav | 851 |
| | Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, | 852 |
| | Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin | 853 |
| | Sokolowsky, Yang Song, Natalie Staudacher, Fe- | 854 |
| | lipe Petroski Such, Natalie Summers, Ilya Sutskever, | 855 |
| | Jie Tang, Nikolas Tezak, Madeleine B. Thompson, | 856 |
| | Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, | 857 |
| | Preston Tuggle, Nick Turley, Jerry Tworek, Juan Fe- | 858 |
| | lipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, | 859 |
| | Chelsea Voss, Carroll Wainwright, Justin Jay Wang, | 860 |
| | Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, | 861 |
| | CJ Weinmann, Akila Welihinda, Peter Welinder, Ji- | 862 |
| | ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, | 863 |
| | Clemens Winter, Samuel Wolrich, Hannah Wong, | 864 |
| | Lauren Workman, Sherwin Wu, Jeff Wu, Michael | 865 |
| | Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qim- | 866 |
| | ing Yuan, Wojciech Zaremba, Rowan Zellers, Chong | 867 |
| | Zhang, Marvin Zhang, Shengjia Zhao, Tianhao | 868 |
| | Zheng, Juntang Zhuang, William Zhuk, and Bar- | 869 |
| | ret Zoph. 2024. Gpt-4 technical report . <i>Preprint</i> , | 870 |
| | arXiv:2303.08774. | 871 |
| | Giovanni Puccetti, Anna Rogers, Chiara Alzetta, Felice | 872 |
| | Dell’Orletta, and Andrea Esuli. 2024. Ai "news" | 873 |
| | content farms are easy to make and hard to detect: A | 874 |
| | case study in italian . <i>Preprint</i> , arXiv:2406.12128. | 875 |
| | Vatsal Raina and Mark Gales. 2024. Question-based | 876 |
| | retrieval using atomic units for enterprise RAG . In | 877 |
| | <i>Proceedings of the Seventh Fact Extraction and VER-</i> | 878 |
| | <i>ification Workshop (FEVER)</i> , pages 219–233, Miami, | 879 |
| | Florida, USA. Association for Computational Lin- | 880 |
| | guistics. | 881 |
| | Mark Rothmel, Tobias Braun, Marcus Rohrbach, and | 882 |
| | Anna Rohrbach. 2024. Infact: A strong baseline | 883 |
| | for automated fact-checking. In <i>Proceedings of the</i> | 884 |
| | <i>Seventh Fact Extraction and VERification Workshop</i> | 885 |
| | <i>(FEVER)</i> , pages 108–112. | 886 |
| | Michael Schlichtkrull, Yulong Chen, Chenxi White- | 887 |
| | house, Zhenyun Deng, Mubashara Akhtar, Rami Aly, | 888 |
| | Zhijiang Guo, Christos Christodoulopoulos, Oana | 889 |
| | Cocarascu, Arpit Mittal, James Thorne, and Andreas | 890 |
| | Vlachos. 2024. The automated verification of textual | 891 |
| | claims (AVeriTeC) shared task . In <i>Proceedings of</i> | 892 |
| | <i>the Seventh Fact Extraction and VERification Work-</i> | 893 |
| | <i>shop (FEVER)</i> , pages 1–26, Miami, Florida, USA. | 894 |
| | Association for Computational Linguistics. | 895 |
| | Herbert Ullrich, Tomáš Mlynář, and Jan Drchal. 2024. | 896 |
| | AIC CTU system at AVeriTeC: Re-framing auto- | 897 |

mated fact-checking as a simple RAG task. In *Proceedings of the Seventh Fact Extraction and VERification Workshop (FEVER)*, pages 137–150, Miami, Florida, USA. Association for Computational Linguistics.

Soroush Vosoughi, Deb Roy, and Sinan Aral. 2018. The spread of true and false news online. *Science*, 359(6380):1146–1151.

Yuxia Wang, Revanth Gangi Reddy, Zain Mujahid, Arnav Arora, Aleksandr Rubashevskii, Jiahui Geng, Osama Mohammed Afzal, Liangming Pan, Nadav Borenstein, Aditya Pillai, et al. 2024. Factcheck-bench: Fine-grained evaluation benchmark for automatic fact-checkers. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 14199–14230.

Moy Yuan and Andreas Vlachos. 2024. Zero-shot fact-checking with semantic triples and knowledge graphs. In *Proceedings of the 1st Workshop on Knowledge Graphs and Large Language Models (KaLLM 2024)*, pages 105–115, Bangkok, Thailand. Association for Computational Linguistics.

A Appendix

A.1 Example claims

1. TotalEnergies is the title sponsor for the TotalEnergies BWF Thomas & Uber Cup Finals 2024.
2. Three separate shark attacks have been reported off the coast of Maine from June 30 to July 4, 2024.
3. The United States has reported a 10% growth in Gross Domestic Product (GDP) for the fiscal year 2024 on May 29, 2024.
4. The TotalEnergies BWF Thomas & Uber Cup Finals 2024 was held at the Chongqing Olympic Sports Center.
5. The opening ceremony of the 2024 Summer Olympics was held at the Bangkok Olympic Stadium on May 29, 2024.
6. The 2024 Australian Open Women’s Final was held at the Margaret Court Arena on January 27, 2024.
7. Spain won over France in the 2024 UEFA European Championship semi-final at the BVB Stadion Dortmund on July 9, 2024.
8. Zheng Qinwen is competing for her first Grand Slam final at the 2024 Australian Open Women’s final at Rod Laver Arena on January 20, 2024.
9. The free rave hosted by Watertown, MA on July 15, 2024 will be held at Arsenal Park.

