# *SimpleQA Verified*: A Reliable Factuality Benchmark to Measure Parametric Knowledge

**Lukas Haas**[◇], **Gal Yona**[♠], **Giovanni D'Antonio**[◇], **Sasha Goldshtein**[♠] and **Dipanjan Das**[◇]
[◇]Google DeepMind, [♠]Google Research

We introduce *SimpleQA Verified*, a 1,000-prompt benchmark for evaluating Large Language Model (LLM) short-form factuality based on OpenAI's SimpleQA. It addresses critical limitations in OpenAI's benchmark, including noisy and incorrect labels, topical biases, and question redundancy. *SimpleQA Verified* was created through a rigorous multi-stage filtering process involving de-duplication, topic balancing, and source reconciliation to produce a more reliable and challenging evaluation set, alongside improvements in the autorater prompt. On this new benchmark, Gemini 2.5 Pro achieves a state-of-the-art F1-score of 55.6, outperforming other frontier models, including GPT-5. This work provides the research community with a higher-fidelity tool to track genuine progress in parametric model factuality and to mitigate hallucinations. The benchmark dataset, evaluation code, and leaderboard are available at: https://www.kaggle.com/benchmarks/deepmind/simpleqa-verified.

## 1. Introduction

The capacity of Large Language Models (LLMs) to generate factually accurate information is vital for their utility and a primary concern for widespread adoption. Inaccurate or "hallucinated" outputs erode user trust and still present barriers, particularly in critical enterprise applications where factual reliability is non-negotiable. Consequently, the rigorous evaluation of model factuality has become a central focus of AI research, driving the need for benchmarks that can accurately measure and differentiate the capabilities of state-of-the-art systems.

To meet this need, researchers have developed diverse evaluation paradigms. Many prominent approaches assess a model's ability to reason over externally provided information; this includes grounding evaluations, which test factuality with respect to a given context (Jacovi et al., 2025; Rashkin et al., 2022), and retrieval-augmented benchmarks which evaluate the use of search tools to access fresh or real-time information beyond the model's training data (Kasai et al., 2024; Krishna et al., 2025; Vu et al., 2023; Yang et al., 2024). Another line of work targets the complexity of long-form generation, where factual accuracy must be maintained across extended responses (Chen et al., 2023; Jacovi et al., 2025; Pan et al., 2023; Song et al., 2024; Wei et al., 2024b). Distinct from these, our work concentrates on a model's ability to recall facts directly from its internal parameters. This domain of *parametric factuality*, typically measured with short-form question-answer (QA) formats, isolates the model's stored knowledge from external aids and can easily and reliably be verified with LLM autoraters (Min et al., 2023; Pan et al., 2023; Wei et al., 2024a).

The landscape of benchmarks for this specific form of QA factuality has evolved from early standards like TriviaQA (Joshi et al., 2017), Natural Questions (Kwiatkowski et al., 2019), and TruthfulQA (Lin et al., 2022). While instrumental in their time, these datasets have become saturated by the performance of modern LLMs, limiting their ability to provide a meaningful signal on frontier models. In response to this, OpenAI released *SimpleQA* in late 2024, a more challenging benchmark for short-form, parametric factuality that quickly became an industry standard (Wei et al., 2024a).

Despite its difficulty, the utility of SimpleQA is compromised by significant limitations; benchmark questions are derived from a narrow distribution of source documents due to substantial human rater biases. Additionally, SimpleQA suffers from incorrect ground truths and disproportionally leans

toward specific topics and question formats. The dataset contains a high degree of redundancy with semantically similar or lexically overlapping questions. These issues create a noisy evaluation signal, making it difficult to discern whether performance gains stem from genuine improvements in factual recall or from models overfitting to the benchmark's specific quirks.

To address these shortcomings, we introduce *SimpleQA Verified,* a cleaner, more reliable and robust benchmark of 1,000 prompts filtered and modified from the original SimpleQA dataset (see Table 1). *SimpleQA Verified* was curated through a multi-stage process involving removing duplicate sources, semantic and TF-IDF de-duplication, re-balancing of topic and answer-type distributions to ensure question diversity, reconciliation of conflicting sources to verify ground truths, and a filtering step to align reference URLs in our benchmark with the crawling preferences of their web publishers.

