

# *One Thousand and One Pairs:* A “novel” challenge for long-context language models

Marzena Karpinska<sup>☾</sup> Katherine Thai<sup>☾</sup> Kyle Lo<sup>✍</sup> Tanya Goyal<sup>☆</sup> Mohit Iyyer<sup>☾</sup>  
<sup>☾</sup>UMass Amherst <sup>✍</sup>Allen Institute for AI <sup>☆</sup>Cornell University  
{mkarpinska,kbthai,miiyer}@umass.edu kylel@allenai.org tanyagoyal@cornell.edu

## Abstract

Synthetic long-context LLM benchmarks (e.g., “needle-in-the-haystack”) test only surface-level retrieval capabilities, but how well can long-context LLMs retrieve, synthesize, and reason over information across book-length inputs? We address this question by creating NOCHA, a dataset of 1,001 *minimally different* pairs of true and false claims about 67 recently-published English fictional books, written by human readers of those books. In contrast to existing long-context benchmarks, our annotators confirm that the largest share of pairs in NOCHA require *global reasoning over the entire book* to verify. Our experiments show that while human readers easily perform this task, it is enormously challenging for all ten long-context LLMs that we evaluate: no open-weight model performs above random chance (despite their strong performance on synthetic benchmarks), while GPT-4o achieves the highest accuracy at 55.8%. Further analysis reveals that (1) on average, models perform much better on pairs that require only sentence-level retrieval vs. global reasoning; (2) model-generated explanations for their decisions are often inaccurate even for correctly-labeled claims; and (3) models perform substantially worse on speculative fiction books that contain extensive world-building. The methodology proposed in NOCHA allows for the evolution of the benchmark dataset and the easy analysis of future models.

## 1 Introduction

The context size of large language models has increased by multiple orders of magnitude over the last year: for instance, Google’s GEMINI PRO 1.5 can process millions of input tokens at once. But can models *truly* utilize and reason over their claimed context? Existing long-context evaluation methods such as finding the “needle in the haystack” (NIAH) (Kamradt, 2023) measure surface-level retrieval capabilities, but do not necessarily assess

performance on the more challenging task of synthesizing distant and underlying information as we show in §4.

We bridge this gap by introducing NOCHA (A **Novel Challenge**), in which LLMs are prompted to verify claims written about recently-published fiction books. Claim verification has been extensively studied in the context of shorter documents (Thorne et al., 2018; Wadden et al., 2020; Fabbri et al., 2022a), but its application to book-length fictional texts presents unique challenges. Firstly, the task necessitates reasoning over both explicit information directly stated in the text and implicit information inferred from the narrative, which is often distributed throughout the entire document. Secondly, the use of recent, fictional context prevents the model from relying solely on parametric knowledge, necessitating the comprehension and interpretation of the long text.

Our data collection process aims to balance efficiency and quality. Rather than pre-selecting a set of books, we follow the approach of Kim et al. (2024) and ask annotators to self-report recently published novels they have read. The annotators then create true/false **narrative minimal pairs** that isolate a single narrative phenomenon present in their novels. Each false claim differs from the true claim in its pair *only* by the inclusion of false information regarding the same event or entity (see Figure 1). This approach offers two key advantages: (1) it minimizes the chances that the model is “correct for the wrong reason,” as it must accurately predict *both* labels in the pair, and (2) it simplifies the process of quality control, as the false claim can be easily verified against its true counterpart, making it easier to identify claims that are either too similar or overly subjective. NOCHA contains 1,001 narrative minimal pairs for 67 books, created at a total cost of \$3,330 USD.

Accurately labeling the minimal pairs in NOCHA often requires not only information retrieval from

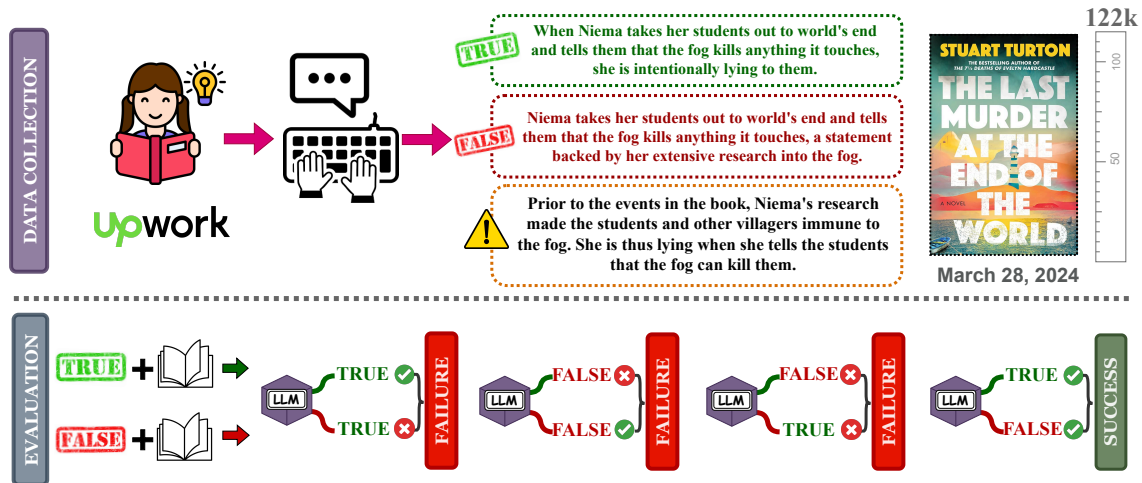


Figure 1: Overview of NOCHA’s data collection and evaluation pipeline. Readers create true/false claim pairs for books published between 2023 and 2024 with written justifications for the false labels. Each model is given the full book as context and evaluates one claim at a time. We measure *pair accuracy*, where the model must identify the true claim as true and the false claim as false to receive credit. This approach helps guard against label bias while also better assessing the true understanding of the text, as both claims pertain to the same events or parts of the story. Our books range from 49k to 336k tokens, and each model is tested on a subset of books that fit within its context.

the context, but also synthesis, reasoning, and inference over information from multiple parts of the document. Experiments on 5 openly-available and 6 closed-source models show that: (1) though all models struggle on NOCHA, no open-weight model performs above random chance, while GPT-4o is the best-performing model overall at 55.8% accuracy; (2) on average, all models perform much better on pairs that require sentence-level retrieval than global reasoning (59.8% vs 41.6%), though their performance on these pairs in NOCHA is much lower (59.8%) than on NIAH (100%) as reported in Hsieh et al. (2024); (3) model-generated explanations for their decisions are often inaccurate, even for correctly labeled claims; and (4) models perform substantially worse on books with extensive world-building, as these require more sophisticated reasoning.

To summarize, our contributions are threefold:

1. **DATA:** We introduce NOCHA, a dataset of 1,001 narrative minimal pairs about *recently* published fictional books, designed to evaluate the long-context language models’ ability to process and reason about long narratives.
2. **ANALYSIS:** We use NOCHA to conduct a comprehensive analysis of 5 openly available and 6 closed-source models, identifying where models struggle, thereby providing valuable insights

improving long-context models.<sup>1</sup>

3. **METHODOLOGY:** We contribute our data collection and evaluation methodology, which balances efficiency and quality, minimizing the possibility of models receiving credit for correct predictions made without full utilization of the context.

## 2 Data & Methods

In this section, we first motivate NOCHA’s design—narrative minimal pairs written by readers of books—and then describe our data collection process and evaluation methodology.

### Issues with existing long-context benchmarks:

A popular method to evaluate long-context models is the “needle-in-a-haystack” task (Kamradt, 2023, NIAH), which involves injecting a sentence-level piece of information at different depths of a document. NIAH was recently extended by RULER (Hsieh et al., 2024) to include various types of needles and other synthetic tasks. While this approach allows control over the exact position of the evidence, it has several limitations: (1) evidence scope is heavily restricted by the length and complexity of the needle, (2) the needle is unrelated to the continuous context, mak-

<sup>1</sup>In this paper, we define long-context models as those with a claimed context window of at least 128k tokens.

ing it easier to retrieve,<sup>2</sup> and (3) the task is synthetic, making it a poor proxy for real-world tasks. Recently, more realistic long-context benchmarks have emerged. The most similar to NOCHA are L-EVAL (An et al., 2024), NOVELQA (Wang et al., 2024a), and  $\infty$ bench (Zhang et al., 2024), which are all literary QA tasks. However, all three use novels that are publicly available online, meaning they are likely part of LLMs’ pretraining data along with the numerous books, articles, and summaries written about them.<sup>3</sup> They also lack the complex global coverage of many of our claims; for instance, while 35% of NOVELQA reportedly consists of multi-hop questions, these involve simpler tasks such as aggregation (e.g., “How many times has Bran jumped off and ran?”) that do not require reasoning over *implicit* information in the text.

## 2.1 Data collection

To address these limitations, we adopt a design that includes (1) human-written examples for claim verification, an important real-world task; (2) recently-published texts to mitigate data contamination issues (Jacovi et al., 2023); (3) fictional texts to prevent over-reliance on parametric knowledge vs. in-context information; and (4) minimal pairs to guard against models being “right for the wrong reasons” (see Figure 2 for example) while also enabling easy verification to ensure dataset quality.

### Collecting a corpus of recently-published fiction:

We collect the 67 books in NOCHA following the approach of Kim et al. (2024), who created the FABLES dataset to evaluate summarization faithfulness with similar objectives in mind (reducing data contamination and reliance solely on parametric knowledge). Concretely, we select books that are: (1) fictional narratives, (2) published in 2023 or 2024,<sup>4</sup> and (3) self-reported as having been read by our annotators. Our resulting dataset comprises 63 new books (33 published in 2023 and 30 in

<sup>2</sup>The context in Hsieh et al. (2024) may match the needle’s topic, but it is not a continuous, logically connected text.

<sup>3</sup>While  $\infty$ bench swaps core entities in the texts, we found that prompting GPT-4-TURBO with excerpts from the “falsified” novels results in the model recognizing the original novel and identifying the changed entities.

<sup>4</sup>We will not release the full NOCHA dataset because (1) the books are copyrighted, and (2) we want to prevent model providers from training on this data and compromising the benchmark. Instead, we additionally annotate a smaller subset of *classic books* that are out of copyright. This subset will serve as a publicly available sample of the data. The authors commit to updating NOCHA with new books and evaluating new long-context models.

	Books 📖 (n=67)		Claim <sub>{Pairs}</sub> 📄 (n=2002) <sub>{1001}</sub>		
	TOKENS	WORDS	TOKENS	WORDS	# CLAIM/BOOK
MEAN	127,324	98,587	23.22	18.26	14.94
ST. DEV.	52,561	39,506	7.62	6.49	8.37
MAX	336,288	257,445	63	57	46
MIN	49,156	38,023	5	4	4

Table 1: Summary statistics for NOCHA.

2024) and four classic novels (see Table 11 for full list of the books). The mean length of books in our dataset is **127k** tokens and **98.5k** words (see Table 1 for statistics).<sup>5</sup>

**Annotators:** We recruited 18 annotators via [Upwork](#) and 5 volunteer annotators (convenience sample) who reported having read books published in 2023 or 2024.<sup>6</sup> All annotators were required to read the guidelines and sign a consent form before starting the task. They were compensated \$3–\$5 USD per pair of claims, based on the agreed rate and the number of annotations, resulting in an hourly rate of approximately \$18–\$30 USD.<sup>7</sup> Overall, we collected 1001 claim pairs (approximately 10-15 per book) at a cost of \$2.8k USD. The protocol was reviewed and deemed *Not Human Subjects Research* by the Institutional Review Board. See §C.2 for details on the annotators and the recruitment process.

**Collecting true/false pairs:** To collect true/false pairs, annotators were instructed to first write a true statement about the events or characters in the book, and then create a corresponding false statement addressing the same aspect of the book such that the verification of one claim as “true” necessitates the verification of the other as “false”. We trained annotators by having them write a true counterpart to a false model-generated claim.<sup>8</sup> Both types of pairs—those written from scratch and those where only one claim was created by the annotator—were included in the final dataset after quality control.<sup>9</sup>

<sup>5</sup>All token counts provided in this paper are based on tiktoken (<https://github.com/openai/tiktoken>) with the cl100k encoding and word counts are determined by splitting the text on whitespaces, unless stated otherwise.

<sup>6</sup>The volunteer group included three of the authors, each of whom read three books of their choice and performed the annotations to better understand the difficulty of the task and the time required to complete it.

<sup>7</sup>The annotators were usually able to create about 6-10 claim pairs per hour.

<sup>8</sup>Model-generated claims were selected from FABLES (Kim et al., 2024) as well as generated by CLAUDE-3-OPUS and GPT-4-TURBO.

<sup>9</sup>Around 30% of the pairs in the data contain one model generated claim. In some cases, these claims were edited in

MODEL	CONTEXT	AVAIL.	CHECKPOINTS	# PARAM
GPT-4o	128k	🔒	gpt-4o-2024-05-13	~19.7
GPT-4-TURBO	128k	🔒	gpt-4-turbo-2024-04-09	~19.7
CLAUDE-3-OPUS	200k	🔒	claude-3-opus-20240229	~19.7
CLAUDE-3.5-SONNET	200k	🔒	claude-3-5-sonnet@20240620	~19.7
GEMINI PRO 1.5	1M	🔒	gemini-1.5-pro-preview-0514	~19.7
GEMINI FLASH 1.5	1M	🔒	gemini-1.5-flash-preview-0514	~19.7
COMMAND R	128k	🔒	c4ai-command-r-v01	35B
COMMAND R+	128k	🔒	c4ai-command-r-plus	104B
GEMMA-10M	10M	🔒	gemma-2b-10m	2B
PHI-3-MINI	128k	🔒	Phi-3-mini-128k-instruct	3.8B
LONGLLAMA	256k	🔒	long_llama_3b_instruct	3B

Table 2: Evaluated models: the upper row displays all closed-source models, while the lower row lists all open-weight models (see §D for details).

The annotators were also asked to write a short explanation as to why these claims were true or false based on the content of the book. We closely monitored this process, providing feedback and seeking clarifications from the annotators.

**Quality control:** To ensure the quality of the claims, we hired annotators who had read the same books to validate 152 claims from six books, one new annotator per book. The annotators were paid \$1.66–\$2.77 USD per claim based on their requested rate, with the total annotation costing \$285 USD. Overall, the annotators agreed on 148 out of 152 labels (Krippendorff’s  $\alpha = 0.960$ ) (see §C.2 for details on this process). This high agreement also allow us to conclude that **human readers are very strong performers on NOCHA’s claim verification task (97.4%)**. Finally, after collecting all pairs, two of the authors reviewed all instances, resolving unclear cases with the annotators and between themselves if a claim appeared to be subjective or incorrect.

**Evidence scope:** To assess the reasoning required to validate a claim pair, we obtained additional annotations on 121 claim pairs from 8 books. These annotations determine whether verification is possible based on (1) one or two sentences (similar to NIAH), (2) a single contiguous passage, or (3) global reasoning over the full book. Annotators were compensated \$2 USD per pair, totaling \$242 USD for the annotations. Overall, our annotations suggest that 12.4% of pairs can be validated with reasoning over one or two sentences, 39.7% require reasoning over a longer contiguous passage, and **the largest fraction, 47.9%, necessitate global reasoning for correct verification.**

### 3 Experiments

We test the long context reasoning capabilities of 5 open-weight models and 6 closed-source models (see Table 2) by prompting each model in a zero-shot manner to verify a single claim given the entire book content as evidence, similar to the summarization setup of Kim et al. (2024). After pilot experiments with different prompts (see §D for a full description), we observe that the best-performing prompt asks the model to first explain its reasoning before making a final decision, similar to chain-of-thought prompting (Wei et al., 2022).

**Model and prompt details:** Each model was prompted with a claim and the entire book using the prompt template in Table 17.<sup>10</sup> All **closed-source models** were accessed via the provider’s API<sup>11</sup> at a total cost of ~\$8k USD. Since we found that **open-weight models** struggle to follow the prompt’s instructions, we also experiment a simplified prompt that just asked for a true or false decision (see Table 18).<sup>12,13</sup> Finally, to measure whether retrieval-augmented language models perform better or worse than the long-context setting, we also experiment with a **BM25+GPT-4o** configuration, in which BM25 (Robertson et al., 1995) is used to rank the most relevant  $k \in \{5, 25, 50\}$  excerpts from the book.<sup>14</sup> We then prompt GPT-4o (our best performing model) with the retrieved excerpts as evidence using the prompt in Table 19.<sup>15</sup>

**Evaluation:** We report the overall PAIRWISE ACCURACY for each model. Models get credit only if they label *both* the true and false claim in a pair correctly; and no points otherwise. We report the number of correctly-verified pairs divided by the

order to decontextualize the claim.

<sup>10</sup>We set the temperature to 0 and generate up to 800 tokens.

<sup>11</sup>We encounter several issues while generating the text. Most notably, GEMINI PRO 1.5 and GEMINI FLASH 1.5 refuse to process some of the books returning a prohibited content error (likely due to copyrights). Overall, both models refused to generate the label in about 48% of cases significantly reducing the number of pairs tested for these models.

<sup>12</sup>Outputs generated using the simplified prompt are denoted by the subscript *simple*.

<sup>13</sup>We encountered some generation issues with LONGLLAMA which failed to produce outputs when prompted with the prompt in Table 17. Hence, we only report the results with the simplified prompt for this model.

<sup>14</sup>Excerpts were an average of 285 words based on whitespace splitting with all paragraph breaks preserved.

<sup>15</sup>50 excerpts equals at most 39% of the total text of any book in our dataset; additional statistics on the total percentage of a book retrieved using this approach are available in Figure 13. The total cost of this experiment was \$330 USD.

MODEL	PAIR ACC <sub>(correct/total)</sub>	COMMON SET ACC
GPT-4o	<b>55.8</b> (344/617)	<b>58.2</b> (206/354)
GPT-4-TURBO	40.2 (248/617)	40.1 (142/354)
CLAUDE-3-OPUS	49.4 (463/937)	50.8 (180/354)
CLAUDE-3.5-SONNET	41.0 (384/937)	40.7 (144/354)
GEMINI PRO 1.5	48.1 (247/514)	48.3 (171/354)
GEMINI FLASH 1.5	34.2 (176/515)	35.0 (124/354)
COMMAND R	19.6 (87/445)	<i>n/a</i>
COMMAND R <sub>simple</sub>	22.5 (100/445)	<i>n/a</i>
COMMAND R+	17.3 (77/445)	<i>n/a</i>
COMMAND R+ <sub>simple</sub>	13.7 (61/445)	<i>n/a</i>
PHI-3-MINI	9.3 (23/247)	<i>n/a</i>
PHI-3-MINI <sub>simple</sub>	14.5 (48/331)	<i>n/a</i>
GEMMA-10M	3.9 (39/1001)	<i>n/a</i>
GEMMA-10M <sub>simple</sub>	7.5 (75/1001)	<i>n/a</i>
LONGLLAMA <sub>simple</sub>	4.9 (61/937)	<i>n/a</i>
BM25+GPT-4o ( <i>k</i> =5)	28.2 (282/1001)	28.9 (102/353)
BM25+GPT-4o ( <i>k</i> =25)	44.1 (441/1001)	46.5 (164/353)
BM25+GPT-4o ( <i>k</i> =50)	49.7 (497/1001)	51.0 (180/353)
RANDOM	25.0 (250/1001)	25.0 (88/353)

Table 3: Percentage of claim pairs identified correctly by each model (see Table 17 for the prompt; and Table 19 for the prompt employed with BM25). “COMMON SET” refers to claim pairs shared among the models. The subscript “SIMPLE” refers to the calls done with simplified prompt (Table 18). Models performing below random are marked in red. This results include also four classic novels, which were most likely in the training data. For results excluding these novels see Table 24.

total number of pairs that the model labeled. See §D for details.

## 4 Results & analysis

Table 3 reports results for individual models as well as for a common set of claims shared among closed-source models.<sup>17</sup> Overall, GPT-4o exhibits the highest accuracy at 55.8% (58.2% for the common set), though it still misidentifies at least one claim in the pair 44.2% of the time (see Table 6 for example failure cases). The second-best performing model is CLAUDE-3-OPUS with an overall accuracy of 49.4%, followed by GEMINI PRO 1.5 with an accuracy of 48.1%. CLAUDE-3.5-SONNET, GPT-4-TURBO and GEMINI FLASH 1.5 perform worse, with accuracies of 41.0%, 40.2% and 34.2%, respectively.<sup>18</sup> All open-weight models perform below random (25%), ranging from 22.5% for COMMAND R (with a simplified prompt) to as low as 3.9% for GEMMA-10M. Because of their poor accuracies, we exclude open-weights models from the following analysis; see §E for details.

<sup>17</sup>Open-weights models are excluded as this would significantly reduce the number of shared pairs to just over 100.

<sup>18</sup>See Figure 9 for example of claim on which all models failed to generate the accurate label.

MODEL	RULER (%) VANILLA NIAH	RULER (%) NIAH SUITE	NoCHA (%)
GPT-4-TURBO	100.0	89.6	40.2
COMMAND R	98.0	84.8	19.6 / 22.5 <sub>simple</sub>

Table 4: Performance on NIAH does not translate to NoCHA accuracy for GPT-4-TURBO and COMMAND R. The table includes published results with 128k tokens from Hsieh et al. (2024) for both the standard NIAH (Table 11) and the NIAH suite, which averages results from 8 NIAH variants (Table 13). Additionally, we provide a comparison of performance on the entire RULER benchmark versus NoCHA for overlapping models in Table 5 for reference.

MODEL	RULER (%) (128k)	NoCHA (%) (~128k)
GEMINI PRO 1.5	94.4	48.1 / 48.3 (~128k)
GPT-4-TURBO	81.2	40.2
COMMAND R	76.0	19.6 / 22.5 <sub>simple</sub>
COMMAND R+	63.1	17.3 / 13.7 <sub>simple</sub>
PHI-3-MINI	43.3	9.3 / 14.5 <sub>simple</sub>

Table 5: Performance on RULER (Hsieh et al., 2024) compared to NoCHA accuracy for overlapping models. The table includes results for the 128k tokens group from <https://github.com/hsiehjackson/RULER>.

**Good performance on “needle in the haystack” does not imply high NoCHA accuracy:** In Table 4, we compare the performance of COMMAND R and GPT-4-TURBO on NoCHA with the results reported in Hsieh et al. (2024, RULER), a more complex variant of NIAH, to estimate how performance on RULER translates to our task. Despite both models achieving high performance on RULER (84.8% and 89.6%, respectively, for the longest tested context of 128k), they underperform on our task, with GPT-4-TURBO achieving only 40.2% accuracy and COMMAND R performing *below* random at 19.6%.<sup>19</sup> This shows that synthetic retrieval-focused benchmarks are insufficient to evaluate global long-document reasoning.