10. The discovery of the new COVID-19 in New Hampshire variant was announced by health officials on May 29, 2024.

11. The Dallas Mavericks won Game 4 of the 2024 NBA Finals in overtime at the American Airlines Center in Dallas, Texas.

12. The 2024 NBA Finals Game 7 was played at the American Airlines Center in Dallas on May 29, 2024.

13. Kobe Bryant has announced his return to the Los Angeles Lakers on May 29, 2024.

14. On July 3, 2024, FIFA has announced that London, United Kingdom, will host the 2026 FIFA World Cup.

15. On May 29, 2024, Stephen Curry was traded from the Golden State Warriors to the Chicago Bulls.

A.2 Example RAG prompt

The prompt structure includes: The original claim. The concatenated evidence snippets from the retrieved results. And a preamble describing the task (e.g., assessing the factual accuracy of the claim). Example Prompt: "The following claim needs to be evaluated for accuracy: 'CLAIM'." "Here are the related search snippets: EVIDENCE-TEXT." "Based on the snippets provided, evaluate whether the claim is accurate or false. " "Provide a clear and reasoned explanation."

A.3 Assessments for State and Event Claims

Table 6 presents a comparison between the evaluation performance of LLMs on state and non-state (event) claims. LLMs are better at assessing state claims than non-state claims, as indicated by the generally lower number of WA and NA cases for state claims and higher number of such cases for non-state claims.

There is also a significant temperature effect. For state claims, which often pertain to more standardized and systemic issues, higher temperatures might enhance the model’s ability to identify patterns and make accurate assessments. These claims are typically based on broader, more consistent data that may not be as sensitive to small fluctuations or variability in the input data. Conversely, higher temperatures introduce greater variability in responses, which impacts non-state claims differently. Non-state claims defined by more dynamic, event-specific details like timing, location, or participants, become harder for the models to verify

| Generator | GPT-4o | | | | GLM | | | |
|------------------|--------|-------|-------|-------|-------|-------|-------|-------|
| | GPT/1 | GPT/0 | GLM/1 | GLM/0 | GPT/1 | GPT/0 | GLM/1 | GLM/0 |
| State | | | | | | | | |
| CA | 29 | 14 | 0 | 1 | 14 | 6 | 2 | 0 |
| CA/CR | 121 | 81 | 145 | 85 | 99 | 68 | 100 | 71 |
| CA/WR | 0 | 2 | 1 | 1 | 0 | 15 | 0 | 4 |
| WA | 8 | 3 | 4 | 3 | 1 | 4 | 12 | 5 |
| NA | 30 | 88 | 38 | 98 | 23 | 44 | 23 | 57 |
| CA (%) | 15.4 | 7.4 | 0 | 0.5 | 10.2 | 4.3 | 1.4 | 0 |
| CA/CR (%) | 64.4 | 43.1 | 77.1 | 45.2 | 42.2 | 49.6 | 72.9 | 51.8 |
| CA/WR (%) | 0 | 1.1 | 0.5 | 0.5 | 0 | 10.9 | 0 | 2.9 |
| WA (%) | 4.2 | 1.6 | 2.1 | 1.6 | 0.7 | 2.9 | 8.7 | 3.6 |
| NA (%) | 15.9 | 46.8 | 20.2 | 52.1 | 16.7 | 32.1 | 16.7 | 41.6 |
| Subtotal | 188 | 188 | 188 | 188 | 137 | 137 | 137 | 137 |
| Non-State | | | | | | | | |
| CA | 15 | 24 | 0 | 0 | 0 | 9 | 3 | 3 |
| CA/CR | 170 | 225 | 131 | 186 | 254 | 281 | 206 | 241 |
| CA/WA | 10 | 3 | 14 | 12 | 31 | 27 | 24 | 20 |
| WA | 34 | 30 | 10 | 9 | 11 | 5 | 17 | 10 |
| NA | 280 | 227 | 354 | 302 | 207 | 181 | 253 | 229 |
| CA (%) | 2.9 | 4.7 | 0 | 0 | 0 | 1.8 | 0.6 | 0.6 |
| CA/CR (%) | 33.4 | 44.2 | 25.7 | 36.5 | 50.5 | 55.9 | 50.9 | 47.9 |
| CA/WA (%) | 1.9 | 0.5 | 2.8 | 2.4 | 6.2 | 5.4 | 4.8 | 4.9 |
| WA (%) | 6.7 | 5.9 | 1.9 | 1.8 | 2.2 | 0.9 | 4.4 | 1.9 |
| NA (%) | 55.0 | 44.6 | 69.5 | 59.3 | 41.1 | 35.9 | 50.3 | 45.5 |
| Subtotal | 509 | 509 | 509 | 509 | 503 | 503 | 503 | 503 |
| Total | 697 | 697 | 697 | 697 | 640 | 640 | 640 | 640 |

Table 6: A comparison of state vs non-state claims.

with confidence under higher temperatures. The randomness at this setting leads to the model producing a broader array of responses, which is beneficial for creativity but not ideal for precision. In fact, the variability might cause the model to contradict itself or lose consistency, particularly when precise details are required to confirm an event. This can explain the higher no assessments for non-state claims under high-temperature settings, as the models struggle with conflicting or incomplete information about specific events.