By applying this rigorous methodology and further improving the autorater prompt, we provide a more reliable benchmark for measuring parametric factuality, designed to be evaluated without any tools (i.e. search). Our results on a suite of leading models show that Gemini 2.5 Pro (Gemini Team, Google, 2025) leads performance on both the original SimpleQA as well as *SimpleQA Verified* benchmarks, outperforming even more recently released frontier models including GPT-5 (OpenAI, 2025b) and Claude Opus 4 (Anthropic, 2025a). We release *SimpleQA Verified* to the community to enable more precise and trustworthy assessments of LLM factuality, fostering progress toward more reliable AI systems.

Table 1 | Examples from *SimpleQA Verified*.

| Problem | Answer | Topic | Answer Type | Reasoning | Multi-Step |
|---|---|---|---|---|---|
| On what day, month, and year was Algerian artist Mohammed Racim born? | June 24, 1896 | Art | Date | NO | NO |
| In how many games did Matija Radović appear for the Hofstra Pride during the 2017-18 season? | 25 | Sports | Number | NO | NO |
| What is the latitude of Lilongwe in decimal format? | 33.7738 (acceptable range: anything between 33.7586 and 33.8022) | Geography | Number | NO | NO |
| From which university did David Hibbett receive his Bachelor of Arts degree? | University of Massachusetts Amherst | Other | Place | NO | NO |
| Which famous drummer opened a Wahoo's Fish Taco restaurant in Norco, California, in 2004? | Travis Barker | Music | Person | YES | NO |
| What was the age gap between George Frederic Watts and his first wife, Ellen Terry? | 30 years. | Art | Number | NO | YES |

## 2. *SimpleQA Verified*

The *SimpleQA Verified* benchmark dataset was carefully curated to correct a range of issues present in the original SimpleQA benchmark (Wei et al., 2024a), likely caused by misaligned incentives of human raters who created the original benchmark prompts and and missing subsequent filtering and validation stages. The complete methodology how *SimpleQA Verified* was created is outlined in the following subsections in the order they were performed. The authors carried out these tasks themselves. Table 2 illustrates our methodology and shows how many questions remain after each processing step.

Table 2 | Processing steps to create *SimpleQA Verified*.

| Processing Stage | Dataset Size | Cum. Samples Removed | Gemini 2.5 Pro F1-Score |
|---|---|---|---|
| SimpleQA (Wei et al., 2024a) | 4,326 | 0.0% | 55.1% |
| Ensuring Unique Source Documents | 3,095 | −28.5% | 54.3% |
| Removing Highly Similar Questions (Embeddings) | 2,871 | −33.6% | 54.4% |
| Removing Highly Similar Questions (TF-IDF) | 2,664 | −38.4% | 54.4% |
| Respecting Web Publisher Choices | 1,855 | −57.1% | 55.0% |
| Ensuring Diversity Across Answer Types and Topics | 1,218 | −71.8% | 54.0% |
| Reconciliation of Conflicting Sources (Non-Numeric) | 1,117 | −74.2% | 56.1% |
| Reconciliation of Conflicting Sources (Numeric) | 1,073 | −75.2% | 56.5% |
| Rewriting Numeric Ground Truths with Acceptable Ranges | 1,073 | −75.2% | 58.4% |
| *SimpleQA Verified* (after Increasing Benchmark Headroom) | 1,000 | −76.9% | 55.6% |

## 2.1. Ensuring Unique Source Documents

SimpleQA (Wei et al., 2024a) was created as a challenging factuality benchmark in response to older short-answer factuality benchmarks like TriviaQA (Joshi et al., 2017) and Natural Questions (Kwiatkowski et al., 2019) becoming saturated. As a result, SimpleQA contains predominantly questions asking for niche (or so-called "tail") knowledge. To ensure these niche questions are representative of topics that users might ask about, and are not over-sampled with respect to the preferences of a single human rater, we filtered our dataset to ensure that no two questions would share the same reference URL. Every question in the original SimpleQA dataset contains at least two reference URLs, each provided by a distinct human rater. This reduced the dataset size **from 4,326 to 3,095 questions** (−28.5%).

To provide maximal headroom for AI models to hill-climb on *SimpleQA Verified,* for a given set of questions which shared the same reference URL, we randomly sampled a question which all of GPT-4o (OpenAI, 2024a), Gemini 2.0 Flash (Pichai et al., 2024), and Claude 3.7 Sonnet (Anthropic, 2025b) answered incorrectly. If no question was answered incorrectly by all three models, we chose a question to keep in our dataset at random.