**NoCHA claims that require global reasoning are particularly difficult to verify:** Table 7 contains

<sup>19</sup>Note that RULER reports the results for gpt-4-1106-preview checkpoints for GPT-4-TURBO, while the current work uses gpt-4-turbo-2024-04-09 checkpoints, presumably a more advanced version of that model.

<sup>19</sup>In fact, after the first victim dies, the narrator mentions: “All I can tell you is what I’ve been telling you so far: the truth. After all, I told you Henry McTavish would be poisoned, didn’t I? Well, not in those words, I suppose. But I did say the inspiration for this book would come from a drink with him.” It is also worth noting that the model incorrectly states that the first victim dies during the panel, when in fact McTavish dies the next day during breakfast.



Everyone on this train is a suspect by Stevenson Benjamin

**TRUE** Ernest, the narrator, drops numerous hints about the killer's identity from the start, he also hints on how the first victim dies as early as chapter two.

**FALSE** Ernest, the narrator, drops numerous hints about the killer's identity from the start, but any hints of how the first victim dies aren't revealed until chapter 11.

**Claude-3.5-Sonnet**

After carefully reviewing the context, I can confirm that Ernest does indeed drop hints about the killer's identity and the first victim's cause of death early in the story. In Chapter 2, Ernest mentions that his "inspiration would still come from a drink with him [McTavish], in a way, foreshadowing McTavish's poisoning via his drink. Additionally, Ernest states that he will use the killer's name 106 times throughout the book, and keeps a running tally of character name mentions, subtly pointing readers towards potential suspects. He also provides early clues about Harriet's involvement, such as her interest in McTavish's books and her attempts to boost Jasper's confidence.

After carefully reviewing the context, I found that Ernest does provide some hints about the killer's identity early on, such as mentioning the name tally and discussing various characters' motives. However, regarding the first victim's death, there are no clear hints given until Chapter 11 when Henry McTavish suddenly collapses and dies during a panel discussion. Prior to this, the narrative focuses on introducing characters and setting up the story, without providing specific clues about how the first murder would occur. Therefore, the statement appears to be accurate.

**TRUE** ✓ **FALSE** ✗



The Tainted Cup by Robert Jackson

**TRUE** Despite her skills as an Apoth, Nuisis is unable to reverse engineer the type of portal opened by the reagents key found in Rona's wooden chest.

**FALSE** By using her skills as an Apoth, Nuisis is able to reverse engineer the type of portal opened by the reagents key found in Rona's wooden chest.

**Gemini Pro 1.5**

The statement is false because the reagents key found in Rona's wooden chest was not given to Nuisis. The reagents key found in Rona's wooden chest was the one Dinkos kept and used to open the Maya estate. The key given to Nuisis was the one found in the empty house near Rona's residence.

The text states that Nuisis is unable to determine the type of portal opened by the reagents key found in Rona's wooden chest. She runs it under tests, but they are inconclusive.

**TRUE** ✗ **FALSE** ✓

Figure 2: Examples of claim pairs where the models failed to validate one of the claims in the pair. Employing narrative minimal pairs helps us avoid awarding the model for cases where it only *appears* to produce the correct answer, when in fact the prediction was made without fully or efficiently utilizing the context. In the top example, the model first correctly identifies hints dropped by the author about how the first victim dies (when verifying the true claim),<sup>16</sup> then incorrectly claims that no such hints exist (when verifying the false claim). In the bottom example, the model first incorrectly claims that the key given to Nuna was not found in Rona's wooden chest (when verifying the true claim), then correctly raises no objections to the fact that the key was found in Rona's wooden chest (when verifying the false claim).

CLAIM & GOLD LABEL	MODEL'S LABEL & JUSTIFICATION	HUMAN EXPLANATION
Wyatt apologizes to Ernest because he believes McTavish left him a negative review.  Everyone on the Train is a Suspect by Benjamin Stevenson  -True	Wyatt's apology to Ernest is mentioned in the context, but it is not because he believes McTavish left him a negative review. <b>Instead, Wyatt's apology is vague and not directly linked to the review.</b> Ernest himself is confused about the reason for the apology, and it is later revealed that Simone, Ernest's agent, was behind the negative review as part of a coded message to McTavish. Therefore, the statement that Wyatt apologizes to Ernest because he believes McTavish left him a negative review is incorrect.  -GPT-4o: <b>False</b>	When Wyatt apologizes, Ernest does not initially understand the reason behind it. However, it is heavily implied shortly after that Wyatt, McTavish's publisher, was apologizing for the bad review, which at that time everyone, including Wyatt thought was written by McTavish himself.
While on Sigrid's boat, Seventeen admits to Kat that he had the same underwater dream that she had.  Assassin Eighteen by John Brownlow  -True	The statement is true based on the context provided. In Part 8, after Kat is rescued from the water by Rachel and Vilmos in the Zodiac, she describes her near-death experience underwater where she felt she and Seventeen were entwined as one. Seventeen then admits "I had the same dream", confirming he experienced the same vision while he was drowning.  -CLAUDE-3-OPUS: <b>True</b>	It is true that Seventeen admits he had the same dream in Part 8 of the book. However, Kat was not rescued by Rachel and Vilmos in the Zodiac during this scene. That happens much later in the book, towards the end of Part 11. This conversation happened when Kat and Seventeen were rescued by Sigrid.

Table 6: **Top row:** Example of a claim where the model generated incorrect label, along with the model-generated label and justification, and human-written explanations for the correct label. **Bottom row:** Example of claim correctly labeled by the model but with wrong justifications. The explanations were provided by the authors, based on examples from books they read. Tags have been removed from the generated content for clarity. For more examples see Table 23.

further analysis of model accuracy on a subset of NoCHA annotated for the scope of evidence (see §2.1 for details). Overall, models perform worst for claim pairs requiring global reasoning (41.6%), followed by reasoning over a longer pas-

sage (47.6%), and, finally, sentence-level evidence (59.8%). While performance on sentence-level evidence is higher than in the other two setups, it is still much lower than the "needle-in-a-haystack" performance reported in Hsieh et al. (2024). This

SCOPE	CORRECT	INCORRECT	ACCURACY (%)
GLOBAL	104	146	41.6
PASSAGE	101	111	47.6
SENTENCE	49	33	59.8

Table 7: Overall accuracy on claim pairs of closed-source models on a subset of data annotated for evidence scope (see §2 for the annotation process). Pairs requiring global reasoning to verify are most difficult for models. Figure 11 shows results by model.

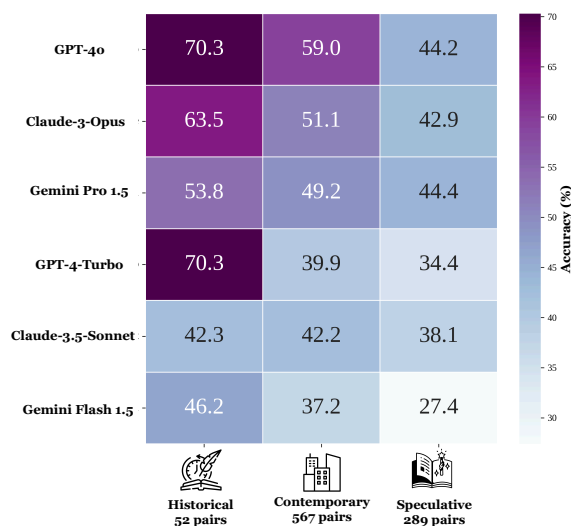


Figure 3: Performance of closed-source models on different types of novels. Note that two novels were excluded from this analysis as they could not clearly be classified in one of these categories. We provide the total number of claims in each category for reference however these numbers will vary slightly between the models due to the context-limitation and GEMINI PRO 1.5/GEMINI FLASH 1.5 refusals.

indicates that NOCHA claim pairs with sentence-level evidence are still much harder to solve than NIAH, possibly due to NIAH’s out-of-context evidence injection.

**Claims about speculative fiction are harder to verify:** Are texts with extensive world-building more difficult for models to process than those that take place in a fictional but realistic world? We categorized 65 out of 67<sup>20</sup> of the books in our dataset into one of three broad categories:

- **Historical:** works set in our world before World War II, without any unrealistic elements.

<sup>20</sup>One book was a nonfiction collection of essays and another split the plot between modern day and the past, rendering it a combination of historical and contemporary fiction.

- **Contemporary:** works set in our world after World War II, without any unrealistic elements.
- **Speculative:** works set in an alternate version of our world, containing both realistic and unrealistic elements, or in a completely invented universe.

The accuracy across all six closed-source models is 56.4% for historical fiction, 46.8% for contemporary fiction, and 38.8% for speculative fiction. Figure 3 shows that this pattern holds for each individual model, with accuracy being highest for historical fiction, followed by contemporary fiction, and lowest for speculative fiction. These results support the intuition that texts set in a realistic version of our world require less “work” from models to reason over than texts set in a universe that is defined within the text, perhaps because the models can rely more on their parametric knowledge.<sup>21</sup>

**Irrelevant context confuses models:** If we know that a claim can be verified by a short span of the book in isolation, does providing the rest of the book affect its accuracy? Figure 4 shows model performance on four *short story collections*, where the model is either given just the short story relevant to the claim (~21k tokens)<sup>23</sup> or the entire collection (~129k tokens). By prompting the model with the same claims, we partially control for confounding variables such as the inherent difficulty of the claim. While the Gemini models are relatively robust to added context, CLAUDE-3-OPUS’s pair accuracy drops as much as 44.5% absolute when given the collection vs. the story, while CLAUDE-3.5-SONNET, GPT-4-TURBO and GPT-4O exhibit smaller but still substantial drops.<sup>24</sup> More generally, it is unclear whether claims about longer books (with potentially more irrelevant context) are harder to verify; see §E for more details.

<sup>21</sup>We acknowledge that there may be confounding factors, such as the relative difficulty of writing challenging claims for different books types and the high average length of speculative fiction books (148k tokens)—though we note that all models perform better on historical (avg. length 133k tokens) than the shorter contemporary books (avg. length 115k tokens).

<sup>23</sup>Mean=8.5k, Min=702, Std.Dev.=4,355.

<sup>24</sup>Notably, all models perform suboptimally on the short stories, suggesting that reasoning tasks may be challenging, even with shorter contexts.

<sup>24</sup>Here we employ the 32k, which can also be extended to 128k using YARN (Peng et al., 2023).

<sup>25</sup>We do not analyze CLAUDE-3.5-SONNET’s explanations because the model was run on June 22, after the completion



Figure 4: Model accuracy on claim pairs for different stories within collections. Accuracy is shown for (1) using the entire collection as context when prompting about a story, and (2) using only the individual story ("story") for the same set of claims. For GPT-4o and GPT-4-TURBO, one book was too long, so only the "story" context performance is presented. GEMINI PRO 1.5 and GEMINI FLASH 1.5 refused to process two books but handled the stories within, so only "story" context performance is available. We also provide performance of MIXTRAL-8X22B (65k) and QWEN-2-72B (32k)<sup>22</sup> on story-level input, for comparison.

### Model-generated explanations are not always accurate:

Do models correctly explain why a claim is true or false? Three of this paper’s authors examined explanations generated for NOCHA books that they read, covering a subset of 7 books (293 claims). Table 8 shows that *no* model consistently produces accurate explanations for every correctly-labeled claim, indicating a potential reliance on flawed or incomplete reasoning (see Table 6 for examples). CLAUDE-3-OPUS demonstrated the highest explanation accuracy, with 16.9% of its explanations containing errors, followed by GPT-4o and GEMINI PRO 1.5 with 21.7% and 28.3% incorrect explanations respectively. GPT-4-TURBO and GEMINI FLASH 1.5 were the worst performing models with incorrect explanations reaching 45.0% and 65.9% respectively.<sup>26</sup> This is problematic when considering that users tend to rely on model explanations when verifying claims, even when the explanation is incorrect (Si et al., 2024). Further discussion

of this analysis. Nonetheless, we have observed mistakes in its explanation, with an example provided in Table 6.

<sup>26</sup>As we run CLAUDE-3.5-SONNET on June 22, we exclude it from this analysis.

MODEL	INCORRECT % (incorrect/total)
GPT-4o	21.7 (15/69)
GPT-4-TURBO	45.0 (27/60)
CLAUDE-3-OPUS	16.9 (13/77)
GEMINI PRO 1.5	28.3 (13/46)
GEMINI FLASH 1.5	65.9 (27/41)

Table 8: The percentage of incorrect justifications provided for correct label by model which generated it in the analyzed subset of data.<sup>25</sup>

SCOPE	BM25+GPT-4o		
	k=5 (%)	k=25 (%)	k=50 (%)
GLOBAL	25.9	29.3	41.4
PASSAGE	22.9	43.8	45.8
SENTENCE	46.7	66.7	73.3

Table 9: Performance of the BM25+GPT-4o pipeline on the subset of data annotated for evidence scope.

on citations of evidence from the source found in model justifications can be found in §E.

### Can BM25 help prioritize important context?

We observe that our BM25-assisted GPT-4o approach with  $k = 50$  excerpts performs better than all models except for GPT-4o with the full book. For  $k = 50$ , an average of only 17% of a book was retrieved by BM25 and fed to GPT-4o in ranked order. Retrieval-based methods struggle on global reasoning due to out-of-order chunks, but they are effective for claims that require sentence- and, to a much lesser extent, passage-level retrieval to verify; thus, there is likely an upper bound to their NOCHA accuracy (see Table 9).

### Models have different predilections for predicting True vs False:

Our pairs were designed so that validating one claim should enable validation of the other. However, we observe in Table 10 that some models tend to predict one label much more frequently than another. This tendency was particularly evident in CLAUDE-3.5-SONNET, GEMINI PRO 1.5, GEMINI FLASH 1.5, and GPT-4-TURBO, which had strong preferences for predicting False, and is in line with the observation reported for GEMINI PRO 1.5 in Levy et al. (2024). In contrast, CLAUDE-3-OPUS exhibited much higher accuracy on True labels (82.2%) compared to False (64.7%). GPT-4o was the only balanced model among the closed-source models, with accuracies of 77.5% for True and 75.9% for False.



MODEL	TRUE <sub>(correct/total)</sub>	FALSE <sub>(correct/total)</sub>
GPT-4o	77.5 <sub>(478/617)</sub>	75.9 <sub>(468/617)</sub>
GPT-4-TURBO	57.2 <sub>(353/617)</sub>	78.8 <sub>(486/617)</sub>
CLAUDE-3-OPUS	82.2 <sub>(770/937)</sub>	64.7 <sub>(606/937)</sub>
CLAUDE-3.5-SONNET	55.3 <sub>(518/937)</sub>	83.5 <sub>(782/937)</sub>
GEMINI PRO 1.5	59.7 <sub>(307/514)</sub>	83.7 <sub>(431/515)</sub>
GEMINI FLASH 1.5	51.7 <sub>(266/515)</sub>	78.1 <sub>(402/515)</sub>
COMMAND R	59.3 <sub>(264/445)</sub>	52.8 <sub>(235/445)</sub>
COMMAND R <sub>simple</sub>	74.4 <sub>(331/445)</sub>	44.3 <sub>(197/445)</sub>
COMMAND R+	58.7 <sub>(261/445)</sub>	49.2 <sub>(219/445)</sub>
COMMAND R+ <sub>simple</sub>	34.6 <sub>(154/445)</sub>	75.3 <sub>(335/445)</sub>
PHI-3-MINI	34.8 <sub>(88/253)</sub>	46.5 <sub>(118/254)</sub>
PHI-3-MINI <sub>simple</sub>	69.5 <sub>(230/331)</sub>	32.0 <sub>(106/331)</sub>
GEMMA-10M	20.2 <sub>(202/1001)</sub>	27.8 <sub>(278/1001)</sub>
GEMMA-10M <sub>simple</sub>	77.5 <sub>(776/1001)</sub>	18.7 <sub>(187/1001)</sub>
LONGLLAMA <sub>simple</sub>	63.8 <sub>(598/937)</sub>	19.7 <sub>(185/937)</sub>
BM25+GPT-4o ( $k=5$ )	33.7 <sub>(337/1001)</sub>	90.0 <sub>(901/1001)</sub>
BM25+GPT-4o ( $k=25$ )	57.5 <sub>(576/1001)</sub>	81.7 <sub>(818/1001)</sub>
BM25+GPT-4o ( $k=50$ )	65.7 <sub>(658/1001)</sub>	79.6 <sub>(797/1001)</sub>
RANDOM	50.0 <sub>(500/1001)</sub>	50.0 <sub>(50/1001)</sub>

Table 10: Model accuracy on True and False claims for all data (see Table 17 for the prompt; and Table 19 for the prompt employed with BM25). “COMMON SET” refers to claim pairs shared among the models. The subscript “SIMPLE” refers to the calls made with the simplified prompt (Table 18).

## 5 Related Work

**Evaluation of long-context models:** Several benchmarks with inputs ranging from at least 8k up to 1M tokens have recently been introduced to evaluate long-context language models. Some of these benchmarks use synthetic tasks that can be generated programmatically or by LLMs, such as NIAH (Tay et al., 2021; Sun et al., 2021; Kamradt, 2023; Li et al., 2023; Hsieh et al., 2024; Levy et al., 2024; Liu et al., 2024; Lee et al., 2024; Laban et al., 2024; Vodrahali et al., 2024). Others contain “realistic” tasks that typically require human annotators to devise, like traditional QA, summarization, and claim verification (Shaham et al., 2023; Bai et al., 2023; Dong et al., 2023; Wang et al., 2024b) sometimes created by utilizing existing datasets (Hudson and Al Moubayed, 2022; Li et al., 2024). The existing long-context evaluation benchmarks with literary tasks all consist of publicly available English (Hudson and Al Moubayed, 2022; Zhang et al., 2024; Wang et al., 2024a; An et al., 2024) or Chinese language books (Yu et al., 2024). All but An et al. (2024) contain generative or multiple choice questions, including true/false questions, but unlike NOCHA, these are not minimal pairs.

**Claim verification:** Our work relates to prior work on claim verification, where claims are veri-

fied against a whole knowledge datastore (Thorne et al., 2018; Wadden et al., 2020; Schuster et al., 2021) or single evidence documents (Maynez et al., 2020; Falke et al., 2019). While earlier methods primarily relied on task-specific natural language inference (Laban et al., 2022; Honovich et al., 2022) or question answering models (Fabbri et al., 2022b; Wang et al., 2020) for claim verification, more recent work has explored using LLMs to evaluate the factuality of long-form model-generated text (Min et al., 2023; Wei et al., 2024; Manakul et al., 2023).

**Minimal pairs:** Minimal pairs, or contrast sets, are pairs of test set instances where a slight but impactful difference between the instances affects the gold label. Gardner et al. (2020) proposed contrast sets as a method for evaluating models on their intended tasks by manually perturbing test set instances. Our narrative minimal pairs were inspired by the DEMETR (Karpinska et al., 2022) dataset of minimal pairs for evaluating machine translation metrics and the BLiMP (Warstadt et al., 2020) benchmark of minimal pairs for evaluating English grammatical phenomena.

## 6 Conclusion and Discussion

We introduce NOCHA, a claim verification dataset designed to evaluate long context LLMs in a realistic task setting. Our design ensures that most claims necessitate global reasoning over extended contexts for verification and cannot be easily “gamed” by relying solely on parametric knowledge or generating correct answers through incorrect reasoning. Our experiments with 11 different long-context LLMs (5 open-weight, 6 closed-source) demonstrate that all models struggle on NOCHA. Furthermore, we reveal a substantial competency gap with human readers who can very easily perform this task.

Importantly, our results show that models that are “state-of-the-art” according to synthetic benchmarks like NIAH actually perform very poorly on our meticulously designed dataset. Nevertheless, we argue that complex synthetic datasets (such as RULER) are useful and complementary to our realistic dataset; they allow for much higher flexibility such as easily adjusting for different context lengths or analyzing the lost-in-the-middle phenomenon. We encourage researchers to use a holistic approach and consider *both* synthetic and realistic tasks when evaluating long-context language models.

## Limitations

The scope of our work is limited to novels published in English and the task of claim verification. It is unclear how the models' performance would translate to other languages, domains, or realistic tasks, and we leave that for future work.

We also acknowledge that this study and the methodology proposed are inherently and possibly prohibitively expensive due to the hiring of expert annotators and the thousands of LLM API calls. Because using the model developer's API is often the only way to access a closed-source model, the extent of our evaluation was limited to whether or not the closed-source models provided a response to our prompts. This resulted in a common set of only 354 pairs that all closed-source models labeled, making it difficult to truly compare the models to each other. Finally, while we commit to periodically updating the dataset and evaluating the models ourselves, we are unable to release the entire dataset due to copyright restrictions and to prevent model providers from training on it.

## Ethical Considerations

The data collection protocol was reviewed by the University of Massachusetts, Amherst Institutional Review Board and received a *Not Human Subjects Research* determination (IRB: #5587). All annotators consented to the use and publication of their annotations, which we will release for the portions of data that is *not* under copyright. Additionally, we ensured annotators received fair compensation for their contributions and respected their preferred rates. All copyrighted books were purchased using the funds that supported this work. We will *not* release the copyrighted portion of the data. Instead, we commit to consistently updating the dataset with newly published novels and evaluating the models ourselves.

## Acknowledgements

We would like to extend our gratitude to the Upwork annotators for their dedicated efforts, and to the anonymous reviewers for their time and valuable feedback. We are also grateful to the members of the UMass NLP lab for their insightful input, with special thanks to Chau Pham, Nader Akoury, Yekyung Kim, and Yapei Chang. Our deep appreciation goes to Byron Wallace for suggesting the project name and to Simeng Sun for the discussions that inspired the initial concept of this work.

This project was partially supported by awards IIS-2202506 and IIS-2312949 from the National Science Foundation (NSF) as well as an award from Adobe.

## References

- Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Qin Cai, Martin Cai, Caio César Teodoro Mendes, Weizhu Chen, Vishrav Chaudhary, Dong Chen, Dongdong Chen, Yen-Chun Chen, Yi-Ling Chen, Parul Chopra, Xiyang Dai, Allie Del Giorno, Gustavo de Rosa, Matthew Dixon, Ronen Eldan, Victor Fragoso, Dan Iter, Mei Gao, Min Gao, Jianfeng Gao, Amit Garg, Abhishek Goswami, Suriya Gunasekar, Emman Haider, Junheng Hao, Russell J. Hewett, Jamie Huynh, Mojan Javaheripi, Xin Jin, Piero Kauffmann, Nikos Karampatziakis, Dongwoo Kim, Mahoud Khademi, Lev Kurilenko, James R. Lee, Yin Tat Lee, Yuanzhi Li, Yunsheng Li, Chen Liang, Lars Liden, Ce Liu, Mengchen Liu, Weishung Liu, Eric Lin, Zeqi Lin, Chong Luo, Piyush Madan, Matt Mazzola, Arindam Mitra, Hardik Modi, Anh Nguyen, Brandon Norick, Barun Patra, Daniel Perez-Becker, Thomas Portet, Reid Pryzant, Heyang Qin, Marko Radmilac, Corby Rosset, Sambudha Roy, Olatunji Ruwase, Olli Saarikivi, Amin Saied, Adil Salim, Michael Santacrose, Shital Shah, Ning Shang, Hiteshi Sharma, Swadheen Shukla, Xia Song, Masahiro Tanaka, Andrea Tupini, Xin Wang, Lijuan Wang, Chunyu Wang, Yu Wang, Rachel Ward, Guanhua Wang, Philipp Witte, Haiping Wu, Michael Wyatt, Bin Xiao, Can Xu, Jiahang Xu, Weijian Xu, Sonali Yadav, Fan Yang, Jianwei Yang, Ziyi Yang, Yifan Yang, Donghan Yu, Lu Yuan, Chengruidong Zhang, Cyril Zhang, Jianwen Zhang, Li Lyna Zhang, Yi Zhang, Yue Zhang, Yunan Zhang, and Xiren Zhou. 2024. [Phi-3 technical report: A highly capable language model locally on your phone](#). *Preprint*, arXiv:2404.14219.
- Mustafa Aljadery, Siddharth Sharma, and Aksh Garg. 2024. [Gemma 2b - 10m context](#). Technical report. Accessed: 2024-05-23.
- Chenxin An, Shansan Gong, Ming Zhong, Mukai Li, Jun Zhang, Lingpeng Kong, and Xipeng Qiu. 2024. L-eval: Instituting standardized evaluation for long context language models. In *ICLR*.
- Anthropic. 2024a. [Claude 3.5 Sonnet](#). Technical report, Anthropic. Accessed: 2024-06-23.
- Anthropic. 2024b. [The Claude 3 Model Family: Opus, Sonnet, Haiku](#). Technical report, Anthropic. Accessed: 2024-05-23.
- Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang,

- and Juanzi Li. 2023. Longbench: A bilingual, multitask benchmark for long context understanding. *arXiv preprint arXiv:2308.14508*.
- Simone Balloccu, Patrícia Schmidová, Mateusz Lango, and Ondrej Dusek. 2024. **Leak, cheat, repeat: Data contamination and evaluation malpractices in closed-source LLMs**. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 67–93, St. Julian’s, Malta. Association for Computational Linguistics.
- Douglas Bates, Martin Mächler, Ben Bolker, and Steve Walker. 2015. **Fitting linear mixed-effects models using lme4**. *Journal of Statistical Software*, 67(1):1–48.
- Kent Chang, Mackenzie Cramer, Sandeep Soni, and David Bamman. 2023. **Speak, memory: An archaeology of books known to ChatGPT/GPT-4**. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7312–7327, Singapore. Association for Computational Linguistics.
- Tong Chen, Akari Asai, Niloofar Miresghallah, Sewon Min, James Grimmermann, Yejin Choi, Hannaneh Hajishirzi, Luke Zettlemoyer, and Pang Wei Koh. 2024. **Copybench: Measuring literal and non-literal reproduction of copyright-protected text in language model generation**. *Preprint*, arXiv:2407.07087.
- Manuel Ciosici, Joe Cecil, Dong-Ho Lee, Alex Hedges, Marjorie Freedman, and Ralph Weischedel. 2021. **Perhaps PTLMs should go to school – a task to assess open book and closed book QA**. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6104–6111, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Cohere. 2024a. **Command R: Retrieval-Augmented Generation at Production Scale**. Technical report, Cohere. Accessed: 2024-05-23.
- Cohere. 2024b. **Introducing Command R+: A Scalable LLM Built for Business**. Technical report, Cohere. Accessed: 2024-05-23.
- Zican Dong, Tianyi Tang, Junyi Li, Wayne Xin Zhao, and Ji-Rong Wen. 2023. **Bamboo: A comprehensive benchmark for evaluating long text modeling capacities of large language models**. *arXiv:2309.13345*.
- Alexander Fabbri, Chien-Sheng Wu, Wenhao Liu, and Caiming Xiong. 2022a. **QAFactEval: Improved QA-based factual consistency evaluation for summarization**. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2587–2601, Seattle, United States. Association for Computational Linguistics.
- Alexander Richard Fabbri, Chien-Sheng Wu, Wenhao Liu, and Caiming Xiong. 2022b. **Qafacteval: Improved qa-based factual consistency evaluation for summarization**. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2587–2601.
- Tobias Falke, Leonardo FR Ribeiro, Prasetya Ajie Utama, Ido Dagan, and Iryna Gurevych. 2019. **Ranking generated summaries by correctness: An interesting but challenging application for natural language inference**. In *Proceedings of the 57th annual meeting of the association for computational linguistics*, pages 2214–2220.
- Matt Gardner, Yoav Artzi, Victoria Basmov, Jonathan Berant, Ben Bogin, Sihao Chen, Pradeep Dasigi, Dheeru Dua, Yanai Elazar, Ananth Gottumukkala, Nitish Gupta, Hannaneh Hajishirzi, Gabriel Ilharco, Daniel Khashabi, Kevin Lin, Jiangming Liu, Nelson F. Liu, Phoebe Mulcaire, Qiang Ning, Sameer Singh, Noah A. Smith, Sanjay Subramanian, Reut Tsarfaty, Eric Wallace, Ally Zhang, and Ben Zhou. 2020. **Evaluating models’ local decision boundaries via contrast sets**. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1307–1323, Online. Association for Computational Linguistics.
- Gemini Team. 2024. **Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context**. *Preprint*, arXiv:2403.05530.
- Or Honovich, Roei Aharoni, Jonathan Herzig, Hagai Taitelbaum, Doron Kukliansy, Vered Cohen, Thomas Scialom, Idan Szepkator, Avinatan Hassidim, and Yossi Matias. 2022. **True: Re-evaluating factual consistency evaluation**. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3905–3920.
- Cheng-Ping Hsieh, Simeng Sun, Samuel Krizan, Shantanu Acharya, Dima Rekish, Fei Jia, Yang Zhang, and Boris Ginsburg. 2024. **Ruler: What’s the real context size of your long-context language models?** *Preprint*, arXiv:2404.06654.
- George Hudson and Noura Al Moubayed. 2022. **MuLD: The multitask long document benchmark**. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 3675–3685, Marseille, France. European Language Resources Association.
- Alon Jacovi, Avi Caciularu, Omer Goldman, and Yoav Goldberg. 2023. **Stop uploading test data in plain text: Practical strategies for mitigating data contamination by evaluation benchmarks**. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5075–5084, Singapore. Association for Computational Linguistics.
- Gregory Kamradt. 2023. **Needle In A Haystack - pressure testing LLMs**. *Github*.

- Sayash Kapoor and Arvind Narayanan. 2023. [Leakage and the reproducibility crisis in machine-learning-based science](#). *Patterns*, 4(9):100804.
- Marzena Karpinska, Nishant Raj, Katherine Thai, Yixiao Song, Ankita Gupta, and Mohit Iyyer. 2022. [DEMETER: Diagnosing evaluation metrics for translation](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9540–9561, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Yekyung Kim, Yapei Chang, Marzena Karpinska, Aparna Garimella, Varun Manjunatha, Kyle Lo, Tanya Goyal, and Mohit Iyyer. 2024. [Fables: Evaluating faithfulness and content selection in book-length summarization](#). *Preprint*, arXiv:2404.01261.
- Philippe Laban, Alexander R. Fabbri, Caiming Xiong, and Chien-Sheng Wu. 2024. [Summary of a haystack: A challenge to long-context llms and rag systems](#). *Preprint*, arXiv:2407.01370.
- Philippe Laban, Tobias Schnabel, Paul Bennett, and Marti A Hearst. 2022. [Summac: Re-visiting nli-based models for inconsistency detection in summarization](#). *Transactions of the Association for Computational Linguistics*, 10:163–177.
- Jinhyuk Lee, Anthony Chen, Zhuyun Dai, Dheeru Dua, Devendra Singh Sachan, Michael Boratko, Yi Luan, Sébastien M. R. Arnold, Vincent Perot, Siddharth Dalmia, Hexiang Hu, Xudong Lin, Panupong Pasupat, Aida Amini, Jeremy R. Cole, Sebastian Riedel, Iftekhar Naim, Ming-Wei Chang, and Kelvin Guu. 2024. [Can long-context language models subsume retrieval, rag, sql, and more?](#) *ArXiv*, abs/2406.13121.
- Russell V. Lenth. 2023. [emmeans: Estimated Marginal Means, aka Least-Squares Means](#). R package version 1.8.8.
- Mosh Levy, Alon Jacoby, and Yoav Goldberg. 2024. [Same task, more tokens: the impact of input length on the reasoning performance of large language models](#). *Preprint*, arXiv:2402.14848.
- Jiaqi Li, Mengmeng Wang, Zilong Zheng, and Muhan Zhang. 2023. [Loogle: Can long-context language models understand long contexts?](#) *Preprint*, arXiv:2311.04939.
- Tianle Li, Ge Zhang, Quy Duc Do, Xiang Yue, and Wenhui Chen. 2024. [Long-context llms struggle with long in-context learning](#). *Preprint*, arXiv:2404.02060.
- Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. [Lost in the middle: How language models use long contexts](#). *Transactions of the ACL*, 12:157–173.
- Potsawee Manakul, Adian Liusie, and Mark JF Gales. 2023. [Selfcheckgpt: Zero-resource black-box hallucination detection for generative large language models](#). *EMNLP*.
- Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. [On faithfulness and factuality in abstractive summarization](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1906–1919, Online. Association for Computational Linguistics.
- Sewon Min, Kalpesh Krishna, Xinxin Lyu, Mike Lewis, Wen-tau Yih, Pang Wei 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, EMNLP 2023, Singapore, December 6-10, 2023*, pages 12076–12100. Association for Computational Linguistics.
- Shinichi Nakagawa, Paul C. D. Johnson, and Holger Schielzeth. 2017. [The coefficient of determination<sup>2</sup> and intra-class correlation coefficient from generalized linear mixed-effects models revisited and expanded](#). *Journal of The Royal Society Interface*, 14(134):20170213.
- OpenAI. 2024a. [Gpt-4 technical report](#). *Preprint*, arXiv:2303.08774.
- OpenAI. 2024b. [Model release blog: GPT-4o](#). Technical report, OpenAI. Accessed: 2024-05-23.
- Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. 2023. [Yarn: Efficient context window extension of large language models](#). *Preprint*, arXiv:2309.00071.
- Adam Roberts, Colin Raffel, and Noam Shazeer. 2020. [How much knowledge can you pack into the parameters of a language model?](#) In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5418–5426, Online. Association for Computational Linguistics.
- Stephen E. Robertson, Steve Walker, Susan Jones, Micheline M. Hancock-Beaulieu, Mike Gatford, et al. 1995. [Okapi at trec-3](#). *NIST Special Publication SP*, 109:109.
- Oscar Sainz, Jon Campos, Iker García-Ferrero, Julen Etxaniz, Oier Lopez de Lacalle, and Eneko Agirre. 2023. [NLP evaluation in trouble: On the need to measure LLM data contamination for each benchmark](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 10776–10787, Singapore. Association for Computational Linguistics.
- Tal Schuster, Adam Fisch, and Regina Barzilay. 2021. [Get your vitamin C! robust fact verification with contrastive evidence](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 624–643, Online. Association for Computational Linguistics.
- Uri Shaham, Maor Ivgi, Avia Efrat, Jonathan Berant, and Omer Levy. 2023. [ZeroSCROLLS: A zero-shot](#)

- benchmark for long text understanding. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 7977–7989, Singapore. Association for Computational Linguistics.
- Chenglei Si, Navita Goyal, Sherry Tongshuang Wu, Chen Zhao, Shi Feng, Hal Daumé III au2, and Jordan Boyd-Graber. 2024. Large language models help humans verify truthfulness – except when they are convincingly wrong. *Preprint*, arXiv:2310.12558.
- Simeng Sun, Kalpesh Krishna, Andrew Mattarella-Micke, and Mohit Iyyer. 2021. Do long-range language models actually use long-range context? In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 807–822, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang, Sebastian Ruder, and Donald Metzler. 2021. Long range arena : A benchmark for efficient transformers. In *International Conference on Learning Representations*.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. FEVER: a large-scale dataset for fact extraction and VERification. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 809–819, New Orleans, Louisiana. Association for Computational Linguistics.
- Szymon Tworowski, Konrad Staniszewski, Mikołaj Pacek, Yuhuai Wu, Henryk Michalewski, and Piotr Miłoś. 2023. Focused transformer: Contrastive training for context scaling. *Preprint*, arXiv:2307.03170.
- Rob van der Goot. 2021. We need to talk about train-dev-test splits. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4485–4494, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Kiran Vodrahalli, Santiago Ontanon, Nilesh Tripuraneni, Kelvin Xu, Sanil Jain, Rakesh Shivanna, Jeffrey Hui, Nishanth Dikkala, Mehran Kazemi, Bahare Fatemi, Rohan Anil, Ethan Dyer, Siamak Shakeri, Roopali Vij, Harsh Mehta, Vinay Ramasesh, Quoc Le, Ed Chi, Yifeng Lu, Orhan Firat, Angeliki Lazariidou, Jean-Baptiste Lespiau, Nithya Attaluri, and Kate Olszewska. 2024. Michelangelo: Long context evaluations beyond haystacks via latent structure queries. *Preprint*, arXiv:2409.12640.
- David Wadden, Shanchuan Lin, Kyle Lo, Lucy Lu Wang, Madeleine van Zuylen, Arman Cohan, and Hannaneh Hajishirzi. 2020. Fact or fiction: Verifying scientific claims. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7534–7550, Online. Association for Computational Linguistics.
- Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020. Asking and answering questions to evaluate the factual consistency of summaries. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5008–5020.
- Cunxiang Wang, Ruoxi Ning, Boqi Pan, Tonghui Wu, Qipeng Guo, Cheng Deng, Guangsheng Bao, Qian Wang, and Yue Zhang. 2024a. Novelqa: A benchmark for long-range novel question answering. *Preprint*, arXiv:2403.12766.
- Minzheng Wang, Longze Chen, Cheng Fu, Shengyi Liao, Xinghua Zhang, Bingli Wu, Haiyang Yu, Nan Xu, Lei Zhang, Run Luo, Yunshui Li, Min Yang, Fei Huang, and Yongbin Li. 2024b. Leave no document behind: Benchmarking long-context llms with extended multi-doc qa. *Preprint*, arXiv:2406.17419.
- Alex Warstadt, Alicia Parrish, Haokun Liu, Anhad Mohananey, Wei Peng, Sheng-Fu Wang, and Samuel R. Bowman. 2020. BLiMP: The benchmark of linguistic minimal pairs for English. *Transactions of the Association for Computational Linguistics*, 8:377–392.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Jerry Wei, Chengrun Yang, Xinying Song, Yifeng Lu, Nathan Hu, Dustin Tran, Daiyi Peng, Ruibo Liu, Da Huang, Cosmo Du, and Quoc V. Le. 2024. Long-form factuality in large language models. *Preprint*, arXiv:2403.18802.
- Linhao Yu, Qun Liu, and Deyi Xiong. 2024. LFED: A literary fiction evaluation dataset for large language models. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 10466–10475, Torino, Italia. ELRA and ICCL.
- Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Khai Hao, Xu Han, Zhen Leng Thai, Shuo Wang, Zhiyuan Liu, and Maosong Sun. 2024.  $\infty$ bench: Extending long context evaluation beyond 100k tokens. *Preprint*, arXiv:2402.13718.

## A Note on dataset versions

Results on NOCHA are being updated as new models are released. The most recent results can be found at <https://novelchallenge.github.io/>. We are also actively collecting new claims/books and plan to update the benchmark periodically. While working on the collection of new claims we put more emphasis on reasoning over longer context and work only with the annotators who proved to

be the best at creating their claims in the initial study. At the time of writing about 10 new books were added to the collection.

## B Issues with current long-context evaluation benchmarks

### B.1 Evaluation may be conducted on the training data

**Data contamination:** One of the most prominent issues in current evaluation practices is the test data. A fundamental rule in evaluating any trained model is to test it on a separate withheld *test* set, avoiding the instances the model was trained on (van der Goot, 2021; Kapoor and Narayanan, 2023). However, this has become increasingly difficult as existing test sets have almost certainly been leaked into the training data, and most available LLMs are either entirely closed-source or only open-weights, meaning we do not have access to the training instances (Sainz et al., 2023; Balloccu et al., 2024). While analyzing existing long-context benchmarks, though this issue is certainly not limited to long-context setups, we noticed that many either utilize existing datasets (e.g., Bai et al., 2023; An et al., 2024) or use old texts due to copyright issues (e.g., Li et al., 2023; Wang et al., 2024a), both of which are likely included in the models' training data. Some researchers have tried to mitigate this issue by transforming the original texts, for instance, by replacing the named entities with others (Zhang et al., 2024; Dong et al., 2023). However, we observe that prompting models with such data (even fragments of books) still leads to recognition of the source materials, an outcome in line with results reported in Chang et al. (2023).

**Issues with closed-book tests:** A popular method used to ensure that the training data do not affect the model's performance is to perform the *closed-book* test (Roberts et al., 2020; Ciosici et al., 2021), used, for instance, in Wang et al. (2024a). Just like a student taking a test without access to the study book, the model is tested without access to the source text. The idea is that if the model has memorized the training data, it will be able to perform well even without access to the source text. However, as was recently pointed out, memorization does not have to be perfect to impact generation (Chen et al., 2024). For instance, even a model that fails a *closed-book* test might have encountered the source text in its training data and

may benefit from this during the test when presented with the task along with the source text as the context. Since little is known about memorization in LLMs, the safest evaluation is to test the model on newly produced data (e.g., novels published past the model's cutoff date). However, even this approach necessitates constant updates to the benchmark.

### B.2 Employing LLMs for evaluation and data creation

**Unreliability of LLM-based evaluation:** Another issue is the way in which model outputs are evaluated. While multiple-choice questions (including True/False evaluations) are relatively straightforward to verify if all plausible answers are correctly marked, they are also likely to be easier for models to answer correctly. Consequently, researchers have developed datasets requiring models to generate long-form answers or perform complex tasks involving long contexts, such as summarization (Hudson and Al Moubayed, 2022; Bai et al., 2023; Zhang et al., 2024, inter alia). These setups are more challenging for the models but also harder to evaluate. Currently, no reliable automatic metric exists for verifying long-form outputs, leading some researchers to use language models, like GPT-4o, for evaluation (e.g., Wang et al., 2024b). This method has several flaws.

First, the model typically evaluates the produced answer against a gold standard, which may penalize valid answers that include additional, relevant information not covered by the gold answer. Second, inconsistent instructions given to the evaluator model can result in variability in its judgments. For example, when the model is asked to assign a score between 0 and 100, with 100 indicating a perfect response (as in Wang et al., 2024b), intermediate scores may fluctuate due to inconsistent deductions for similar errors across different evaluations. Even when using a simpler scale, such as a 5-point Likert scale, the model's evaluations can still be unreliable. Unlike a human expert, a language model does not "build experience" through repeated evaluations of similar outputs and cannot refine its understanding or calibration of the scale over time. Human evaluators, in contrast, assess responses not only in relation to a predefined standard but also by drawing on their experience with a variety of previous responses, making their judgments more consistent and dependable. This inherent limitation in language model-based evaluation underscores

the difficulty in achieving reliable automated assessment for complex tasks.

**Creating test data with language models:** Since employing human annotators to create test instances (e.g., questions about the documents) is time- and resource-intensive (i.e., the annotators have to read long texts),<sup>27</sup> some researchers use language models to generate these instances. However, this approach has its own issues. For instance, using models such as GPT-4o to generate questions and answers (as in Wang et al., 2024b) may result in flawed test data or test data that is inherently easier for a model from the same family to solve.<sup>28</sup> This becomes even more complex when weaker models are used to generate test examples for relatively complex tasks. For instance, Li et al. (2023) employs GPT-3.5-TURBO to generate summaries for a cloze task, where the tested model is supposed to fill in masked named entities based on the source text. However, GPT-3.5-TURBO has been shown to be a poor summarizer, even when summarizing smaller chunks of text (Kim et al., 2024).

## C Dataset

In this section of the appendix we provide more details on our corpus (§C.1) and human annotation efforts (§C.2).

### C.1 Corpus

This section provides detailed information about the corpus collected for this study. Table 11 lists all the novels included in the corpus,<sup>29</sup> with the genre distribution illustrated in Figure 5.<sup>30</sup> Table 12 lists all stories included in the story collections. Additionally, we provide the the statistics for books and claims by year of publication in Table 14. Finally, we report list the statistics for stories from story collections in Table 13.

### C.2 Human annotation efforts

In this section of the appendix we provide additional details about the data collection pipeline.

<sup>27</sup>Reading an average book takes at least 8-10 hours.

<sup>28</sup>In fact, on a benchmark where Q&A pairs were generated by GPT-4o, the model achieved one of the highest scores among all tested models (Wang et al., 2024b).

<sup>29</sup>All books in the corpus were purchased by the researchers.

<sup>30</sup>Please note that multiple genre tags are allowed per book, as a book can belong to more than one genre, such as both romance and mystery.

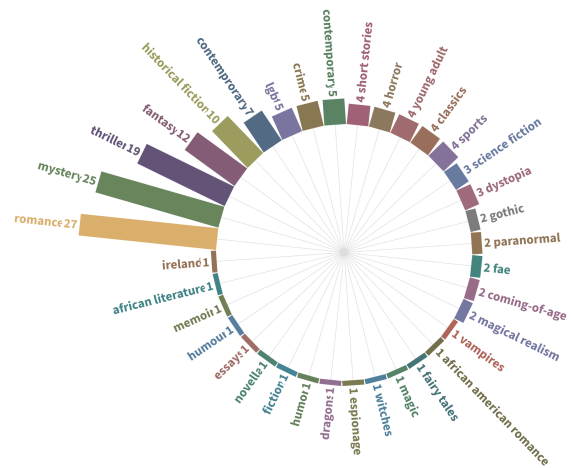


Figure 5: Genre distribution in NOCHA. As a book can belong to multiple genres, such as fantasy and romance, we allow up to three labels per book.