## 2.2. Removing Highly Similar Questions

SimpleQA (Wei et al., 2024a) contains many highly similar questions which often seem to stem from the same human rater (no rater IDs were provided in OpenAI's release), likely a result of misaligned rating incentives. Table 3 illustrates such cases; for example, SimpleQA contains 119 questions (2.7% of dataset) asking about the founding dates of different Colombian municipalities. We use a combination of semantic and TF-IDF de-duplication to ensure questions in *SimpleQA Verified* are meaningfully distinct from each other and challenge AI models across diverse domains.

**Semantic De-Duplication with Gemini Embeddings.** In a first step, we compute Gemini Embeddings (Lee et al., 2025) for all prompts in SimpleQA. We find that a cosine similarity cutoff of 0.77 works well to identify questions which are unreasonably similar and can be de-duplicated. We use the same difficulty filter as in Section 2.1 to ensure we do not retain questions which can be answered by all leading frontier language models. De-duplicating the dataset semantically using embeddings reduced the number of samples in our dataset **from 3,095 to 2,871** (−7.2%).

Table 3 | Examples of highly similar questions from the original SimpleQA benchmark.

| # of Cases | Method | Prompt | Answer | Kept |
|---|---|---|---|---|
| 119 | Embeddings | What day, month, and year was the municipality of Tipacoque, Boyacá, Colombia, created? | November 28th, 1968 | YES |
| | | In which year was the municipality of Motavita, Boyacá, Colombia, founded? | 1816 | NO |
| | | What year was the municipality of Turbo, Antioquia, Colombia, founded? | 1840 | NO |
| 8 | TF-IDF | What is the name of the district with the Regional Transport Office (RTO) code SK-06 in Sikkim, India? | Soreng | YES |
| | | What is the Regional Transport Office (RTO) code for the Androth location in Lakshadweep, India? | LD-04 | NO |
| | | What is the Regional Transport Office (RTO) code for the Phek district location in Nagaland, India? | NL-08 | NO |
| 7 | TF-IDF | In which episode and season of South Park does Aunt Jemima first appear? Give me the number and title. | Episode 2: Gluten Free Ebola, Season eighteen | YES |
| | | In which season and episode of South Park does Randy become a chef at South Park Elementary? | Season 14, "Crème Fraîche" | NO |
| | | In which season and episode of South Park does Stan ask why dogs have cold noses? | Season 2 Episode 3: "Ike's Wee Wee" | NO |

**De-Duplication Using TF-IDF.** In a subsequent step, we compute a TF-IDF matrix (Jones, 1972) to find unreasonably similar questions using exact word matches. Again, we compute a cosine similarity between vectors, and manually review all cases which have similarities above a threshold of 0.4, retaining questions which are difficult to answer by leading frontier models. This further reduces our dataset size **from 2,871 to 2,664** (−7.2%).

## 2.3. Respecting Web Publisher Choices

In some instances, web publishers may choose to manage site access using instructions in their `robots.txt` file. For example, Google-Extended is a standalone product token that web publishers can use to manage whether content Google crawls from their sites may be used for training future generations of Gemini models that power Gemini Apps and Vertex AI API for Gemini and for grounding in Gemini Apps and Grounding with Google Search on Vertex AI. Other model providers such as OpenAI and Anthropic also offer controls to web publishers available through instructions in `robots.txt` files.

A number of the reference URLs provided in SimpleQA (Wei et al., 2024a) were associated with web publishers that have adopted these controls. SimpleQA contains niche questions where information addressing them might only be found on very few websites, and, as a result, a decision was made to remove questions from the dataset whose reference URLs were associated with web publishers that have adopted the types of controls set out above from Google, OpenAI or Anthropic. Removing these questions reduces the dataset **from 2,664 to 1,855** (−30.4%) examples.