**Annotators:** The annotations were done by 18 annotators recruited on Upwork (all female) and 5 volunteer annotators, including three of the authors (four females and one male). All but four annotators are native English speakers from either the US or UK. All annotators hold university degrees, with one being a former college English professor and an author, and one holding a Ph.D. in English literature. As frequent readers, some annotators received advanced reader copies of the books, enabling the annotation of works, which were unpublished at the time of data collection.

**Collecting True/False claim pairs:** Before starting the annotation task, all annotators were required to read the guidelines (Figure 6) and sign a consent form (Figure 7). The collected claims underwent a rigorous review process, being checked at least three times by two of the authors: initially during the writing stage and later independently during the final quality control phase. During this phase, we also proofread all claims and contacted annotators if anything was unclear.

Collecting each pair of annotations typically required additional communication with the annotators, resulting in approximately 160 hours of work from each of the two authors involved in this process. Based on self-reported time and records in spreadsheets, we estimate that the annotators were able to produce between 6 and 10 claim pairs per hour.<sup>31</sup>

<sup>31</sup>The annotators we also asked to provide written justification as to why the true claim is true and the false claim is false relating to the events in the book.

TITLE	AUTHOR	GENDER	GENRE	PUB. DATE	TOKENS	WORDS	# PAIRS	LANG.
<i>A Haunting on the Hill</i>	Elizabeth Hand	F	horror, Gothic, paranormal	Oct 3, 2023	116,411	87,391	20	AmE
<i>Assassin Eighteen</i>	John Brownlow	M	thriller, crime	Apr 23, 2024	126,011	96,830	15	CanE
<i>Come and Get It</i>	Kiley Reid	F	contemporary	Jan 30, 2024	141,484	104,373	30	AmE
<i>Curse of the Soul Collector</i>	Cara Blaine	F	fantasy	Sep 20, 2023	85,012	67,495	15	AmE
<i>Death Comes to Marlow</i>	Robert Thorogood	M	mystery, fiction, crime	Jan 5, 2023	113,544	87,844	15	BrE
<i>Divine Rivals</i>	Rebecca Ross	F	romance, young adult, historical fiction	Apr 4, 2023	135,800	106,371	8	AmE
<i>Everyone on the train is a suspect</i>	Benjamin Stevenson	M	mystery	Oct 17, 2023	108,673	82,296	7	AusE
<i>First Lie Wins</i>	Ashley Elston	F	thriller, mystery	Jan 2, 2024	117,547	97,151	15	AmE
<i>Fourth Wing</i>	Rebecca Yarros	F	fantasy, romance, dragons	May 2, 2023	242,683	183,789	9	AmE
<i>Funny Story</i>	Emily Henry	F	romance, contemporary	Apr 23, 2024	138,892	104,775	14	AmE
<i>Helpless</i>	Kelby C. Hughes	F	fantasy	Mar 25, 2024	65,677	53,381	15	AmE
<i>Home Is Where the Bodies Are</i>	Jeneva Rose	F	thriller, mystery, horror	Apr 30, 2024	97,390	76,866	14	AmE
<i>House of Flame and Shadow</i>	Sarah J. Maas	F	fantasy, romance, fae	Jan 30, 2024	334,581	248,192	15	AmE
<i>How to solve your own murder</i>	Kristen Perrin	F	mystery, thriller, crime	Mar 26, 2024	130,414	104,156	10	AmE
<i>I Hope This Doesn't Find You</i>	Ann Liang	F	romance, young adult	Feb 6, 2024	105,934	81,432	9	AusE
<i>Inheritance</i>	Nora Roberts	F	paranormal, romance, mystery	Nov 21, 2023	170,979	127,511	40	AmE
<i>Leaving</i>	Roxana Robinson	F	romance, contemporary	Feb 13, 2024	133,317	101,039	10	AmE
<i>Long Island</i>	Colm Toibin	M	historical fiction, Ireland	May 7, 2024	103,361	84,339	10	IrE
<i>Love, Lies, and Cherry Pie</i>	Jackie Lou	F	contemporary, romance	May 7, 2024	112,179	86,994	15	AmE
<i>Monstrous Alterations</i>	Christopher Barzak	M	short stories, fantasy, LGBT	Sep 8, 2023	73,094	59,717	10	AmE
<i>Only for the week</i>	Natasha Bishop	F	African American romance	May 11, 2023	85,969	68,517	5	AmE
<i>Pet</i>	Catherine Chidgey	F	thriller, mystery	July 13, 2023	124,189	93,294	46	NZE
<i>Random in Death</i>	J. D. Robb	F	mystery, romance, crime	Jan 23, 2024	132,107	97,172	30	AmE
<i>Romantic Comedy</i>	Curtis Sittenfeld	F	romance, contemporary	Apr 4, 2023	115,004	88,775	12	AmE
<i>Roux</i>	Tamika Christy	F	romance, historical fiction, LGBT	Jan 9, 2024	121,364	89,221	15	AmE
<i>Ruthless Vows</i>	Rebecca Ross	F	fantasy, romance, young adult	Dec 26, 2023	161,337	127,090	15	AmE
<i>Safe and Sound</i>	Laura McHugh	F	thriller, mystery	Apr 23, 2024	103,054	79,243	10	AmE
<i>Same Time Next Year</i>	Tessa Bailey	F	romance, sports, novella	Apr 18, 2023	49,156	38,023	35	AmE
<i>She Is a Haunting</i>	Trang Thanh Tran	F	horror, young adult, LGBT	Feb 28, 2023	105,378	80,291	41	AmE
<i>Six scorched roses</i>	Carissa Broadbent	F	fantasy, romance, vampires	Mar 21, 2023	51,779	40,062	6	AmE
<i>The Agency for Scandal</i>	Laura Wood	F	historical fiction, romance, mystery	Jan 5, 2023	115,382	90,741	18	AmE
<i>The Atonement Murders</i>	Jennifer Ruff	F	mystery, thriller	Apr 14, 2023	104,258	82,134	4	AmE
<i>The Beautiful and the Wild</i>	Peggy Townsend	F	thriller, mystery	Nov 7, 2023	92,424	75,908	16	AmE
<i>The book of love</i>	Kelly Link	F	mystery, magical realism, thriller	Feb 13, 2024	272,343	209,950	10	AmE
<i>The Bootleggers Daughter</i>	Nadine Nettman	F	mystery, thriller	May 1, 2024	97,386	75,696	15	AmE
<i>The eye of the bedlam bride</i>	Matt Dinman	M	fantasy, science fiction, humor	Jul 2, 2023	336,288	257,445	15	AmE
<i>The Future</i>	Naomi Alderman	F	science fiction, dystopia, fantasy	Nov 7, 2023	157,019	122,342	14	BrE
<i>The Glass Woman</i>	Alice McLroy	F	historical fiction, Gothic, mystery	Dec 19, 2023	99,113	78,795	15	BrE
<i>The Guest</i>	Emma Cline	F	mystery, thriller	May 16, 2023	89,042	68,735	20	AmE
<i>The Hanging City</i>	Charlie Holmberg	F	fantasy, romance, magic	Aug 1, 2023	140,622	106,182	8	AmE
<i>The Heiress</i>	Rachel Hawkins	F	mystery, thriller	Jan 9, 2024	94,406	73,954	15	AmE
<i>The Husbands</i>	Holly Gramazio	F	romance, magical realism	Apr 2, 2024	130,017	100,432	15	AusE
<i>The Last Murder at the End of the World</i>	Stuart Turton	M	mystery, thriller, dystopia	May 21, 2024	122,364	91,280	15	BrE
<i>The Library of Borrowed Hearts</i>	Lucy Gilmore	F	contemporary, romance	Apr 30, 2024	132,205	105,259	15	AmE
<i>The Limits: A Novel</i>	Nell Freudenberger	F	contemporary	Apr 9, 2024	131,766	102,133	15	AmE
<i>The Marriage Act</i>	John Marrs	M	thriller, science fiction, dystopia	Jan 19, 2023	134,493	106,341	16	BrE
<i>The Promise of Tomorrow</i>	Mary Ellen Taylor	F	contemporary, romance	June 1, 2024	116,035	89,682	15	AmE
<i>The Prospects</i>	KT Hoffman	F	romance, LGBT, sports	Apr 9, 2024	122,720	94,591	15	AmE
<i>The Resort</i>	Sarah Ochs	F	thriller, mystery	Feb 6, 2024	124,200	99,627	15	AmE
<i>The Spy Coast</i>	Tess Gerritsen	F	mystery, thriller, espionage	Nov 1, 2023	126,906	98,326	8	AmE
<i>The Tainted Cup</i>	Robert Jackson Bennett	M	mystery, thriller	Feb 6, 2024	162,289	119,350	10	AmE
<i>The Teacher</i>	Freida McFadden	F	thriller, mystery	Feb 6, 2024	110,260	88,104	15	AmE
<i>The White Lady</i>	Jacqueline Winspear	F	historical fiction, mystery	Mar 21, 2023	124,585	97,094	4	BrE
<i>Viciously Yours</i>	Jaime Applegate	F	fantasy, romance, fae	Jan 23, 2024	85,118	65,986	7	AmE
<i>Weyward</i>	Emilia Hart	F	historical fiction, fantasy, witches	Feb 2, 2023	126,501	98,842	9	AusE
<i>White Cat, Black Dog</i>	Kelly Link	F	short stories, horror, fairy tales	Mar 2, 2023	92,826	74,154	10	AmE
<i>Wildfire</i>	Hannah Grace	F	romance, sports, contemporary	Oct 3, 2023	138,441	109,641	5	BrE
<i>Yellowface</i>	R.F. Kuang	F	thriller, contemporary, mystery	May 25, 2023	113,164	87,648	12	AmE
<i>You Should Be So Lucky</i>	Cat Sebastian	F	romance, LGBT, sports	May 7, 2024	140,422	109,764	15	AmE
<i>You, Again</i>	Kate Goldbeck	F	romance, contemporary	Sep 12, 2023	128,445	96,681	7	AmE
<i>Yours truly</i>	Abby Jimenez	F	romance, contemporary	April 11, 2023	134,609	105,968	10	AmE
<i>Quietly Hostile</i>	Samantha Irby	F	essays, humour, memoir	May 16, 2023	95,842	76,424	14	AmE
<i>Anne of Green Gables</i>	L. M. Montgomery	F	classics, coming-of-age, historical	Jan 1, 1908	129,908	102,366	15	CanE
<i>Little Women</i>	Louisa May Alcott	F	classics, coming-of-age, historical, romance	Sep 30, 1868	235,118	185,930	15	AmE
<i>The Adventures of Sherlock Holmes</i>	Arthur Conan Doyle	M	classics, short stories, mystery, crime	Oct 14, 1892	129,293	104,434	18	BrE
<i>The Great Gatsby</i>	F. Scott Fitzgerald	M	classics, historical fiction, romance	Apr 10, 1925	61,689	48,187	15	AmE

Table 11: List of novels included in NOCHA. The upper row presents all recent fictional books included in the data. The middle row displays the essay collections, while the lower row shows classical books included in the data. The genre is provided as listed on GoodReads (<https://www.goodreads.com/>). The language is indicated based on the author’s native language or, for non-native English authors, the primary language of the country where they spent most of their time. The token count is reported as per tiktoken tokenization with cl100k encoding while the word count was determined by a whitespace split.

Overall, we collected about 15 claim pairs per book. While we initially aimed to collect more claim pairs per book, we observed that creating meaningful and challenging pairs becomes significantly more difficult beyond the first 10-15 pairs, though this number may vary slightly depending on the book. For some books, we collected fewer than 10 pairs due to the unavailability of annotators

to create more pairs.

**Advantages of a minimal pair design** We employ minimal pairs for two main reasons. Firstly, it ensures data quality, allowing us to verify the false claim against its true counterpart easily and identify cases where both claims were too similar (i.e., the false claim could potentially be true) or subjective.



ID	COLLECTION 📖	AUTHOR 🧑	STORY TITLE	# PAIRS	STORY TOKENS	COLLECTION TOKENS	LOCATION	PUB. DATE 📅
1	A Kind of Madness	Uche Okonkwo	Nwunye Belgium	2	10299	80946	<i>beginning</i>	Apr 16, 2024
2	A Kind of Madness	Uche Okonkwo	Shadow	3	8549	80946	<i>beginning</i>	Apr 16, 2024
3	A Kind of Madness	Uche Okonkwo	Debris	3	1822	80946	<i>beginning</i>	Apr 16, 2024
4	A Kind of Madness	Uche Okonkwo	Long Hair	3	3464	80946	<i>beginning</i>	Apr 16, 2024
5	A Kind of Madness	Uche Okonkwo	Animals	3	8193	80946	<i>beginning</i>	Apr 16, 2024
6	A Kind of Madness	Uche Okonkwo	Milk and Oil	3	11829	80946	<i>middle</i>	Apr 16, 2024
7	A Kind of Madness	Uche Okonkwo	The Harvest	2	7145	80946	<i>middle</i>	Apr 16, 2024
8	A Kind of Madness	Uche Okonkwo	Eden	2	7408	80946	<i>middle</i>	Apr 16, 2024
9	A Kind of Madness	Uche Okonkwo	The Girl Who Lied	3	12834	80946	<i>end</i>	Apr 16, 2024
10	A Kind of Madness	Uche Okonkwo	Burning	3	9524	80946	<i>end</i>	Apr 16, 2024
1	The Adventures of Sherlock Holmes	Arthur Conan Doyle	A Scandal in Bohemia	2	10715	129293	<i>beginning</i>	Oct 14, 1892
2	The Adventures of Sherlock Holmes	Arthur Conan Doyle	The Red Headed League	1	11369	129293	<i>beginning</i>	Oct 14, 1892
3	The Adventures of Sherlock Holmes	Arthur Conan Doyle	A Case of Identity	1	8715	129293	<i>beginning</i>	Oct 14, 1892
4	The Adventures of Sherlock Holmes	Arthur Conan Doyle	The Boscombe Valley Mystery	3	11748	129293	<i>beginning</i>	Oct 14, 1892
5	The Adventures of Sherlock Holmes	Arthur Conan Doyle	The Five Orange Pips	2	9076	129293	<i>beginning</i>	Oct 14, 1892
6	The Adventures of Sherlock Holmes	Arthur Conan Doyle	The Man with the Twisted Lip	1	11401	129293	<i>middle</i>	Oct 14, 1892
7	The Adventures of Sherlock Holmes	Arthur Conan Doyle	The Adventure of Blue Carbuncle	1	9982	129293	<i>middle</i>	Oct 14, 1892
8	The Adventures of Sherlock Holmes	Arthur Conan Doyle	The Adventure of the Speckled Band	2	12070	129293	<i>middle</i>	Oct 14, 1892
9	The Adventures of Sherlock Holmes	Arthur Conan Doyle	The Adventures of the Engineers Thumb	1	10227	129293	<i>middle</i>	Oct 14, 1892
10	The Adventures of Sherlock Holmes	Arthur Conan Doyle	The Adventure of the Noble Bachelor	1	10086	129293	<i>end</i>	Oct 14, 1892
11	The Adventures of Sherlock Holmes	Arthur Conan Doyle	The Adventure of the Beryl Coronet	2	11776	129293	<i>end</i>	Oct 14, 1892
12	The Adventures of Sherlock Holmes	Arthur Conan Doyle	The Red Adventure of the Copper Beeches	1	12331	129293	<i>end</i>	Oct 14, 1892
1	Quietly Hostile	Samantha Irby	i like it!	0	1783	95842	<i>beginning</i>	May 16, 2023
2	Quietly Hostile	Samantha Irby	the last normal day	1	4188	95842	<i>beginning</i>	May 16, 2023
3	Quietly Hostile	Samantha Irby	david matthews's greatest romantic hits	1	4563	95842	<i>beginning</i>	May 16, 2023
4	Quietly Hostile	Samantha Irby	chub street diet	1	4805	95842	<i>beginning</i>	May 16, 2023
5	Quietly Hostile	Samantha Irby	my firstborn dog	1	4672	95842	<i>beginning</i>	May 16, 2023
6	Quietly Hostile	Samantha Irby	body horror!	1	4798	95842	<i>beginning</i>	May 16, 2023
7	Quietly Hostile	Samantha Irby	two old nuns having amzing [sic] lesbian sex	1	7138	95842	<i>beginning</i>	May 16, 2023
8	Quietly Hostile	Samantha Irby	qvc, ilysm	1	5974	95842	<i>beginning</i>	May 16, 2023
9	Quietly Hostile	Samantha Irby	superfan!!!!!!	1	11397	95842	<i>middle</i>	May 16, 2023
10	Quietly Hostile	Samantha Irby	i like to get high at night and think about whales	0	702	95842	<i>middle</i>	May 16, 2023
11	Quietly Hostile	Samantha Irby	oh, so you actually don't wanna make a show about a horny fat bitch with diarrhea? okay!	1	12351	95842	<i>middle</i>	May 16, 2023
12	Quietly Hostile	Samantha Irby	what if i died like elvis	1	9068	95842	<i>middle</i>	May 16, 2023
13	Quietly Hostile	Samantha Irby	shit happens	1	4960	95842	<i>end</i>	May 16, 2023
14	Quietly Hostile	Samantha Irby	food fight	0	1508	95842	<i>end</i>	May 16, 2023
15	Quietly Hostile	Samantha Irby	o brother, where art thou?	1	4865	95842	<i>end</i>	May 16, 2023
16	Quietly Hostile	Samantha Irby	how to look cool in front of teens?	1	5723	95842	<i>end</i>	May 16, 2023
17	Quietly Hostile	Samantha Irby	we used to get dressed up to go to red lobster	1	5558	95842	<i>end</i>	May 16, 2023
18	Quietly Hostile	Samantha Irby	please invite me to your party	0	2049	95842	<i>end</i>	May 16, 2023
1	White Cat, Black Dog	Kelly Link	The White Cat's Divorce	2	11697	92826	<i>beginning</i>	Mar 2, 2023
2	White Cat, Black Dog	Kelly Link	Prince Hat Underground	2	19724	92826	<i>beginning</i>	Mar 2, 2023
3	White Cat, Black Dog	Kelly Link	The White Road	2	11770	92826	<i>middle</i>	Mar 2, 2023
4	White Cat, Black Dog	Kelly Link	The Girl Who Did Not Know Fear	1	8606	92826	<i>middle</i>	Mar 2, 2023
5	White Cat, Black Dog	Kelly Link	The Game of Smash and Recovery	0	6869	92826	<i>middle</i>	Mar 2, 2023
6	White Cat, Black Dog	Kelly Link	The Lady and the Fox	1	12504	92826	<i>middle</i>	Mar 2, 2023
7	White Cat, Black Dog	Kelly Link	Skinder's Veil	2	21734	92826	<i>end</i>	Mar 2, 2023

Table 12: List of stories included in the collections. We provide token counts for both the entire collection and individual story. The location was determined by dividing collection length into three parts, where the stories which begin in the first part are marked as *beginning*, stories which begin in the second part are marked as *middle*, and stories beginning in the third part are marked as *end*.

	STORIES
MEAN	8,501.5
MAX	21,734
MIN	702
STD. DEV.	4,355.7

Table 13: Statistics of stories in story collections (in tokens). Note that there is 7–18 stories per collection.

Secondly, we require models to correctly label *both* claims in a pair, which minimizes the chances of crediting them for correct predictions made "for the wrong reason" without proper utilization of the context.

Indeed, the results presented in the paper indicate that merely having a balanced set of labels is insufficient for accurately evaluating the models (see §4 and examples in Figure 2). For example, using balanced but unpaired data might lead us to conclude that GPT-4o achieves an overall accuracy

of 76.7% (77.5% for True claims and 75.9% for False claims; see Table 10). However, these results are misleading because the model fails to validate *both* claims correctly in approximately 20% of cases, resulting in an actual accuracy of 55.8%. This is significantly lower than the human accuracy of 96.9% reported in §2, highlighting a substantial performance gap for the best model to bridge.

**Quality control:** To ensure the quality of the annotations, we reannotated a subset of the data, consisting of 128 claims made about 6 books (76 claim pairs). For this task, we hired annotators who had read the same books.<sup>32</sup> All but one of these annotators were also part of the original annotation team but had worked on different books previously.<sup>33</sup> They were provided with *the same*

<sup>32</sup>We hired one new annotator per book which resulted in two annotations per book, the original annotator (author of the claims and labels) and the new annotator.

<sup>33</sup>One person who did not do the original annotation was a volunteer who have also read of the annotated books. This

	Books 📖			Claim [Pairs] 📄		
	2024 (n=33)	2023 (n=30)	Classics (n=4)	2024 (n=898 [449])	2023 (n=978 [489])	Classics (n=126 [63])
TOKENS (tiktoken)						
Mean	129,526	123,908	139,002	23.51	22.75	24.71
St. dev.	52,112	52,176	71,630	7.17	7.94	7.84
Max	334,581	336,288	235,118	63	57	45
Min	65,677	49,156	61,689	5	8	10
WORDS (whitespace)						
Mean	99,753	96,117	110,229	18.44	17.85	20.06
St. dev.	38,463	39,388	56,790	6.15	6.68	6.97
Max	248,192	257,445	185,930	57	46	39
Min	53,381	38,023	48,187	4	5	7

Table 14: Number of tokens and words across books and claims for books published in 2024, 2023, and classics. The number of claim pairs is indicated in [] square brackets.

## Guidelines

### 1. Why are we collecting this data?

We are a university research lab based at the University of Massachusetts, Amherst led by Mohit Iyer. You were most likely contacted by Katherine or Marzena. Currently, we are investigating the ability of large language models (e.g., ChatGPT, Claude) to verify claims based on fictional books.

We are interested in creating a dataset of **True** and **False** claim pairs which will be tricky for the model to verify. We observe that the models struggle especially with claims which are **False** as a whole but **contain some degree of truth**. Below you will find an example of such a claim.

**Claim:** Martha marries Jonathan Strong, an art broker, **despite** her romantic ties with Patrick.

**Explanation:** Martha marries Jonathan Strong, but it is **before** she becomes romantically involved with Patrick.

**Evaluation:** **False**

**NOTE:** If your book is part of a larger series, please make sure that all your claims can be evaluated **ONLY** after reading this concrete book.

(a) Guidelines: Page 1

### 2. What is your task?

Follow the link you have received on the Upwork and create **True-False** claim pairs based on the books you have read. Here are some restrictions:

- (1) Each claim should be roughly **one-sentence long**. Avoid claims spanning multiple sentences, if possible. Length like the claim below is totally fine (preferable to a shorter claim):

"Sally and Noah's relationship faces hurdles including a public break-up between Danny and Annabel, the serious illness of Sally's former partner Jerry, and leaked private photos of Sally and Noah."

- (2) The **true** claim should be **indisputably true**.
- (3) The **false** claim should be **indisputably false as a whole**. It should also contain some **true statement(s)** in it. It should be sufficiently long and **hard to verify**. **Please avoid random entity replacement, this does not make a hard claim.**
- (4) **False** and **true** claims should form a **minimal pair** (twin pair), that is they should be about the same event(s)/state(s)/part(s) of the book, and **differ only by the false information injected into the true claim.**

(b) Guidelines: Page 2

Figure 6: Guidelines provided to the annotators for the annotation task. The annotators were also provided additional examples and guidance during the data collection process.

instructions as the prompt given to the model (Table 17), with the only modification being the phrase "based on the context provided" changed to "based on your book." All annotations were performed using GoogleSpreadsheets.

Overall, the annotators agreed with the original annotator on 148 out of 152 labels (72 out of 76 claim pairs). We then reviewed the two pairs where the annotators disagreed. In one pair, the disagreement was due to a sentence-level detail about whether the main characters were engaged for two or three years. The characters were together for three years but engaged for only two, a detail mentioned in the book only once. The sec-

person read the book for this task specifically and performed their annotations right after finishing the book.

ond annotator *incorrectly* annotated the claim of three years as true and the two years claim as false leading to disagreement. In the second case, the claim was about "The Great Gatsby," where the initial annotator confused parts of the book with parts of the movie. Importantly, none of the core test books have movies made based on them as they were published within the last year, so this situation was specific to this classic novel. We have since reviewed all the classic claims again to ensure this issue does not recur.

**Evidence scope:** For this task, we asked four annotators to annotate the scope of the evidence for two of their books each, resulting in annotations for 121 claim pairs from 8 books (approximately 15 claim pairs per book; see Figure 8 for instructions

SCOPE	PAIRS (%)	TRUE (%)	FALSE (%)
GLOBAL	47.9%	38.0%	47.1%
PASSAGE	39.7%	43.8%	39.7%
SENTENCE	12.4%	18.2%	13.2%

Table 15: Percentage of pairs by scope in the annotated subset of data (*global*, *passage*, *sentence*). Additionally, we report the percentage of *claims* with specific labels for both True and False claims.

### Collecting a dataset for claim verification

**Purpose of the task:** The goal of this research is to evaluate how well large language models (such as ChatGPT) can verify claims made about long input (fictional books).

You will be asked to (1) read the guidelines, (2) access the google spreadsheet shared with you, and (3) create True-False claim pairs about the book you are familiar with. We will also ask you to provide a short (1-2 sentences) explanation of what is incorrect in the False claim.

**No personally identifiable information will be collected or utilized for our analysis.**

By signing this consent, I acknowledge that:

- I voluntarily agree to participate in this research study.
- I understand that I will be paid \$45 for creating 15 claim pairs and \$100 for creating 30 claim pairs, unless agreed otherwise via Upwork.
- I have been informed of the purpose and nature of the study and I have had the opportunity to ask questions about the study. I understand that I also have the right to ask questions during the task.
- I understand that participation involves:
  - Creating True-False claim pairs about the book I have read.
  - Writing a short note explaining what is incorrect in the False claim.
- I understand that all information I provide for this study will be treated confidentially.
- I understand that in any report on the results of this research **my identity will remain anonymous**.

Please sign below if you have read the above terms and fully agree with them.

---

I have read and understood the guidelines. I understand that I can ask any additional questions now or during the task via UpWork. \*

Yes

No

---

I agree that all the data I create (True and False claims, comments, validation and other annotations) can be used in academic experiments. \*

Yes

No

---

Signature \*

Short answer text

Figure 7: Consent form which the annotators were asked to sign via GoogleForms before collecting the data.

For each claim, please tell us how much reasoning is needed to correctly validate it as either True or False. Overall, we are interested to see how much understanding/reasoning is needed to label each pair correctly (True as True and False as False). Please choose one of the following options:

(1) **sentence** - use this category if the claim can be directly validated by one or two consecutive sentences within the text. The content of these sentences alone should provide all the necessary information to affirm or refute the claim.

**Example:** the following fragment "Hi Sonya! Welcome to the manor. I'm Trey, Ace's son." can fully validate the following claim: "Trey, Ace's son, welcomes Sonya to the manor."

(2) **passage** - use this category if the claim requires a bit more text for validation, specifically up to three-four paragraphs. This could be a section where several sentences together provide the context needed to verify the claim.

**Example:** the book contains a short conversation between Sonya and Trey during which we learn everything we need to validate the following claim: "Trey, Ace's son, welcomes Sonya to the manor."

(3) **longer context** -- this category is appropriate when the information required to validate the claim is scattered throughout the text. Different aspects of the claim may be supported by different sections of the book, and piecing these together is necessary for validation.

Figure 8: Instructions given to the annotators for the annotation of claim scope.

given to the annotators). The annotators first annotated the scope for each claim in the pair, and then the maximum scope was taken for the given pair.<sup>34</sup> Overall, 12.4% of pairs were labeled as requiring only one to two sentences to validate, 39.7% of pairs were labeled as requiring a longer passage, and 47.9% were labeled as requiring global reasoning. The distribution of labels for specific claims is comparable (see Table 15).

**Note on evidence location:** Designing a realistic experiment to verify complex claims with evidence at different depths of a book is challenging. Unlike the "needle-in-a-haystack" experiment, where the "needle" can be easily placed in different parts of the (often unrelated) context, our task involves realistic claims about the book's content, making it more problematic to control for various confounders. For instance, a claim about events at the beginning of the book might be inherently eas-

<sup>34</sup>For instance, if one claim in a pair is labeled as "passage" and the other as "global," the resulting scope for the entire pair would be "global." This is because we require the models to validate *both* claims in the pair correctly, meaning the models would have to reason over both "passage" and "global" contexts for this specific pair. It is worth noting that the labels for both claims in the pair matched 87.6% of the time.

ier to verify. Similarly, it is difficult to assert with certainty that the only piece of evidence necessary to validate the claim is present at a specific depth in the book, as there could be corroborating evidence prior or later that aids verification (e.g., hints about the killer’s identity), which the human reader may (reasonably) not pay attention to. Hence, we do not attempt to annotate the evidence location. Instead, we use story collections, where we can precisely identify which story the claim pertains to and where the story is located within the collection. We present the results of this experiment both in the main body of this paper (§4) and later in this appendix (§E).

## D Methods

In this section we provide additional details about our evaluation methodology.

**Models:** We evaluate 6 closed-source and 5 open-source long-context models (i.e., models with context window of at least 128k tokens):

- **🔒 Closed-source models:** GPT-4o (OpenAI, 2024b), GPT-4-TURBO (OpenAI, 2024a), CLAUDE-3-OPUS (Anthropic, 2024b), CLAUDE-3.5-SONNET (Anthropic, 2024a), GEMINI PRO 1.5 (Gemini Team, 2024), GEMINI FLASH 1.5 (Gemini Team, 2024). All models were accessed via provider’s API.<sup>35</sup>
- **🏠 Open-weight models:** COMMAND R (Cohere, 2024a), COMMAND R+ (Cohere, 2024b), GEMMA-10M (Aljadery et al., 2024), PHI-3-MINI (Abdin et al., 2024), and LONGLLAMA (Tworkowski et al., 2023). COMMAND R and COMMAND R+ we accessed via cohere API. Other models were run on up to three A100 80GB GPUs. For comparison, in our stories experiment, we also tested MIXTRAL-8X22B (65k context window) and QWEN-2-72B (32k context window). While their context windows are too short to process entire books, they can handle stories from our collections. Both models were run using the together.ai API.
- **📖 Retrieval pipeline:** We also implement a retrieval pipeline (BM25+GPT-4o), where we

<sup>35</sup>CLAUDE-3-OPUS’s generations were done partially employing anthropic API and partially utilizing vertex-ai due to the rate limit. All CLAUDE-3.5-SONNET generations were done utilizing vertex-ai.

first retrieve the evidence (top 5, 25, and 50 paragraphs) and then prompt GPT-4o for evaluation.

**Inference:** We generate the labels by prompting the models with the entire book and one claim at time for verification. All closed-source models were prompted using the provider’s API, while open-weight models were run on up to three A100 80GB GPUs, excluding COMMAND R and COMMAND R+, which were accessed using the cohere API. Additionally, we implement a BM25+GPT-4o retrieval pipeline where we first retrieve the evidence from the book and then verify the claim based on the retrieved evidence with GPT-4o.<sup>37</sup>

**Prompts:** All models were prompted with the prompt presented in Table 17. As we notice that the open-source models did not follow the instructions well, we also tested them with a simplified prompt as shown in Table 18. Finally, for BM25+GPT-4o we employed the prompt in Table 19. All generations are restricted to 800 tokens.<sup>38</sup>

During a pilot study, we first tested three different types of prompts on a subset of data (10 books, 176 claim pairs):

- **Answer-only:** prompting the model for the answer only (Table 20)
- **Answer-then-explanation:** prompting the model for the answer followed by an explanation (Table 21)
- **Explanation-then-answer:** prompting the model for the explanation followed by the answer (Table 17)<sup>39</sup>

<sup>37</sup>The estimated cost for running all data with each model is as follows: GPT-4o \$640 USD, GPT-4-TURBO \$1,280 USD, CLAUDE-3-OPUS \$3,350 USD, CLAUDE-3.5-SONNET \$670 USD, GEMINI PRO 1.5 1,385 USD, COMMAND R+ \$450 USD, COMMAND R \$74 USD, BM25+GPT-4o pipeline \$330 USD. Note that test runs resulted in small additional costs.

<sup>38</sup>For some open-source models we noticed that the generation are plagued by repetitions and reduced this limit to 600 tokens.

<sup>39</sup>We gradually refined this prompt after the first pilot run on a few examples showed that models sometimes generate unreasonably long explanations, leaving no space for the actual answer. To mitigate this issue, we added a request to provide the explanation "in at most one paragraph." Although we still observed some cases where the model’s explanation was too long to generate the <answer></answer> tags, this occurred in only twice for CLAUDE-3-OPUS and once for GPT-4-TURBO. Importantly, the answer/label was always present in the explanation itself (early answer), so it was correctly extracted for that model.

MODEL	CONTEXT	AVAIL.	CHECKPOINTS	# PARAM
GPT-4O	128k	🔒	gpt-4o-2024-05-13	~(?)
GPT-4-TURBO	128k	🔒	gpt-4-turbo-2024-04-09	~(?)
CLAUDE-3-OPUS	200k	🔒	claude-3-opus-20240229	~(?)
CLAUDE-3.5-SONNET	200k	🔒	claude-3-5-sonnet@20240620	~(?)
GEMINI PRO 1.5	1M	🔒	gemini-1.5-pro-preview-0514	~(?)
GEMINI FLASH 1.5	1M	🔒	gemini-1.5-flash-preview-0514	~(?)
COMMAND R	128k	🔒	c4ai-command-r-v01	35B
COMMAND R+	128k	🔒	c4ai-command-r-plus	104B
GEMMA-10M	10M	🔒	gemma-2b-10m	2B
PHI-3-MINI	128k	🔒	Phi-3-mini-128k-instruct	3.8B
LONGLLAMA	256k	🔒	long_llama_3b_instruct	3B
MIXTRAL-8X22B	65k	🔒	Mixtral-8x22B-Instruct-v0.1	141B
QWEN-2-72B	32k <sup>36</sup>	🔒	Qwen2-72B-Instruct	72B

Table 16: Evaluated models: the upper part displays all closed-source models, while the lower part lists all open-weight models. We also provide details of two shorter context models, which were tested on story-length inputs for comparison.

As testing these prompts on all models would be prohibitively expensive, we conducted these experiments with GPT-4O and CLAUDE-3-OPUS.<sup>40</sup> We observed that CLAUDE-3-OPUS’s accuracy remained constant across all three setups (39.77%), although the specific pairs it got correct varied. GPT-4O achieved the highest accuracy when prompted for the explanation followed by the answer (50%), compared to 48.3% for the answer followed by the explanation, and 49.43% for the answer-only setup. For further experiments, we employ the *explanation-then-answer* approach, as it yielded the best results despite the small differences between these three methods.

**Label extraction:** As we prompt the models to generate answers following a structured template, we first attempt to extract the answer from the first encountered `<answer></answer>` tags.<sup>41</sup> If this is not possible (e.g., the model did not generate the tags or produced output longer than a simple true/false within the tags), we apply the following steps:

1. Replace any occurrence of “true or false” and the text of the claim itself with an empty string.<sup>42</sup>

<sup>40</sup>The total cost of this experiment, including iterative refinements to the prompt text, was \$3k USD.

<sup>40</sup>This version can also be extended to 128k using YARN (Peng et al., 2023).

<sup>41</sup>In cases where the text between the tags is longer than simply true/false, we follow the replacement pipeline and finally extract the label from the text between the tags.

<sup>42</sup>“True or false” often occurs when the model repeats the

Evaluation Template (Main)
<p>You are provided with a context and a statement. Your task is to carefully read the context and then determine whether the statement is true or false.</p> <p>Answer TRUE if the statement is true in its entirety based on the context provided. Answer FALSE if any part of the statement is false based on the context provided.</p> <p>&lt;context&gt;book text&lt;/context&gt; &lt;statement&gt;claim&lt;/statement&gt;</p> <p>&lt;question&gt;Based on the context provided, is the above statement TRUE or FALSE?&lt;/question&gt;</p> <p>First provide an explanation of your decision-making process in at most one paragraph, and then provide your final answer. Use the following format: &lt;explanation&gt;YOUR EXPLANATION&lt;/explanation&gt; &lt;answer&gt;YOUR ANSWER&lt;/answer&gt;</p>

Table 17: Prompt template used for Evaluation. This prompt was employed to evaluate *all* models. All open source models were also evaluated using prompt in Table 18.

2. Replace “not true” with “false.”
3. Extract “true” if only “true” is present, “false” if only “false” is present, and the first occurrence if both are present in the generation.

If neither “true” nor “false” is present, we count it as an automatic failure.

original question. We additionally replace the claim itself, as it can also contain words like true/false (e.g., “true friend”).

Evaluation Template (Simplified)
You are provided with a context and a statement. Your task is to carefully read the context and then determine whether the statement is true or false.
Answer TRUE if the statement is true in its entirety based on the context provided.
Answer FALSE if any part of the statement is false based on the context provided.
<context>book text</context>
<statement>claim</statement>
<question>Based on the context provided, is the above statement TRUE or FALSE?</question>

Table 18: Simplified prompt template used for Evaluation of open-source models only.

**Early response:** We employ the *explanation-then-answer* template to allow the model to reason about the answer before providing its final choice. However, we observe that models sometimes still return the answer first, within the <explanation></explanation> tags, which is in line with the observations in Levy et al. (2024).<sup>43</sup> Overall, CLAUDE-3-OPUS and CLAUDE-3.5-SONNET are the most affected, with 68.4% and 55.7% of explanations containing either "true" or "false" in the first sentence, respectively (see Table 22).

**Generation issues:** We encounter several issues during model prompting. Most notably, GEMINI PRO 1.5 and GEMINI FLASH 1.5 return a "prohibited content" API error for about 48.6% of claims, significantly reducing the number of claims (and claim pairs) processed by these models despite their large claimed context window. We also observe that open-source models struggle to follow the assigned template. Only COMMAND R returns the answer with proper <answer></answer> tags 100% of the time. This compares to COMMAND R+ at 66.3%, PHI-3-MINI at 62.7%, and GEMMA-10M at 52.5%. Additionally, we notice that generations by GEMMA-10M are plagued with code output, often not including any instances of "true" or "false," while LONGLLAMA returns empty outputs for the prompt in Table 17 (i.e., the non-simplified version of the prompt).

<sup>43</sup>We identify instances of an early response by segmenting the text between the <explanation></explanation> tags into sentences. In the first sentence, we remove occurrences of "true or false," "to determine if the statement is true," and similar phrases. We then check if the first sentence contains "true" or "false." Note that this excludes phrases like "this statement is incorrect." Additionally, we acknowledge that phrases such as "true nature," which may appear in the first sentence, may slightly affect this percentage as they are irrelevant to the model's answer.

Evaluation Template (BM25)
You are provided with excerpts of context and a statement. Your task is to carefully read the excerpts and then determine whether the statement is true or false.
Answer TRUE if the statement is true in its entirety based on the excerpts provided.
Answer FALSE if any part of the statement is false based on the excerpts provided.
<excerpt_1>excerpt_1</excerpt_1>
<excerpt_2>excerpt_2</excerpt_2>
<excerpt_(i)>...</excerpt_(i)>
<excerpt_k>excerpt_k</excerpt_k>
<statement>claim</statement>
<question>Based on the excerpts provided, is the above statement TRUE or FALSE?</question>
First provide an explanation of your decision-making process in at most one paragraph, and then provide your final answer. Use the following format:
<explanation>YOUR EXPLANATION</explanation>
<answer>YOUR ANSWER</answer>

Table 19: Prompt used for the BM25+GPT-4o Evaluation pipeline.

## E Results

In this section of the appendix we provide more results of our analysis. We also provide more examples of claims where the models have generated incorrect label or incorrect justification in Table 23.

Table 24 reports the accuracy of each model excluding the classic novels. Overall, GPT-4o still performs the best from all the models, although its accuracy drops slightly from 55.8% to 55.1% when classic novels are excluded.

**Are newer books harder to verify?** Our data includes (1) classic books, which were almost certainly in the models' training data, (2) books published in 2023, which might have been included in the training data, and (3) books published in 2024, which were likely not included in the training data.<sup>44</sup> Hence, we look at models' performance by the publication year (see Table 25). While we do not observe a large difference in performance for books published in 2023 vs. 2024, GPT-4o and GEMINI FLASH 1.5's accuracy drops slightly for the newer books, from 56.7% to 53.5% and from 37.4% to 30.7%, respectively. On the other hand, the performance of CLAUDE-3-OPUS, CLAUDE-3.5-SONNET, GEM-

<sup>44</sup>Given the time required to collect, clean the data, and train the model, we hypothesize that at least books published in May/June 2024 were not included in the models' training data.

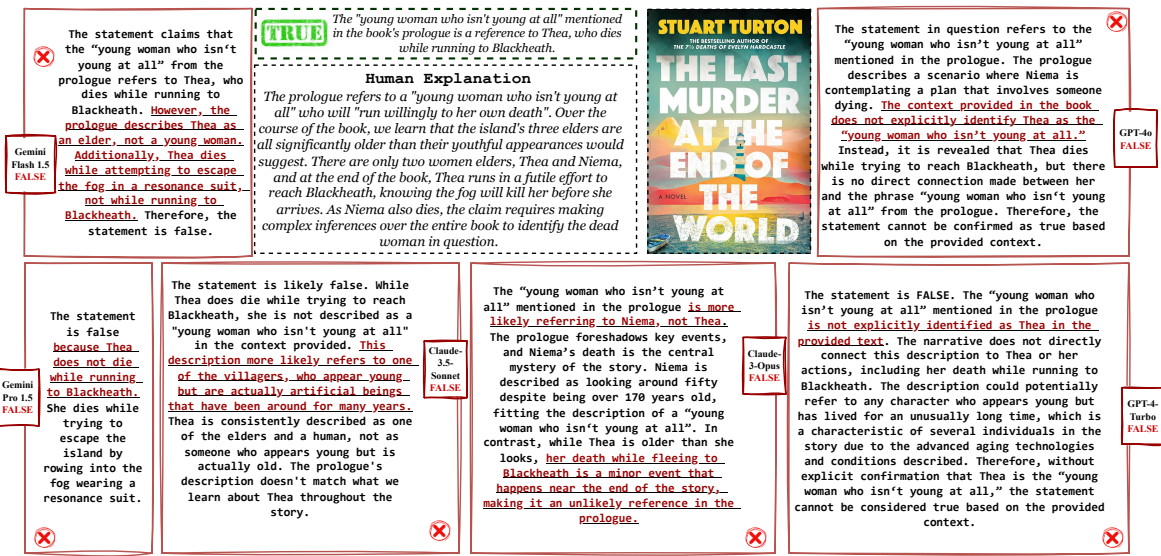


Figure 9: Example of a True claim for which all models predicted incorrect label. We provide the claim along with the annotator's explanation for clarity.

**Evaluation Template (Answer-only)**

---

You are provided with a context and a statement. Your task is to carefully read the context and then determine whether the statement is true or false.

Answer TRUE if the statement is true in its entirety based on the context provided.  
 Answer FALSE if any part of the statement is false based on the context provided.

<context>book text</context>

<statement>claim</statement>

<question>Based on the context provided, is the above statement TRUE or FALSE?</question>

First provide an explanation of your decision-making process in at most one paragraph, and then provide your final answer. Use the following format:

<answer>YOUR ANSWER</answer>

---

Table 20: Prompt used for the initial pilot study when the model is prompted to return only the answer.

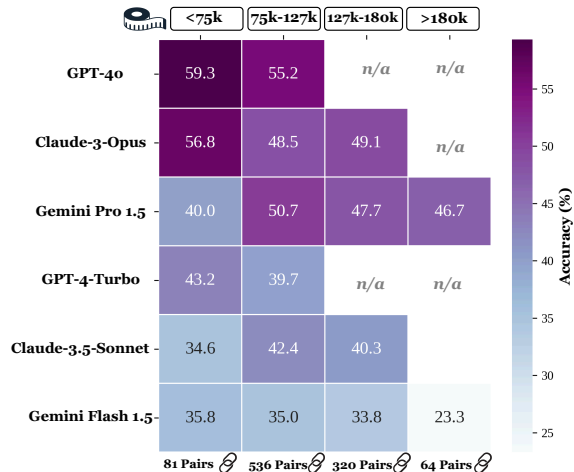


Figure 10: Model performance across different book lengths. Token counts are provided as per tiktoken. Length buckets were determined by taking one standard deviation from the mean on both sides, with additional buckets for values above and below this range. The number of valid pairs in each bucket is provided below.

INI PRO 1.5, and GPT-4-TURBO seems to be higher for books published in 2024 than in 2023. Overall, these differences may be due to the inherent difficulty of the claims for a given model rather than the publication year itself. We do observe that both GPT-4o and GPT-4-TURBO perform much better on "The Great Gatsby" (the only classic book, which fits their context window) at 73.3% and 66.7%, respectively. However, no definitive claims can be made as the set of claims is too small and related to one book only.