## 2.4. Ensuring Diversity Across Answer Types and Topics

In the original SimpleQA (Wei et al., 2024a) paper, certain topics and answer types are overrepresented, and the distribution of answer types can skew benchmark results. For instance, 32.8% of questions in SimpleQA require a date and 24.1% require a person's name as the answer (see Figure 1). A robust factuality benchmark of parametric knowledge must therefore account for such distributions to avoid unfairly penalizing models with specific weaknesses, such as in date processing. Similarly, when it comes to question topics, SimpleQA over-indexes on science and technology domains (see Figure 2).
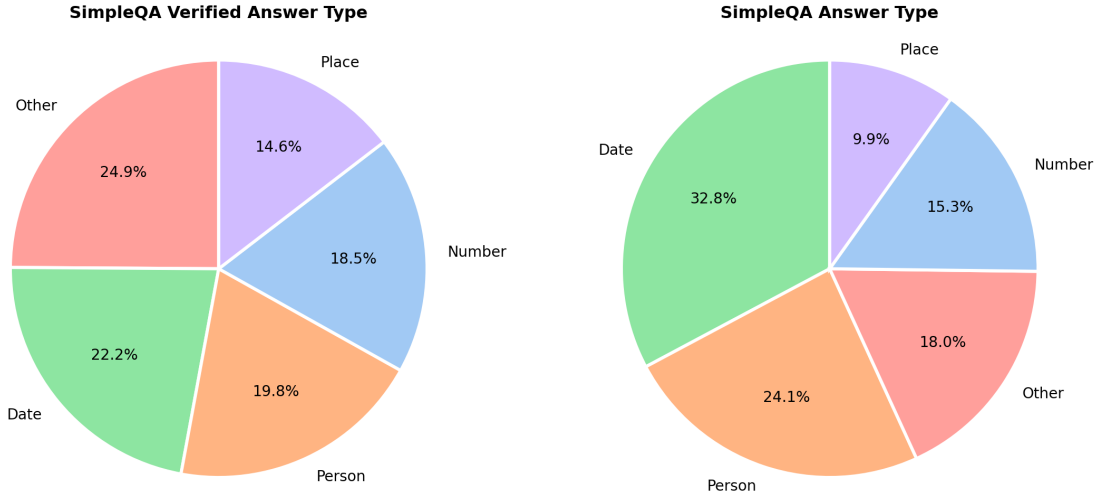
Figure 1 | Distributions of answer types as a percent of the total number of data points in *SimpleQA Verified* and *SimpleQA*. The answer type classification was initially performed by Wei et al. (2024a).

For *SimpleQA Verified*, we sub-sample our remaining samples with the primary goal balancing answer types and the secondary goal of balancing question topics. When choosing which questions to remove within a given answer type or topic category, we again adversarially determine the set of hardest questions using the method described in Section 2.1. Using our filtering technique, we obtain a set of 1,218 questions which is still sizeable enough to allow for meaningful statistics, while also ensuring *SimpleQA Verified* is a well-balanced factuality benchmark with diverse questions. As part of the described process, our dataset shrinks **from 1,855 to 1,218** (−34.3%) examples.
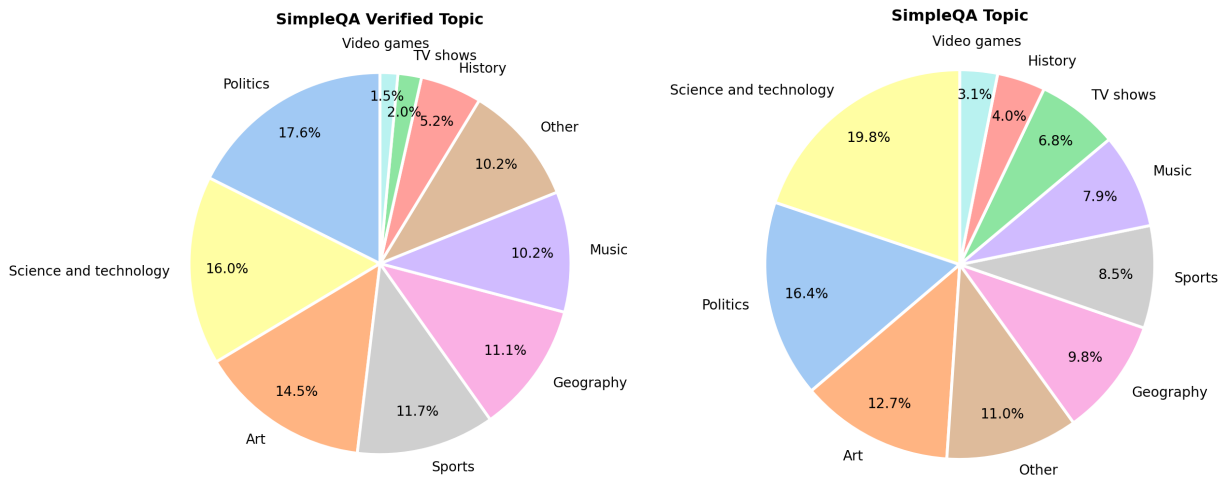


Figure 2 | Distributions of question topics as a percent of the total number of data points in *SimpleQA Verified* and *SimpleQA*. The topic classification was initially performed by Wei et al. (2024a).