Evaluation Template (Answer and Explanation)	
You are provided with a context and a statement. Your task is to carefully read the context and then determine whether the statement is true or false.	
Answer TRUE if the statement is true in its entirety based on the context provided.	
Answer FALSE if any part of the statement is false based on the context provided.	
<context>book text</context>	
<statement>claim</statement>	
<question>Based on the context provided, is the above statement TRUE or FALSE?</question>	
First provide an explanation of your decision-making process in at most one paragraph, and then provide your final answer. Use the following format:	
<answer>YOUR ANSWER</answer>	
<explanation>YOUR EXPLANATION</explanation>	

Table 21: Prompt used for the initial pilot study when the model is prompted to return the answer and then the explanation.

MODEL	PERCENTAGE	COUNT (PRESENT/TOTAL)
GPT-4o	1.1%	13/1234
GPT-4-TURBO	53.2%	657/1234
CLAUDE-3-OPUS	68.4%	1281/1874
CLAUDE-3.5-SONNET	55.7%	1044/1874
GEMINI PRO 1.5	44.0%	453/1029
GEMINI FLASH 1.5	24.3%	250/1030

Table 22: Percentage and count of explanations where the first sentence contains an early response by the model. Open-source models are excluded as they typically do not generate structured responses or explanations.

**It is unclear if claims about longer books are harder to verify:** We categorized NoCHA books into four buckets based on length: (1) up to 75k tokens, (2) 75k to 127k tokens, (3) 127k to 180k tokens, and (4) above 180k tokens.<sup>45</sup> We observed a slight drop in performance between the first and second buckets for CLAUDE-3-OPUS (56.8% to 48.5%), GPT-4o (59.3% to 55.2%), and GPT-4-TURBO (43.2% to 39.7%). However, these differences are small and could be influenced by factors such as the number of claims in each bucket or the complexity of the book itself (i.e., longer books tend to have more complex narratives). No such difference was observed for GEMINI FLASH 1.5 and CLAUDE-3.5-SONNET, which performance is similar or slightly better, albeit still low, for longer books. GEMINI PRO 1.5 performed worst for the shortest

<sup>45</sup>We selected these numbers based on the mean length of our books (127k tokens), adjusted by adding or subtracting one standard deviation. Additional buckets were included beyond these limits on both ends.

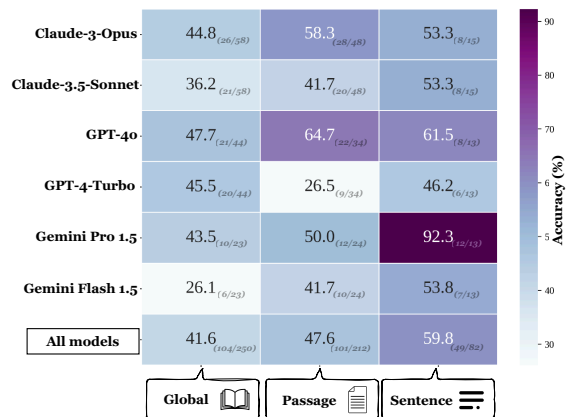


Figure 11: Performance of different closed-source models based on the scope of evidence. The count of correctly identified pairs and the total count of pairs are provided in brackets for reference.

books (40.0%), with performance for longer books varying from 46.7% to 50.7%. See Table 26 and Figure 10 for the details. We also report Pearson correlation between the book length and model’s accuracy for each model (see Table 27).

**Impact of evidence scope:** We report the performance of all models on the subset of data annotated for evidence scope in Table 28. Note that the number of examples varies between models due to the restrictions of the models’ context windows and GEMINI PRO 1.5/GEMINI FLASH 1.5 API errors. Additionally, the performance of closed-source models is visualized in Figure 11.

**Impact of irrelevant context:** As mentioned previously, creating *complex* claims with evidence at different book depths is challenging, especially when aiming to create global claims that require the model to reason over a longer context. Instead, we utilize short stories to investigate whether claims made about the same story are harder to identify when the model is prompted with the entire collection versus just the story itself. The accuracies of all models for the whole collection versus individual stories are presented in Table 29. The story-only setup is indicated with a subscript "story." To see these results on common set of claim pairs, that is claim pairs fully processed by all models, see Table 30. Additionally, we report the performance by story location, either the beginning, middle or the end of the collection, relative to the performance when prompted with that story only (see



CLAIM & GOLD LABEL	MODEL'S LABEL & JUSTIFICATION	HUMAN EXPLANATION
<p>Wyatt apologizes to Ernest because he believes McTavish left him a negative review.</p> <p>Everyone on the Train is a Suspect by Benjamin Stevenson</p> <p style="text-align: center;"><b>-True</b></p>	<p>Wyatt's apology to Ernest is mentioned in the context, but it is not because he believes McTavish left him a negative review. <b>Instead, Wyatt's apology is vague and not directly linked to the review.</b> Ernest himself is confused about the reason for the apology, and it is later revealed that Simone, Ernest's agent, was behind the negative review as part of a coded message to McTavish. Therefore, the statement that Wyatt apologizes to Ernest because he believes McTavish left him a negative review is incorrect.</p> <p>-GPT-4o: <b>False</b></p>	<p>When Wyatt apologizes, Ernest does not initially understand the reason behind it. However, it is heavily implied shortly after that Wyatt, McTavish's publisher, was apologizing for the bad review, which at that time everyone, including Wyatt thought was written by McTavish himself.</p>
<p>Emily Sparrow was pregnant with John's child and disappeared in 1966.</p> <p>How to Solve Your Own Murder by Kristen Perrin</p> <p style="text-align: center;"><b>-True</b></p>	<p>The text states that Emily Sparrow was pregnant, but it is never explicitly stated who the father of the child is. While it is implied that the father could be John, it is also implied that it could be Rutherford Gravesdown (Ford).</p> <p>-GEMINI PRO 1.5: <b>False</b></p>	<p>Throughout the book, the reader is taken back to 1966, where it is revealed that Emily was pregnant, and the father's identity was uncertain—either Lord Gravesdown or John. The baby, Laura, was later adopted and had a daughter, Annie. In the latter part of the book, set in the present day, Annie meets John, who confesses that he is her grandfather, thus confirming that he was the father of Emily's baby.</p>
<p>None of the writers on the train went to the Gorge excursion as they had to stay to participate in a festival event.</p> <p>Everyone on the Train is a Suspect by Benjamin Stevenson</p> <p style="text-align: center;"><b>-False</b></p>	<p>The statement is true based on the context. The passage states that while the non-festival tourists went on a scheduled day trip sailing down Katherine Gorge, "In front of our carriages were forty or so black fold-out chairs set up in the red dirt for the festival attendees." <b>This indicates that the writers stayed behind for a festival event while the other passengers went on the Gorge excursion.</b></p> <p>-CLAUDE-3-OPUS: <b>True</b></p>	<p>While all writers invited to the festival stayed for the event, Juliet, who was also an accomplished writer and, in fact, we later learn she was originally invited as well, went on the trip.</p>
<p>Lost Bride Manor is haunted by various spirits, including the benevolent ghost of Clover, Sonya's grandmother and Collin's mother, who communicates through music.</p> <p>Inheritance by Nora Roberts</p> <p style="text-align: center;"><b>-True</b></p>	<p>The statement is mostly true, but contains one key inaccuracy. The context does indicate that Lost Bride Manor is haunted by various spirits, including Clover, who communicates through music. Clover is described as Sonya's grandmother and Collin's mother, which is correct. However, the statement characterizes Clover as "benevolent", but the context does not provide enough information to definitively label her as such. <b>While Clover seems to be helpful at times, playing music and aiding Sonya, the text does not explicitly state that she is entirely benevolent.</b> Given this small but significant discrepancy, the statement cannot be considered entirely true.</p> <p>-CLAUDE-3.5-SONNET: <b>False</b></p>	<p>Clover is a ghost haunting the manor where Sonya moves in. From the first day, Clover, who we learn is the ghost of Sonya's grandmother, helps Sonya by playing music to cheer her up, adding to her courage, and keeping her company. Clover also tries to protect Sonya from the evil spirit that also haunts the manor. Although it is not directly stated in the text, Clover's kind and well-meaning nature towards Sonya is evident through her actions.</p>
<p>Ernest, the narrator, drops numerous hints about the killer's identity from the start, but any hints of how the first victim dies aren't revealed until chapter 11.</p> <p>Everyone on the Train is a Suspect by Benjamin Stevenson</p> <p style="text-align: center;"><b>-False</b></p>	<p>To determine the accuracy of the statement, I reviewed the context for any early hints about the killer's identity and the method of the first victim's death. The prologue and early chapters contain several clues about the characters' motives and potential involvement in the murder. Specifically, <b>the prologue includes a detailed description of a potential murder scene, mentioning "the bloodied handprint on the bathroom door" and "the crimson, dripping tip of the murder weapon."</b> This indicates that hints about the method of the first victim's death are indeed revealed before chapter 11.</p> <p>-GPT-4o: <b>False</b></p>	<p>It is true that Ernest, the narrator, drops some hints about how the first victim dies from the start. However, the passage cited by the model is irrelevant. It is taken from the Prologue, where Ernest, a mystery novel writer, talks about elements that typically appear in a mystery novel's prologue. He mentions a "bloodied handprint" and "the crimson, dripping tip of the murder weapon," which are common in mystery novels prologues. However, in this novel, the first victim was poisoned, and Ernest hints at this by saying that the inspiration for this book "came from a drink" with the victim.</p>
<p>Levi promises Grace that he will kill anyone who tries to hurt her, and Grace calls him to help her dispose of Norman's body.</p> <p>Safe and Sound by Laura McHugh</p> <p style="text-align: center;"><b>-False</b></p>	<p>In the context, Levi does promise Grace that he would kill anyone who tried to hurt her again. However, when Grace kills Norman... she calls Tyson, Levi's brother. The passage states: "She... wanted to make it look like she couldn't have survived. So she could just go missing..." <b>It then says "Tyson grabbed me, clamped his hands around my neck, and didn't let go. 'I dumped a body for you.'" This shows that it was Tyson, not Levi, who Grace called to help dispose of Norman's body.</b></p> <p>-GPT-4o: <b>False</b></p>	<p>The model's justification begins correctly, but the first passage cited is unrelated to the claim and clearly a lie told by Tyson in the context of the story. Additionally, the text "Tyson grabbed me, clamped his hands around my neck, and didn't let go." never appears in the book. Tyson does say "I dumped a body for you." to Grace—which proves the claim is false—but he does not choke her before saying so or at any other point in the book.</p>
<p>While on Sigrid's boat, Seventeen admits to Kat that he had the same underwater dream that she had.</p> <p>Assassin Eighteen by John Brownlow</p> <p style="text-align: center;"><b>-True</b></p>	<p>The statement is true based on the context provided. In Part 8, <b>after Kat is rescued from the water by Rachel and Vilmos in the Zodiac, she describes her near-death experience underwater where she felt she and Seventeen were entwined as one.</b> Seventeen then admits "I had the same dream", confirming he experienced the same vision while he was drowning.</p> <p>-CLAUDE-3-OPUS: <b>True</b></p>	<p>It is true that Seventeen admits he had the same dream in Part 8 of the book. However, Kat was not rescued by Rachel and Vilmos in the Zodiac during this scene. That happens much later in the book, towards the end of Part 11. This conversation happened when Kat and Seventeen were rescued by Sigrid.</p>
<p>Sonya befriends Trey's sister Anna, a potter, and his cousin Owen Poole, who becomes involved in solving the mystery of the lost brides.</p> <p>Inheritance by Nora Roberts</p> <p style="text-align: center;"><b>-False</b></p>	<p>The statement is mostly true, but contains one key inaccuracy. The context shows that Sonya does befriend Trey's sister Anna, who is a potter, and helps redesign her website. Sonya also meets and interacts with Owen Poole, <b>Trey's cousin, who becomes involved in investigating the haunted manor and the mystery of the lost brides.</b> However, the statement is incorrect in saying that Owen "becomes involved in solving the mystery of the lost brides." While Owen does help investigate strange occurrences at the manor and even enters the mirror with Sonya at the end, he is not described as actively trying to solve the mystery of the lost brides specifically. Sonya is the main character focused on that mystery. Therefore, while most of the statement is accurate, this one detail makes it not entirely true.</p> <p>-CLAUDE-3.5-SONNET: <b>False</b></p>	<p>It is true that Sonya befriends Trey's sister, Anna, and Owen. However, Owen is not Trey's cousin; he is Sonya's cousin, whom she did not know about. Owen and Trey then try to help Sonya solve the mystery of the lost brides.</p>

Table 23: **Top half:** Examples of claims where the model generated incorrect labels, along with the model-generated label and justification, and human-written explanations for the correct label of each claim. **Bottom half:** Examples of claims correctly labeled by the model but with wrong justifications. The explanations were provided by the authors, based on examples from books they read. Tags have been removed from the generated content for clarity.

MODEL	PAIR ACC <sub>(correct/total)</sub>	COMMON SET ACC
GPT-4O	<b>55.3</b> <sub>(333/602)</sub>	<b>57.5</b> <sub>(195/339)</sub>
GPT-4-TURBO	39.5 <sub>(238/602)</sub>	38.9 <sub>(132/339)</sub>
CLAUDE-3-OPUS	49.4 <sub>(439/889)</sub>	49.9 <sub>(169/339)</sub>
CLAUDE-3.5-SONNET	42.0 <sub>(373/889)</sub>	41.3 <sub>(140/339)</sub>
GEMINI PRO 1.5	48.1 <sub>(224/466)</sub>	48.4 <sub>(164/339)</sub>
GEMINI FLASH 1.5	34.5 <sub>(161/467)</sub>	35.1 <sub>(119/339)</sub>
COMMAND R	18.8 <sub>(81/430)</sub>	<i>n/a</i>
COMMAND R <sub>simple</sub>	22.1 <sub>(95/430)</sub>	<i>n/a</i>
COMMAND R+	17.4 <sub>(75/430)</sub>	<i>n/a</i>
COMMAND R+ <sub>simple</sub>	13.3 <sub>(57/430)</sub>	<i>n/a</i>
PHI-3-MINI	9.7 <sub>(23/237)</sub>	<i>n/a</i>
PHI-3-MINI <sub>simple</sub>	14.2 <sub>(45/316)</sub>	<i>n/a</i>
GEMMA-10M	4.2 <sub>(39/938)</sub>	<i>n/a</i>
GEMMA-10M <sub>simple</sub>	7.7 <sub>(72/938)</sub>	<i>n/a</i>
LONGLLAMA <sub>simple</sub>	5.1 <sub>(45/889)</sub>	<i>n/a</i>
BM25+GPT-4o ( <i>k</i> =5)	27.8 <sub>(261/938)</sub>	28.1 <sub>(95/338)</sub>
BM25+GPT-4o ( <i>k</i> =25)	44.0 <sub>(413/938)</sub>	45.9 <sub>(155/338)</sub>
BM25+GPT-4o ( <i>k</i> =50)	49.5 <sub>(464/938)</sub>	50.3 <sub>(170/338)</sub>
RANDOM	25.0 <sub>(234/938)</sub>	25.0 <sub>(85/339)</sub>

Table 24: Model accuracy on claim pairs for all data excluding *classic* novels (see Table 17 for the prompt; and Table 19 for the prompt employed with BM25). “COMMON SET” refers to claim pairs shared among the models. The subscript “SIMPLE” refers to the calls done with simplified prompt (Table 18).

Table 31).<sup>46</sup>

**CLAUDE-3-OPUS cites evidence more often than other models:** We employ a simple heuristic by looking for quotation marks in the evidence to identify how often models cite excerpts from the source. Qualitatively, we observed that CLAUDE-3-OPUS cites the source more often than other closed-source models, and the citations we verified were always present in the source. This quantitative analysis confirms our observation, as 49.6% of CLAUDE-3-OPUS’s responses contain quotation marks, followed by 18.6% for both GPT-4O and GPT-4-TURBO, 11.5% for GEMINI PRO 1.5, and 10.8% for GEMINI FLASH 1.5. We also examined whether responses containing citations are more likely to be correct but did not find such a relation (see Table 32).

**Effect of genre:** We classify the books into three genres: historical, contemporary, and speculative. Table 33 presents results by genre for all models.

<sup>46</sup>Note that we do not provide an in-depth analysis and discussion of accuracy relative to the story’s location due to the small number of annotations at different depths.

**Length of the justifications:** Figure 12 provides lengths of the justification provided by each model for the main prompt template (Table 17). We report the lengths in words by accuracy for each model. Note that GEMMA-10M never produces real explanation, rather sometimes repeats part of the original prompt, i.e., <explanation>YOUR ANSWER</explanation>.

**Length of the retrieved chunks:** Figure 13 shows the average percentage of the book retrieved by BM25 for varying values of *k*, i.e., for top 5, 25, and 50 chunks.

**Performance by book:** We report the performance of each model by the book title in Figure 14. Empty cells indicate cases where the model did not process the book GEMINI PRO 1.5 and GEMINI FLASH 1.5) or could not process the book because of its context window. See Table 11 for the number of claim pairs written for each book.

**Statistical analysis:** We conducted statistical analysis by fitting generalized linear mixed-effects models using the `glmer()` function in R (Bates

MODEL	CLASSICS <sub>(correct/total)</sub>	2023 <sub>(correct/total)</sub>	2024 <sub>(correct/total)</sub>
GPT-4O	73.3 <sub>(11/15)</sub>	56.7 <sub>(194/342)</sub>	53.5 <sub>(139/260)</sub>
GPT-4-TURBO	66.7 <sub>(10/15)</sub>	39.5 <sub>(135/342)</sub>	39.6 <sub>(103/260)</sub>
CLAUDE-3-OPUS	50.0 <sub>(24/48)</sub>	48.8 <sub>(227/465)</sub>	50.0 <sub>(212/424)</sub>
CLAUDE-3.5-SONNET	22.9 <sub>(11/48)</sub>	41.9 <sub>(195/465)</sub>	42.0 <sub>(178/424)</sub>
GEMINI PRO 1.5	47.9 <sub>(23/48)</sub>	45.7 <sub>(121/265)</sub>	51.2 <sub>(103/201)</sub>
GEMINI FLASH 1.5	31.2 <sub>(15/48)</sub>	37.4 <sub>(99/265)</sub>	30.7 <sub>(62/202)</sub>
COMMAND R	40.0 <sub>(6/15)</sub>	18.5 <sub>(51/275)</sub>	19.4 <sub>(30/155)</sub>
COMMAND R <sub>simple</sub>	33.3 <sub>(5/15)</sub>	21.5 <sub>(59/275)</sub>	23.2 <sub>(36/155)</sub>
COMMAND R+	13.3 <sub>(2/15)</sub>	18.5 <sub>(51/275)</sub>	15.5 <sub>(24/155)</sub>
COMMAND R+ <sub>simple</sub>	26.7 <sub>(4/15)</sub>	14.5 <sub>(40/275)</sub>	11.0 <sub>(17/155)</sub>
PHI-3-MINI	10.0 <sub>(1/10)</sub>	8.4 <sub>(12/143)</sub>	10.6 <sub>(10/94)</sub>
PHI-3-MINI <sub>simple</sub>	20.0 <sub>(3/15)</sub>	13.1 <sub>(25/191)</sub>	16.0 <sub>(20/125)</sub>
GEMMA-10M	0.0 <sub>(0/63)</sub>	3.7 <sub>(18/489)</sub>	4.7 <sub>(21/449)</sub>
GEMMA-10M <sub>simple</sub>	4.8 <sub>(3/63)</sub>	9.8 <sub>(48/489)</sub>	5.3 <sub>(24/449)</sub>
LONGLLAMA <sub>simple</sub>	2.1 <sub>(1/48)</sub>	6.5 <sub>(30/465)</sub>	3.5 <sub>(15/424)</sub>
BM25+GPT-4O ( $k=5$ )	33.3 <sub>(21/63)</sub>	28.8 <sub>(141/489)</sub>	26.7 <sub>(120/449)</sub>
BM25+GPT-4O ( $k=25$ )	44.4 <sub>(28/63)</sub>	44.6 <sub>(218/489)</sub>	43.4 <sub>(195/449)</sub>
BM25+GPT-4O ( $k=50$ )	52.4 <sub>(33/63)</sub>	47.2 <sub>(231/489)</sub>	51.9 <sub>(233/449)</sub>

Table 25: Models accuracy on claim pairs by the publication year.