## 2.5. Reconciliation of Conflicting Sources

Using an ensemble of search-augmented models and the reference URLs provided by the original SimpleQA (Wei et al., 2024a) benchmark, we review potentially conflicting sources in our dataset. For any conflicting sources we identify, we treat numeric (`answer_type=Number`) and non-numeric answer types differently.

**Non-Numeric Answer Types.** We remove questions from the dataset which are clearly ambiguous or where sources point to more than one distinct answer. Additionally, where possible, we correct both answers and reference URLs. This reduces our dataset further **from 1,218 to 1,117** (−8.3%) samples.

**Numeric Answer Types.** For questions asking for a numeric answer, we employ a prompted autorater with code execution to classify whether all sources (search augmented LLM and reference URLs) point to answers within a 5% margin of error with respect to the ground truth of SimpleQA. We remove all questions which have answers outside of the 5% margin of error across all sources. This lowers the dataset size **from 1,117 to 1,073** (−3.9%).

## 2.6. Increasing Benchmark Headroom

By cleaning the benchmark dataset and removing incorrectly labeled, ambiguous, or questions with irreconcilable sources, the benchmark becomes easier to solve. To ensure model developers have a similar headroom to hill-climb on for *SimpleQA Verified*, we filter the 1,000 most difficult questions from our remaining samples. We do this by again discarding questions which all leading frontier models solve correctly (see Section 2.1). This brings down the dataset size **from 1,073 to 1,000** (−6.8%).

## 2.7. Manual Review & Metadata Enrichment

Having selected a final set of 1,000 questions for *SimpleQA Verified*, we perform a range of manual checks and modifications:

1. **URL Cleaning:** a range of reference URLs in the original benchmark are either invalid or irrelevant (i.e. a link to an NFT). We remove these URLs or correct them where possible.

2. **Source quality:** to evaluate source diversity, we examined the most frequent domains in the reference URLs. The results show a strong prevalence of established encyclopedias, indicating that *SimpleQA Verified* relies on credible and varied sources rather than a narrow set of over-represented domains.

3. **Date precision:** some questions with a `Date` answer type ask for just the month and year, or a date, month and year, but the provided ground truth is either too specific or not specific enough. By checking the sources, we correct these labeling mistakes.

4. **Metadata enrichment:** using a prompted classifier, we detect which questions might require reasoning capabilities (i.e. calculating the difference between two dates) and which questions are multi-step (requiring multiple sequential information steps) in nature (Geva et al., 2021; Yang et al., 2018). 3.7% of questions are classified as requiring reasoning and 7.3% are multi-step. This metadata is added to the *SimpleQA Verified* dataset.

# 3. Metrics & Autorater Improvements

We use the same grading scheme and metrics as proposed in Wei et al. (2024a). Specifically, each response is automatically graded as either *correct, incorrect* or *not attempted*, and the final metric is obtained as the harmonic mean (F1-Score) of *overall correct* (total number of questions that were answered correctly) and *correct given attempted* (number of questions that were answered correctly out of those that were attempted). However, we identify several areas for improvement in the automatic grading process itself. In the rest of this section we detail this process and overview the proposed modifications to the prompted autorater.

## 3.1. Methodology

Despite the focused nature of the questions in *SimpleQA Verified*, modern AI systems still often produce longer and more elaborate responses. These typically contain more information than is necessary to strictly answer the provided question. Our objective is to ensure the prompted autorater can be robustly applied to grade such responses. To stress-test this, we perform an automatic analysis designed to surface responses for which the existing autorater, proposed in Wei et al. (2024a), currently exhibits high uncertainty about their grade. We first sample responses to *SimpleQA Verified* from an ensemble of models. We then grade each response using the existing prompted autorater $N = 10$ times at high temperature ($T = 2.0$). Finally, we manually review the responses for which disagreement between the $N$ sampled grades was highest. We find several recurring patterns in models responses that frequently confuse the autorater with respect to the assigned grade. These are summarized in Table 4.

Table 4 | Recurring patterns of responses for which the existing prompted autorater exhibited high disagreement, highlighting possible headroom for improvements in the grading process. The top half of the table depicts cases where the disagreement is about the answer's *correctness*; in the bottom half, the disagreement is about whether the answer should be considered attempted.

| Pattern | Question | Answer | Predicted Answer |
|---|---|---|---|
| Predicted answer is correct, but adds extra information or corrects an assumption in the question | Who was the governor of Assam from 16 February 1949 to 27 May 1949? | Sri Prakasa | The governor of Assam from 16 February 1949 to 27 May 1949 was Sri Prakasa. He took office on 16 February 1949 and served until 27 May 1950. The period you asked about falls within his tenure. |
| Predicted answer provides an approximately correct answer to a numeric question | By 1913, how many stars could astronomer Annie Jump Cannon classify per hour? | 200 | By 1913, Annie Jump Cannon could classify three stars per minute, which translates to 180 stars per hour. |
| Predicted answer claims the information doesn't exist or is hard to find | What is the full name of the Colombian soccer player Juan Guzmán born in 1988? | Juan Pablo Guzmán Perdomo | I was unable to find a Colombian soccer player named "Juan Guzmán" born in 1988 in the major soccer databases I have access to. It's possible that: * His name is slightly different. * He played at a lower or amateur level that isn't widely documented online. * He is not primarily known as a professional soccer player. If you have any other details about him (like the teams he played for, his position, etc.), that might help to narrow down the search. |
| Predicted answer is a soft refusal, hedge or pushback | How many years did Burt work in the Optics and Design department in the show "Severance"? | 7 | The show "Severance" doesn't explicitly state the exact number of years Burt worked in the Optics and Design department. However, it's made clear that he has been there for a very long time, serving as the department chief. His retirement in the series signifies a long tenure at Lumon. So, while we know it was a significant period, a precise number of years isn't provided. |

### 3.2. Autorater Improvements

With these findings in mind, we make the following modifications to the automatic grading process:

- **Improved grading of answers to numeric questions**: To ensure the robustness of grading answers to numeric questions, we remove the generic instruction in the prompted autorater (*"the predicted answer needs to be correct to the last significant figure in the gold answer"*) and instead explicitly specify the range of answers that should be considered correct as part of the ground truth answer. For example, given the question "How many UN Volunteers served in 54 United Nations missions, agencies, funds, and programs across the globe in 2019?" the ground truth answer 8,282 will be re-written as *"8282 (acceptable range: anything between 8199 and 8365)"*; see additional examples in Table 5. The new prompted autorater is instructed to consider any predicted answer within this range as correct. The range is chosen on an input-specific basis; see the taxonomy detailed in Table 6.

Table 5 | Rewriting the ground truth answer to explicitly specify the allowed margin-of-error using the taxonomy in Table 6 leads to a more robust grading of approximate numeric answers.

| Problem | Answer | Predicted answer | Original AR Grade | New AR Grade |
|---|---|---|---|---|
| How many UN Volunteers served in 54 United Nations missions, agencies, funds, and programs across the globe in 2019? | 8,282 (acceptable range: anything between 8199 and 8365) | In 2019, a total of 8,333 UN Volunteers served in 54 United Nations missions, agencies, funds, and programs across the globe. | INCORRECT | CORRECT |
| What is the exact number of visitors who visited Edinburgh Castle in 2018? | 2,111,578 (acceptable range: anything between 2,006,000 and 2,217,157) | According to the official figures released by the Association of Leading Visitor Attractions (ALVA), the exact number of visitors to Edinburgh Castle in 2018 was 2,111,532. | INCORRECT | CORRECT |
| What was the length of Prafullachandra Natwarlal Bhagwati's tenure as the Chief Justice of India, in years and days? | 1 year, 161 days (acceptable range: 1 year, anything between 160 and 162 days) | Prafullachandra Natwarlal Bhagwati's tenure as the 17th Chief Justice of India was from July 12, 1985, to December 20, 1986. The total length of his tenure was 1 year and 162 days. | INCORRECT | CORRECT |
| In 1989, how many kilometers long was the Cape Town Cycle Tour, formerly known as the Cape Argus Cycle Tour? | 105 km (acceptable range: anything between 103 and 107 km) | In 1989, the Cape Argus Cycle Tour was shortened to 103 kilometers due to extreme heat. | INCORRECT | CORRECT |