MODEL	BELOW 75K	75K-127K	127K-180K	ABOVE 180K
GPT-4O	59.3 <sub>(48/81)</sub>	55.2 <sub>(296/536)</sub>	<i>n/a</i>	<i>n/a</i>
GPT-4-TURBO	43.2 <sub>(35/81)</sub>	39.7 <sub>(213/536)</sub>	<i>n/a</i>	<i>n/a</i>
CLAUDE-3-OPUS	56.8 <sub>(46/81)</sub>	48.5 <sub>(260/536)</sub>	49.1 <sub>(157/320)</sub>	<i>n/a</i>
CLAUDE-3.5-SONNET	34.6 <sub>(28/81)</sub>	42.4 <sub>(227/536)</sub>	40.3 <sub>(129/320)</sub>	<i>n/a</i>
GEMINI PRO 1.5	40.0 <sub>(32/80)</sub>	50.7 <sub>(139/274)</sub>	47.7 <sub>(62/130)</sub>	46.7 <sub>(14/30)</sub>
GEMINI FLASH 1.5	35.8 <sub>(29/81)</sub>	35.0 <sub>(96/274)</sub>	33.8 <sub>(44/130)</sub>	23.3 <sub>(7/30)</sub>
COMMAND R	25.9 <sub>(21/81)</sub>	18.1 <sub>(66/364)</sub>	<i>n/a</i>	<i>n/a</i>
COMMAND R <sub>simple</sub>	25.9 <sub>(21/81)</sub>	21.7 <sub>(79/364)</sub>	<i>n/a</i>	<i>n/a</i>
COMMAND R+	27.2 <sub>(22/81)</sub>	15.1 <sub>(55/364)</sub>	<i>n/a</i>	<i>n/a</i>
COMMAND R+ <sub>simple</sub>	27.2 <sub>(22/81)</sub>	10.7 <sub>(39/364)</sub>	<i>n/a</i>	<i>n/a</i>
LONGLLAMA <sub>simple</sub>	9.9 <sub>(8/81)</sub>	5.4 <sub>(29/536)</sub>	2.8 <sub>(9/320)</sub>	<i>n/a</i>
PHI-3-MINI	13.9 <sub>(5/36)</sub>	8.5 <sub>(18/211)</sub>	<i>n/a</i>	<i>n/a</i>
PHI-3-MINI <sub>simple</sub>	13.6 <sub>(11/81)</sub>	14.8 <sub>(37/250)</sub>	<i>n/a</i>	<i>n/a</i>
GEMMA-10M	4.9 <sub>(4/81)</sub>	4.1 <sub>(22/536)</sub>	3.4 <sub>(11/320)</sub>	3.1 <sub>(2/64)</sub>
GEMMA-10M <sub>simple</sub>	9.9 <sub>(8/81)</sub>	7.6 <sub>(41/536)</sub>	6.9 <sub>(22/320)</sub>	6.2 <sub>(4/64)</sub>
BM25+GPT-4O ( $k=5$ )	29.6 <sub>(24/81)</sub>	28.0 <sub>(150/536)</sub>	26.6 <sub>(85/320)</sub>	35.9 <sub>(23/64)</sub>
BM25+GPT-4O ( $k=25$ )	49.4 <sub>(40/81)</sub>	42.7 <sub>(229/536)</sub>	41.6 <sub>(133/320)</sub>	60.9 <sub>(39/64)</sub>
BM25+GPT-4O ( $k=50$ )	51.9 <sub>(42/81)</sub>	47.9 <sub>(257/536)</sub>	48.1 <sub>(154/320)</sub>	68.8 <sub>(44/64)</sub>

Table 26: Models accuracy on claim pairs by the book length in tokens (tiktoken).

et al., 2015).<sup>47</sup> The response variable was pairwise **accuracy**, a binary categorical variable ("cor-

<sup>47</sup>All models were fitted using the closed-source models data only, as all open-weights models performed below random.

rect" or "incorrect"). Pair IDs and annotators were modeled as random effects, with various predictors (fixed effects) included in different models:

1. **Model** - a 6-level categorical variable (GPT-

MODEL	CORRELATION	<i>p</i> -VALUE
GPT-4O	-0.02	0.9150
GPT-4-TURBO	0.05	0.7722
CLAUDE-3-OPUS	-0.09	0.4682
CLAUDE-3.5-SONNET	0.03	0.8252
GEMINI PRO 1.5	0.03	0.8773
GEMINI FLASH 1.5	-0.13	0.4698
COMMAND R	-0.38	0.0400*
COMMAND R <sub>simple</sub>	-0.22	0.2329
COMMAND R+	-0.16	0.3865
COMMAND R+ <sub>simple</sub>	-0.53	0.0029**
PHI-3-MINI	-0.57	0.0091**
PHI-3-MINI <sub>simple</sub>	0.12	0.5972
GEMMA-10M	-0.11	0.3959
GEMMA-10M <sub>simple</sub>	-0.09	0.4858
LONGLLAMA <sub>simple</sub>	-0.14	0.2760
BM25+GPT-4O ( <i>k</i> =5)	0.11	0.3774
BM25+GPT-4O ( <i>k</i> =25)	0.21	0.0849
BM25+GPT-4O ( <i>k</i> =25)	0.27	0.0263*

Table 27: Pearson correlation between length and accuracy for different models. Significant correlations are marked with an asterisk.

4O, GPT-4-TURBO, CLAUDE-3-OPUS, CLAUDE-3.5-SONNET, GEMINI PRO 1.5, GEMINI FLASH 1.5), with analysis restricted to pairs processed by all models. GPT-4O was set as the reference level (intercept). See Table 34 and Table 35 for the results.;

2. **Length group** - a 4-level categorical variable ("below 75k", "75k-127k", "127k-180k", "above 180k"). "Below 75k" category was set as the reference level (intercept). See Table 36 and Table 37 for the results;
3. **Year** - a 3-level categorical variable ("classics", "2023", "2024") with "classics" set as the reference level. See Table 38 for the results;
4. **Genre** - a 3-level categorical variable ("historical", "contemporary", "speculative") with "speculative" set as the reference level. See Table 39 and Table 40 for the results;
5. **Scope** - a 3-level categorical variable ("sentence", "passage", "global") with "sentence" set as the reference level. See Table 41 and Table 42 for the results.

All models were fitted using the bobyqa optimizer with a binomial link function. We chose mixed-effects models for two main reasons: (1) to account for repeated measures, as each model validates multiple pairs and each pair is validated by multiple models, and (2) to flexibly model pair IDs

and annotators as random effects, partially controlling for the inherent difficulty of the pairs unrelated to the predictors.

For these mixed-effects models, we report two types of  $R^2$  values:

- *Marginal*  $R^2$ , which indicates the proportion of variance explained by the fixed effects (predictors) alone.
- *Conditional*  $R^2$ , which represents the proportion of variance explained by both the fixed and random effects (Nakagawa et al., 2017).

We further conducted a post-hoc analysis using the emmeans package (Lenth, 2023) in R with Tukey adjustments for multiple comparisons. To obtain probabilities, we first converted log-odds to odds ratios by exponentiating the estimates, and then converted the odds ratios to probabilities which are reported in the post-hoc tables.

MODEL	SENTENCE <sub>(correct/total)</sub>	PASSAGE <sub>(correct/total)</sub>	GLOBAL <sub>(correct/total)</sub>
GPT-4O	61.5 (8/13)	64.7 (22/34)	47.7 (21/44)
GPT-4-TURBO	46.2 (6/13)	26.5 (9/34)	45.5 (20/44)
CLAUDE-3-OPUS	53.3 (8/15)	58.3 (28/48)	44.8 (26/58)
CLAUDE-3.5-SONNET	53.3 (8/15)	41.7 (20/48)	36.2 (21/58)
GEMINI PRO 1.5	92.3 (12/13)	50.0 (12/24)	43.5 (10/23)
GEMINI FLASH 1.5	53.8 (7/13)	41.7 (10/24)	26.1 (6/23)
COMMAND R	25.0 (2/8)	23.1 (6/26)	3.7 (1/27)
COMMAND R <sub>simple</sub>	37.5 (3/8)	11.5 (3/26)	22.2 (6/27)
COMMAND R+	25.0 (2/8)	19.2 (5/26)	18.5 (5/27)
COMMAND R+ <sub>simple</sub>	12.5 (1/8)	7.7 (2/26)	11.1 (3/27)
PHI-3-MINI	0.0 (0/4)	12.5 (2/16)	0.0 (0/16)
PHI-3-MINI <sub>simple</sub>	0.0 (0/6)	15.8 (3/19)	0.0 (0/21)
GEMMA-10M <sub>simple</sub>	6.7 (1/15)	2.1 (1/48)	1.7 (1/58)
GEMMA-10M	0.0 (0/15)	8.3 (4/48)	1.7 (1/58)
LONGLLAMA <sub>simple</sub>	6.7 (1/15)	0.0 (0/48)	3.4 (2/58)
BM25+GPT-4o ( $k=5$ )	46.7 (7/15)	22.9 (11/48)	25.9 (15/58)
BM25+GPT-4o ( $k=25$ )	66.7 (10/15)	43.8 (21/48)	29.3 (17/58)
BM25+GPT-4o ( $k=50$ )	73.3 (11/15)	45.8 (22/48)	41.4 (24/58)

Table 28: Models' accuracy (%) by the scope of evidence on the annotated subset of data. Counts are provided as subscripts.

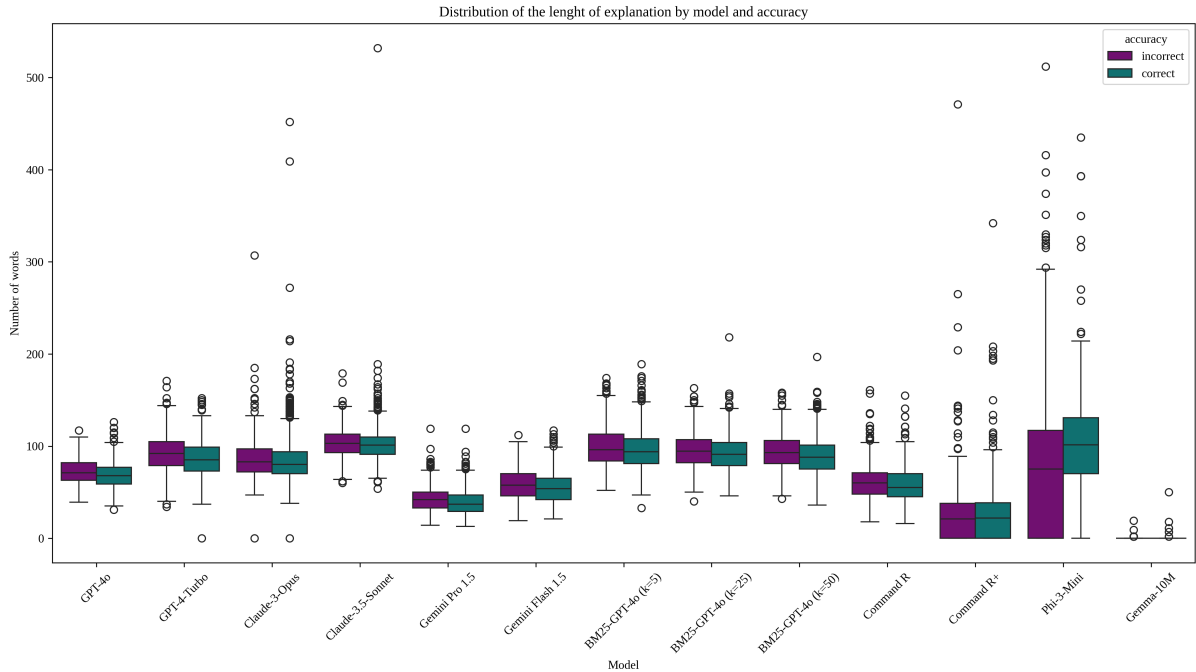


Figure 12: Boxplots of the length of justification (in words) provided for each model for correct (green) and incorrect (purple) predictions.

MODEL	ACC ( <b>PAIR</b> ) <sub>(correct/total)</sub>	ACC ( <b>TRUE</b> ) <sub>(correct/total)</sub>	ACC ( <b>FALSE</b> ) <sub>(correct/total)</sub>
GPT-4O	56.9 (29/51)	82.4 (42/51)	72.5 (37/51)
GPT-4-TURBO	39.2 (20/51)	58.8 (30/51)	76.5 (39/51)
CLAUDE-3-OPUS	37.7 (26/69)	75.4 (52/69)	62.3 (43/69)
CLAUDE-3.5-SONNET	39.1 (27/69)	52.2 (36/69)	84.1 (58/69)
GEMINI PRO 1.5	52.3 (23/44)	68.2 (30/44)	84.1 (37/44)
GEMINI FLASH 1.5	40.9 (18/44)	61.4 (27/44)	77.3 (34/44)
GPT-4O <sub>story</sub>	65.2 (45/69)	76.8 (53/69)	88.4 (61/69)
GPT-4-TURBO <sub>story</sub>	50.7 (35/69)	59.4 (41/69)	87.0 (60/69)
CLAUDE-3-OPUS <sub>story</sub>	63.8 (44/69)	79.7 (55/69)	81.2 (56/69)
CLAUDE-3.5-SONNET <sub>story</sub>	50.7 (35/69)	62.3 (43/69)	88.4 (61/69)
GEMINI PRO 1.5 <sub>story</sub>	51.5 (35/68)	58.8 (40/68)	77.9 (53/68)
GEMINI FLASH 1.5 <sub>story</sub>	42.0 (29/69)	50.7 (35/69)	75.4 (52/69)
COMMAND R+	11.8 (6/51)	68.6 (35/51)	31.4 (16/51)
COMMAND R+ <sub>simple</sub>	19.6 (10/51)	74.5 (38/51)	39.2 (20/51)
COMMAND R	23.5 (12/51)	80.4 (41/51)	39.2 (20/51)
COMMAND R <sub>simple</sub>	31.4 (16/51)	82.4 (42/51)	45.1 (23/51)
PHI-3-MINI	7.9 (3/38)	45.0 (18/40)	35.9 (14/39)
PHI-3-MINI <sub>simple</sub>	9.8 (5/51)	80.4 (41/51)	29.4 (15/51)
COMMAND R <sub>story</sub>	39.1 (27/69)	81.2 (56/69)	53.6 (37/69)
COMMAND R+ <sub>story</sub>	43.5 (30/69)	81.2 (56/69)	59.4 (41/69)
PHI-3-MINI <sub>story</sub>	15.4 (10/65)	26.9 (18/67)	83.3 (55/66)
PHI-3-MINI <sub>simple-story</sub>	24.2 (16/66)	62.1 (41/66)	56.7 (38/67)
MIXTRAL-8X22B <sub>story</sub>	56.6 (39/69)	75.4 (52/69)	81.2 (56/69)
QWEN-2-72B <sub>story</sub>	60.9 (42/69)	79.7 (55/69)	81.2 (56/69)
📚 BM25+GPT-4O ( $k=5$ )	37.7 (26/69)	44.9 (31/69)	91.3 (63/69)
📚 BM25+GPT-4O ( $k=25$ )	53.6 (37/69)	62.3 (43/69)	85.5 (59/69)
📚 BM25+GPT-4O ( $k=50$ )	56.5 (39/69)	71.0 (49/69)	82.6 (57/69)
RANDOM	25.0 (17/69)	50.0 (35/69)	50.0 (35/69)

Table 29: **Pairwise** accuracy, accuracy on **True**, and accuracy on **False** for claims made about collections of stories (80k-129k tokens) versus individual stories from the collection (700-21k tokens, average 8.5k). All texts fitting within the model’s context window were listed. For the 128k token models, the accuracy for the entire collection excludes one book with 129k tokens, which is included in the accuracy for individual stories, as none of the stories exceeded the context window. In the case of GEMINI PRO 1.5 and GEMINI FLASH 1.5, two collections were not processed due to being identified as copyrighted content ("prohibited content" API error). However, both models processed the stories from these collections, likely due to the length difference, though some claims were still refused due to disruptive content. We also provide accuracy on stories for MIXTRAL-8X22B and QWEN-2-72B, which have context window of 65k and 32k respectively, for comparison. Subscript "story" denotes results obtained with prompting the models with individual stories as the context.




MODEL	ACC	ACC	ACC
	PAIR <sub>(correct/total)</sub>	TRUE <sub>(correct/total)</sub>	FALSE <sub>(correct/total)</sub>
GPT-4o	48.0 <sub>(12/25)</sub>	76.0 <sub>(19/25)</sub>	68.0 <sub>17/25</sub>
GPT-4-TURBO	44.0 <sub>(11/25)</sub>	64.0 <sub>16/25</sub>	76.0 <sub>19/25</sub>
CLAUDE-3-OPUS	32.0 <sub>(8/25)</sub>	88.0 <sub>22/25</sub>	44.0 <sub>11/25</sub>
CLAUDE-3.5-SONNET	52.0 <sub>13/25</sub>	68.0 <sub>17/25</sub>	80.0 <sub>20/25</sub>
GEMINI PRO 1.5	56.0 <sub>14/25</sub>	80.0 <sub>20/25</sub>	76.0 <sub>19/25</sub>
GEMINI FLASH 1.5	48.0 <sub>12/25</sub>	76.0 <sub>19/25</sub>	68.0 <sub>17/25</sub>
GPT-4o <sub>story</sub>	72.0 <sub>18/25</sub>	84.0 <sub>21/25</sub>	88.0 <sub>22/25</sub>
GPT-4-TURBO <sub>story</sub>	56.0 <sub>14/25</sub>	68.0 <sub>17/25</sub>	84.0 <sub>21/25</sub>
CLAUDE-3-OPUS <sub>story</sub>	72.0 <sub>18/25</sub>	84.0 <sub>21/25</sub>	84.0 <sub>21/25</sub>
CLAUDE-3.5-SONNET <sub>story</sub>	60.0 <sub>15/25</sub>	68.0 <sub>17/25</sub>	92.0 <sub>23/25</sub>
GEMINI PRO 1.5 <sub>story</sub>	64.0 <sub>16/25</sub>	72.0 <sub>18/25</sub>	88.0 <sub>22/25</sub>
GEMINI FLASH 1.5 <sub>story</sub>	48.0 <sub>12/25</sub>	68.0 <sub>17/25</sub>	76.0 <sub>19/25</sub>
MIXTRAL-8X22B <sub>story</sub>	64.0 <sub>16/25</sub>	80.0 <sub>20/25</sub>	84.0 <sub>21/25</sub>
QWEN-2-72B <sub>story</sub>	64.0 <sub>16/25</sub>	80.0 <sub>20/25</sub>	84.0 <sub>21/25</sub>
 BM25+GPT-4o ( $k=5$ )	52.0 <sub>13/25</sub>	64.0 <sub>16/25</sub>	88.0 <sub>22/25</sub>
 BM25+GPT-4o ( $k=5$ )	68.0 <sub>17/25</sub>	76.0 <sub>19/25</sub>	92.0 <sub>23/25</sub>
 BM25+GPT-4o ( $k=5$ )	72.0 <sub>18/25</sub>	88.0 <sub>22/25</sub>	84.0 <sub>21/25</sub>

Table 30: **Pairwise** accuracy, accuracy on **True**, and accuracy on **False** for claims made about collections of stories (80k-129k tokens) versus individual stories from the collection (700-21k tokens, average 8.5k) on **common set** of claim pairs (i.e., pairs which were processed by all the models in all shown configurations). We also provide accuracies for MIXTRAL-8X22B and QWEN-2-72B on the same set of claims. These models were prompted with individual stories as the context which makes their results comparable with other models marked with subscript "story."

📍 STORY LOCATION WITHIN COLLECTION			
MODEL	📄 BEGINNING	📄 MIDDLE	📄 END
📄 GPT-4o	64.0 <sub>(16/25)</sub>	50.0 <sub>(7/14)</sub>	50.0 <sub>(6/12)</sub>
📄 GPT-4o <sub>story</sub>	76.0 <sub>(19/25)</sub>	64.3 <sub>(9/14)</sub>	66.7 <sub>(8/12)</sub>
📄 GPT-4-TURBO	36.0 <sub>(9/25)</sub>	35.7 <sub>(5/14)</sub>	50.0 <sub>(6/12)</sub>
📄 GPT-4-TURBO <sub>story</sub>	60.0 <sub>(15/25)</sub>	42.9 <sub>(6/14)</sub>	58.3 <sub>(7/12)</sub>
📄 CLAUDE-3-OPUS	38.2 <sub>(13/34)</sub>	31.6 <sub>(6/19)</sub>	43.8 <sub>(7/16)</sub>
📄 CLAUDE-3-OPUS <sub>story</sub>	73.5 <sub>(25/34)</sub>	68.4 <sub>(13/19)</sub>	37.5 <sub>(6/16)</sub>
📄 CLAUDE-3.5-SONNET	38.2 <sub>(13/34)</sub>	31.6 <sub>(6/19)</sub>	50.0 <sub>(8/16)</sub>
📄 CLAUDE-3.5-SONNET <sub>story</sub>	61.8 <sub>(21/34)</sub>	47.4 <sub>(9/19)</sub>	31.2 <sub>(5/16)</sub>
📄 GEMINI PRO 1.5	52.2 <sub>(12/23)</sub>	70.0 <sub>(7/10)</sub>	40.0 <sub>(4/10)</sub>
📄 GEMINI PRO 1.5 <sub>story</sub>	56.5 <sub>(13/23)</sub>	70.0 <sub>(7/10)</sub>	40.0 <sub>(4/10)</sub>
📄 GEMINI FLASH 1.5	39.1 <sub>(9/23)</sub>	27.3 <sub>(3/11)</sub>	60.0 <sub>(6/10)</sub>
📄 GEMINI FLASH 1.5 <sub>story</sub>	34.8 <sub>(8/23)</sub>	63.6 <sub>(7/11)</sub>	30.0 <sub>(3/10)</sub>
📄 COMMAND R	20.0 <sub>(5/25)</sub>	21.4 <sub>(3/14)</sub>	33.3 <sub>(4/12)</sub>
📄 COMMAND R <sub>story</sub>	36.0 <sub>(9/25)</sub>	35.7 <sub>(5/14)</sub>	33.3 <sub>(4/12)</sub>
📄 COMMAND R+	20.0 <sub>(5/25)</sub>	0.0 <sub>(0/14)</sub>	8.3 <sub>(1/12)</sub>
📄 COMMAND R+ <sub>story</sub>	48.0 <sub>(12/25)</sub>	50.0 <sub>(7/14)</sub>	41.7 <sub>(5/12)</sub>
📄 PHI-3-MINI	6.2 <sub>(1/16)</sub>	8.3 <sub>(1/12)</sub>	14.3 <sub>(1/7)</sub>
📄 PHI-3-MINI <sub>story</sub>	18.8 <sub>(3/16)</sub>	33.3 <sub>(4/12)</sub>	0.0 <sub>(0/7)</sub>
📄 PHI-3-MINI <sub>simple</sub>	12.5 <sub>(3/24)</sub>	15.4 <sub>(2/13)</sub>	0.0 <sub>(0/12)</sub>
📄 PHI-3-MINI <sub>simple-story</sub>	37.5 <sub>(9/24)</sub>	30.8 <sub>(4/13)</sub>	25.0 <sub>(3/12)</sub>
📄 MIXTRAL-8X22B <sub>story</sub>	61.8 <sub>(21/34)</sub>	52.6 <sub>(10/19)</sub>	50.0 <sub>(8/16)</sub>
📄 QWEN-2-72B <sub>story</sub>	61.8 <sub>(21/34)</sub>	63.2 <sub>(12/19)</sub>	56.2 <sub>(9/16)</sub>

Table 31: Models’ performance on stories at different depths of the story collection. Each collection is divided into three parts based on the number of tokens: (1) beginning (first third), (2) middle (second third), and (3) end (last third). The subscript "story" (marked with 📄) refers to outputs obtained by prompting with the story (rather than the collection) as the context, for comparison. We also provide results for MIXTRAL-8X22B and QWEN-2-72B for comparison.

MODEL	QUOTATIONS (%)		
	ALL RESPONSES	CORRECT RESPONSES	INCORRECT RESPONSES
GPT-4O	18.4% <sub>227/1234</sub>	18.8% <sub>178/946</sub>	17.0% <sub>49/288</sub>
GPT-4-TURBO	18.8% <sub>232/1234</sub>	19.5% <sub>164/839</sub>	17.2% <sub>68/395</sub>
CLAUDE-3-OPUS	49.8% <sub>934/1874</sub>	49.8% <sub>685/1376</sub>	50.0% <sub>249/498</sub>
CLAUDE-3.5-SONNET	59.6% <sub>1116/1874</sub>	62.2% <sub>808/1300</sub>	53.7% <sub>308/574</sub>
GEMINI PRO 1.5	11.6% <sub>119/1029</sub>	11.4% <sub>84/738</sub>	12.0% <sub>35/291</sub>
GEMINI FLASH 1.5	10.8% <sub>111/1030</sub>	10.5% <sub>70/668</sub>	11.3% <sub>41/362</sub>

Table 32: Percentage of responses with identified quotations by model. Separate values are reported for claims labeled correctly and incorrectly. Note that the percentages are reported by claim, as each explanation is generated at the claim level. The counts (quotations/total) are provided in subscript. We do not report these numbers for open-weights models as the generations often do not follow the requested output format.



MODEL	GENRE		
	HISTORICAL	CONTEMPORARY	SPECULATIVE
GPT-4O	70.3 <sub>(26/37)</sub>	59.0 <sub>(229/388)</sub>	44.2 <sub>(72/163)</sub>
GPT-4-TURBO	70.3 <sub>(26/37)</sub>	39.9 <sub>(155/388)</sub>	34.4 <sub>(56/163)</sub>
CLAUDE-3-OPUS	63.5 <sub>(33/52)</sub>	51.1 <sub>(290/567)</sub>	42.9 <sub>(124/289)</sub>
CLAUDE-3.5-SONNET	42.3 <sub>(22/52)</sub>	42.2 <sub>(239/567)</sub>	38.1 <sub>(110/289)</sub>
GEMINI PRO 1.5	53.8 <sub>(28/52)</sub>	49.2 <sub>(127/258)</sub>	44.4 <sub>(84/189)</sub>
GEMINI FLASH 1.5	46.2 <sub>(24/52)</sub>	37.2 <sub>(96/258)</sub>	27.4 <sub>(52/190)</sub>
COMMAND R	30.3 <sub>(10/33)</sub>	17.8 <sub>(46/259)</sub>	20.1 <sub>(28/139)</sub>
COMMAND R <sub>simple</sub>	21.2 <sub>(7/33)</sub>	22.8 <sub>(59/259)</sub>	21.6 <sub>(30/139)</sub>
COMMAND R+	18.2 <sub>(6/33)</sub>	21.6 <sub>(56/259)</sub>	10.1 <sub>(14/139)</sub>
COMMAND R+ <sub>simple</sub>	15.2 <sub>(5/33)</sub>	13.1 <sub>(34/259)</sub>	15.1 <sub>(21/139)</sub>
PHI-3-MINI	10.0 <sub>(1/10)</sub>	9.5 <sub>(14/147)</sub>	9.0 <sub>(8/89)</sub>
PHI-3-MINI <sub>simple</sub>	20.0 <sub>(3/15)</sub>	15.8 <sub>(29/183)</sub>	12.6 <sub>(15/119)</sub>
GEMMA-10M	1.5 <sub>(1/67)</sub>	3.2 <sub>(18/567)</sub>	4.4 <sub>(15/338)</sub>
GEMMA-10M <sub>simple</sub>	4.5 <sub>(3/67)</sub>	7.2 <sub>(41/567)</sub>	7.7 <sub>(26/338)</sub>
LONGLLAMA <sub>simple</sub>	0.0 <sub>(0/52)</sub>	5.1 <sub>(29/567)</sub>	4.2 <sub>(12/289)</sub>
BM25+GPT-4o ( <i>k</i> =5)	35.8 <sub>(24/67)</sub>	28.6 <sub>(162/567)</sub>	26.0 <sub>(88/338)</sub>
BM25+GPT-4o ( <i>k</i> =25)	50.7 <sub>(34/67)</sub>	42.9 <sub>(243/567)</sub>	44.4 <sub>(150/338)</sub>
BM25+GPT-4o ( <i>k</i> =50)	56.7 <sub>(38/67)</sub>	50.8 <sub>(288/567)</sub>	45.9 <sub>(155/338)</sub>

Table 33: Model performance by genre: *historical* (pre-WWII), *contemporary* (post-WWII), and *speculative* (fantasy/SF/ghosts).

model <- glmer(accuracy ~ model + (1 id) + (1 annotator), common_set_data)				
PREDICTORS	ODDS RATIOS	CI (95%)	p-VALUE	
(Intercept)	1.73	1.18 – 2.55	0.005	**
GPT-4-TURBO	0.33	0.23 – 0.48	<0.001	***
CLAUDE-3-OPUS	0.64	0.45 – 0.91	0.014	*
CLAUDE-3.5-SONNET	0.34	0.24 – 0.50	<0.001	***
GEMINI PRO 1.5	0.55	0.38 – 0.78	0.001	***
GEMINI FLASH 1.5	0.24	0.17 – 0.35	<0.001	***
RANDOM EFFECTS				
$\sigma^2$ (residual variance)	3.29			
$\tau_{00}$ (id) (variance of random intercepts)	2.54			
$\tau_{00}$ (annotator) (variance of random intercepts)	0.19			
ICC	0.45			
N (ID)	354			
N (ANNOTATOR)	15			
OBSERVATIONS	2124			
R <sup>2</sup> (marginal)	0.036			
R <sup>2</sup> (conditional)	0.473			

Table 34: Summary of generalized linear mixed model with **model** as the predictor of **accuracy**: model <- glmer(accuracy ~ model + (1|id) + (1|annotator), data). GPT-4o was set as the reference level (intercept). Note that while we observe significant differences between the models' performance, the marginal R<sup>2</sup> is low, suggesting that model type alone does not explain the majority of the variance in the data. See Table 35 for post-hoc analysis.

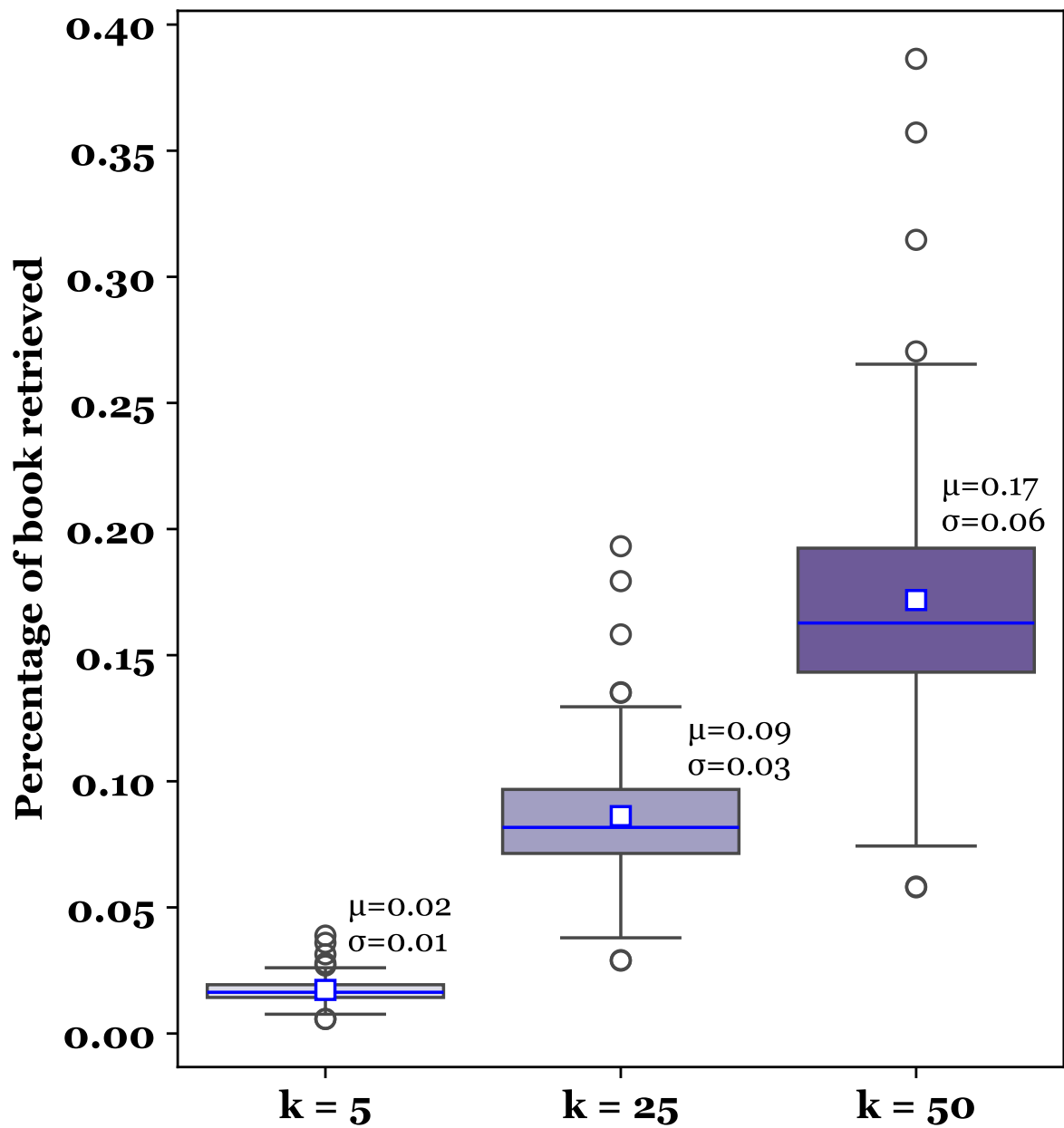


Figure 13: Average percentage of the book retrieved by BM25 for varying values of top k. For example, k = 5 means that the top 5 ranked excerpts according to BM25 were fed to GPT-4o as context for a claim.

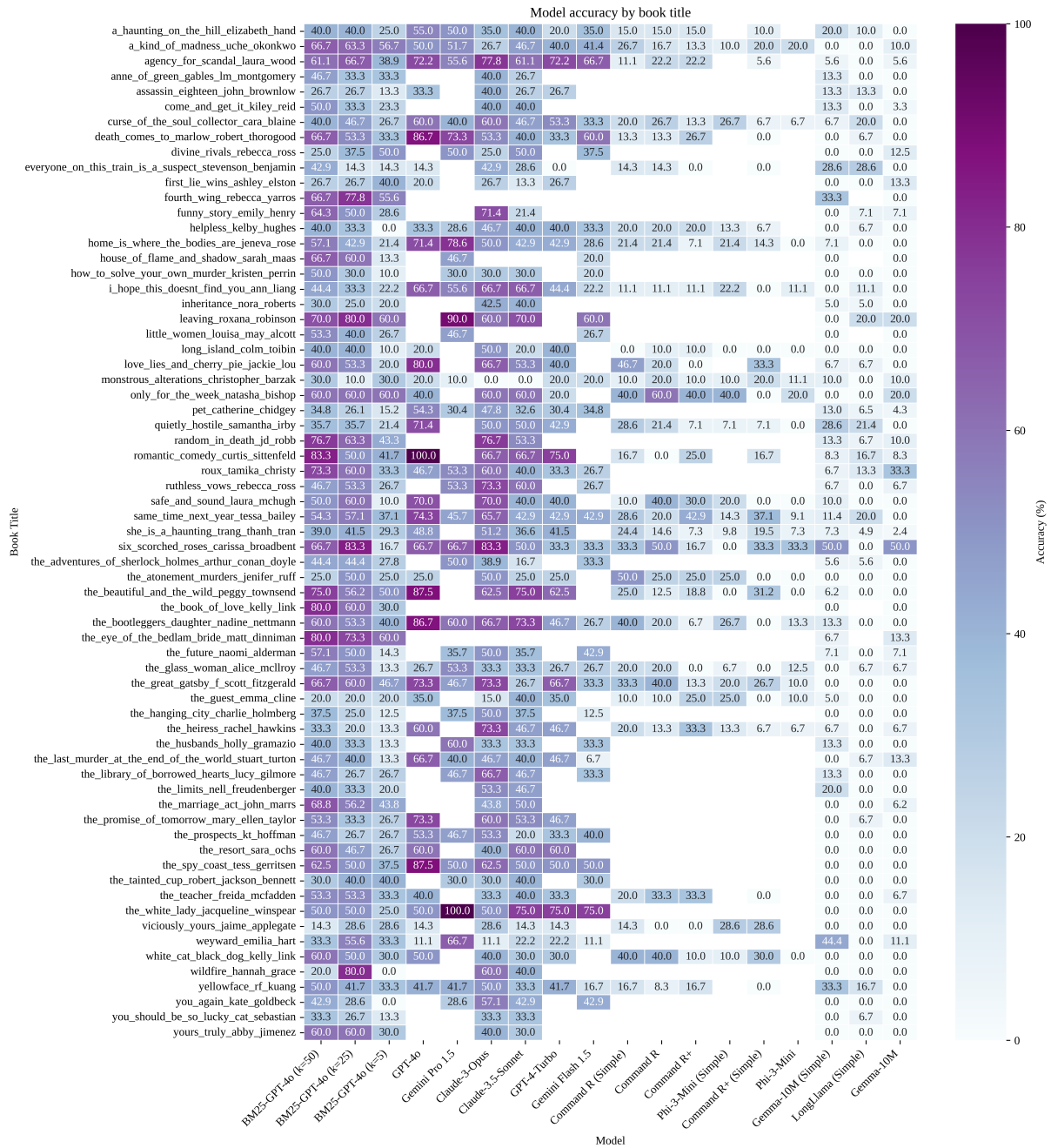


Figure 14: Heatmap of model performance by book title. Empty cells indicate cases where the model did not process the book (GEMINI PRO 1.5 and GEMINI FLASH 1.5) or could not process the book because of its context window.

CONTRAST	ESTIMATE	SE	ODDS RATIO	PROBABILITY	<i>p</i> -VALUE
GPT-4o - CLAUDE-3-OPUS	0.4480	0.1820	1.5652	0.6102	0.1355
GPT-4o - CLAUDE-3.5-SONNET	1.0654	0.1849	2.9021	0.7437	<0.001 ***
GPT-4o - GEMINI FLASH 1.5	1.4229	0.1888	4.1492	0.8058	<0.001 ***
GPT-4o - GEMINI PRO 1.5	0.6013	0.1822	1.8244	0.6459	0.0124 *
GPT-4o - GPT-4-TURBO	1.1005	0.1852	3.0056	0.7503	<0.001 ***
CLAUDE-3-OPUS - CLAUDE-3-OPUS	0.6174	0.1820	1.8541	0.6496	0.0090 **
CLAUDE-3-OPUS - GEMINI FLASH 1.5	0.9749	0.1855	2.6509	0.7261	<0.001 ***
CLAUDE-3-OPUS - GEMINI PRO 1.5	0.1533	0.1801	1.1656	0.5382	0.9578
CLAUDE-3-OPUS - GPT-4-TURBO	0.6524	0.1822	1.9202	0.6576	0.0046 **
CLAUDE-3.5-SONNET - GEMINI FLASH 1.5	0.3575	0.1852	1.4297	0.5884	0.3833
CLAUDE-3.5-SONNET - GEMINI PRO 1.5	-0.4642	0.1816	0.6287	0.3860	0.1084
CLAUDE-3.5-SONNET - GPT-4-TURBO	0.0350	0.1826	1.0356	0.5088	1.0000
GEMINI FLASH 1.5 - GEMINI PRO 1.5	-0.8216	0.1849	0.4397	0.3054	<0.001 ***
GEMINI FLASH 1.5 - GPT-4-TURBO	-0.3224	0.1853	0.7244	0.4201	0.5050
GEMINI PRO 1.5 - GPT-4-TURBO	0.4992	0.1819	1.6474	0.6223	0.0667

Table 35: Post-hoc comparisons of models for **accuracy** (Table 34) using Tukey adjustments for multiple comparisons. The probability values refer to the likelihood that the first model in each contrast is more accurate than the second model (i.e., a value of 0.5 suggests that both models are comparable in terms of accuracy).

model <- glmer(accuracy ~ length_group + (1 pairID) + (1 annotator), data)					
PREDICTORS		ODDS RATIOS	CI (95%)	<i>p</i> -VALUE	
(Intercept)		0.56	0.34 – 0.94	0.028	*
75K–127K		1.45	0.87 – 2.43	0.154	
127K–180K		1.26	0.74 – 2.14	0.391	
ABOVE 180K		0.75	0.28 – 1.96	0.551	
RANDOM EFFECTS					
$\sigma^2$ (residual variance)		3.29			
$\tau_{00}$ (id) (variance of random intercepts)		2.15			
$\tau_{00}$ (annotator) (variance of random intercepts)		0.24			
ICC		0.42			
N (ID)		967			
N (ANNOTATOR)		22			
OBSERVATIONS		4137			
$R^2$ (marginal)		0.004			
$R^2$ (conditional)		0.423			

Table 36: Summary of generalized linear mixed model with **length group** as the predictor of **accuracy**: model <- glmer(accuracy ~ length\_group + (1|pairID) + (1|annotator), data). The "below 75k" group was set as the reference level (intercept). Note that the marginal  $R^2$  is very low, indicating that length group alone does not explain the majority of the variance in the data. See Table 37 for post-hoc analysis.

CONTRAST	ESTIMATE	SE	ODDS RATIO	PROBABILITY	<i>p</i> -VALUE
BELOW 75K - 127K-180K	-0.2312	0.2698	0.7936	0.4424	0.8267
BELOW 75K - 75K-127K	-0.3739	0.2620	0.6880	0.4076	0.4823
BELOW 75K - ABOVE 180K	0.2938	0.4925	1.3416	0.5729	0.9331
127K-180K - 75K-127K	-0.1427	0.1652	0.8670	0.4644	0.8235
127K-180K - ABOVE 180K	0.5251	0.4461	1.6906	0.6283	0.6415
75K-127K - ABOVE 180K	0.6677	0.4464	1.9498	0.6610	0.4399

Table 37: Post-hoc comparisons of length groups for **accuracy** (Table 36) using Tukey adjustments for multiple comparisons. The probability values refer to the likelihood that the verification for the first group in each contrast is more accurate than the verification for the second group.

model <- glmer(accuracy ~ year + (1 pairID) + (1 annotator), data)			
PREDICTORS	ODDS RATIOS	CI (95%)	p-VALUE
(Intercept)	0.65	0.35 – 1.20	0.170
2023	1.17	0.64 – 2.15	0.321
2024	1.17	0.63 – 2.16	0.344
RANDOM EFFECTS			
$\sigma^2$ (residual variance)	3.29		
$\tau_{00}$ (id) (variance of random intercepts)	2.16		
$\tau_{00}$ (annotator) (variance of random intercepts)	0.22		
ICC	0.42		
N (ID)	967		
N (ANNOTATOR)	22		
OBSERVATIONS	4137		
$R^2$ (marginal)	0.000		
$R^2$ (conditional)	0.420		

Table 38: Summary of generalized linear mixed model with **year** as the predictor of **accuracy**: model <- glmer(accuracy ~ year + (1|pairID) + (1|annotator), data). The "classics" year group was set as the reference level (intercept). Note that the marginal  $R^2$  is 0.000, indicating that year alone does not explain the majority of the variance in the data.

model <- glmer(accuracy ~ genre + (1 id) + (1 annotator), data = filtered_data)				
PREDICTORS	ODDS RATIOS	CI (95%)	p-VALUE	
(Intercept)	0.55	0.40 – 0.74	<0.001	***
CONTEMPORARY	1.53	1.11 – 2.11	0.010	*
HISTORICAL	2.18	1.20 – 3.96	0.010	*
RANDOM EFFECTS				
$\sigma^2$ (residual variance)	3.29			
$\tau_{00}$ (id) (variance of random intercepts)	2.10			
$\tau_{00}$ (annotator) (variance of random intercepts)	0.16			
ICC	0.41			
N (ID)	938			
N (ANNOTATOR)	22			
OBSERVATIONS	3991			
$R^2$ (marginal)	0.010			
$R^2$ (conditional)	0.412			

Table 39: Summary of generalized linear mixed model with **genre** as the predictor of **accuracy**: model <- glmer(accuracy ~ genre + (1|id) + (1|annotator), data = filtered\_data). The "speculative" genre was set as the reference level (intercept). We also excluded two books which fell into "historical and contemporary" and "essays" categories. Note that the marginal  $R^2$  is low, indicating that genre alone does not explain the majority of the variance in the data. The post-hoc analysis for this model is presented in [Table 40](#).

CONTRAST	ESTIMATE	SE	ODDS RATIO	PROBABILITY	<i>p</i> -VALUE
SPECULATIVE - CONTEMPORARY	-0.4234	0.1648	0.6548	0.3957	0.0275 *
SPECULATIVE - HISTORICAL	-0.7802	0.3042	0.4583	0.3143	0.0278 *
CONTEMPORARY - HISTORICAL	-0.3568	0.2874	0.6999	0.4117	0.4287

Table 40: Post-hoc comparisons of genres for **accuracy** (Table 39) using Tukey adjustments for multiple comparisons. The probability values refer to the likelihood that the validation of books of the first genre in each contrast is more accurate than the validation of books of the second genre (i.e., a value of 0.5 suggests that both genres are comparable in terms of accuracy).

```
model <- glmer(accuracy ~ scope + (1|id) + (1|annotator), data = filtered_data)
```

PREDICTORS	ODDS RATIOS	CI (95%)	<i>p</i> -VALUE
(Intercept)	2.12	0.74 – 6.11	0.163
PASSAGE	0.35	0.14 – 0.87	0.024 *
GLOBAL	0.32	0.13 – 0.78	0.012 *

RANDOM EFFECTS	
$\sigma^2$ (residual variance)	3.29
$\tau_{00}$ (id) (variance of random intercepts)	1.19
$\tau_{00}$ (annotator) (variance of random intercepts)	0.53
ICC	0.34
N (ID)	121
N (ANNOTATOR)	4
OBSERVATIONS	544

$R^2$ (marginal)	0.030
$R^2$ (conditional)	0.363

Table 41: Summary of generalized linear mixed model with **scope** as the predictor of **accuracy**: `model <- glmer(accuracy ~ scope + (1|id) + (1|annotator), data = filtered_data)`. The "sentence" scope was set as the reference level (intercept). Note that the marginal  $R^2$  is low, indicating that scope alone does not explain the majority of the variance in the data. See Table 42 for the post-hoc analysis.

CONTRAST	ESTIMATE	SE	ODDS RATIO	PROBABILITY	<i>p</i> -VALUE
SENTENCE - GLOBAL	1.1322	0.4506	3.1026	0.7563	0.0321 *
SENTENCE - PASSAGE	1.0488	0.4642	2.8542	0.7405	0.0617 .
GLOBAL - PASSAGE	-0.0834	0.3139	0.9199	0.4792	0.9618

Table 42: Post-hoc comparisons of scopes for **accuracy** (Table 41) using Tukey adjustments for multiple comparisons. The probability values refer to the likelihood that the validation of claim pairs in the first scope in each contrast is more accurate than the validation of claim pairs in the second scope (i.e., a value of 0.5 suggests that both scopes are comparable in terms of accuracy).