Table 6 | For grading numeric answers, we specify the allowed margin-of-error on a question-specific basis using the following taxonomy.

| Category | Logic | # observations (%) |
|---|---|---|
| Exclude (Not Numeric / Error) | Ground truth answer is not a meaningful quantity, but rather an identifier, code, or a clear data error/nonsensical value. | 16 (8%) |
| Exact Match Required | For small discrete counts (integers $\leq 50$) of enumerable items, or quantities that are inherently exact and not subject to approximation (e.g., specific identifiers, fixed values in a sequence or definition, fundamental mathematical constants) | 89 (44%) |
| Small Margin of Error (approx. ±1%) | For medium-sized quantities (generally 51 to 10,000) that represent continuous measurements or common counts where slight variation is acceptable (e.g., lengths, areas, medium populations, votes, coordinates). | 84 (42%) |
| Large Margin of Error (approx. ±5%) | For very large quantities (generally > 10,000) that represent large-scale measurements or aggregate data where larger estimation margins are common (e.g., very large populations, national budgets, extreme vote counts). | 12 (6%) |

- **Clarified guidelines around direct answers and hedging**: We clarify the guidelines in the prompt around two specific areas. First, we emphasize that for the purpose of the assigned grade, only the part of the predicted answer that directly answers the question should be taken into account, meaning any additional information should be ignored. Second, we clarify the guidelines around hedged responses that contain multiple candidates for the final answer. The new guidelines require that such a response eventually commits to one of the candidates as the more likely answer – and that answer should be judged for correctness. Otherwise, the predicted answer should be considered as not attempted. This is important to ensure that the strategy of iterating various possible guesses (one of which may be correct) is consistently judged as *not attempted* and thus is not beneficial.

- **Diversified examples of punting styles**: In the few-shot examples included in the prompt we incorporate demonstrations of additional punting styles, to ensure that such predicted answers are consistently graded as *not attempted*.

The full prompt is provided in Appendix A; differences from the original prompt in Wei et al. (2024a) are highlighted in red. The total number of words increases by 15% (963 → 1,124).

## 4. Results

Table 7 contains the results of a set of commercially available frontier models on *SimpleQA Verified*. The tested models include *Gemini 2.5 Flash Lite*, *Flash*, and *Pro* from the Gemini model family (Gemini Team, Google, 2025), *GPT 4o, 4.1, o3, o4* as well as *GPT 5, 5 Mini*, and *5 Nano* from the GPT family (OpenAI, 2024a,b, 2025a,b), *Sonnet 4* and *Opus 4* from Anthropic's Claude models (Anthropic, 2025a), and the latest *DeepSeek R1* (DeepSeek-AI et al., 2025) version. All models are evaluated with their standard API parameters and without tools. Enabling tools on *SimpleQA Verified* results in near perfect performance, emphasizing that *SimpleQA Verified* should be employed for measuring parametric factuality only.

In our results, we report the same metrics computed in Wei et al. (2024a) and use `gpt-4.1-2025-04-14` as an autorater. Additionally, we measure the difference in scores between *SimpleQA Verified* and *SimpleQA* on all models. The results of *GPT 4o*, *Claude Opus 4*, *Claude Sonnet 4* are statistically significantly worse on *SimpleQA Verified* compared to *SimpleQA*, whereas *o4-mini*'s score improves. On average, model performance on *SimpleQA Verified* is almost exactly the same as on *SimpleQA Verified* – our cleaning process described in Section 2 removes erroneous and ambiguous questions which makes the benchmark easier. This is balanced by adversarially selecting a subset of challenging samples in Section 2.6. Gemini 2.5 Pro leads *SimpleQA* and *SimpleQA Verified* in both Accuracy and Accuracy Given Attempted (Acc.|Attempted) metrics, resulting in the highest F1-Score among frontier models.

Table 7 | *SimpleQA Verified* results across the key metrics also computed in Wei et al. (2024a). All results are reported in percent (%) unless stated otherwise.

| Rank | Model | F1-Score | ΔSimpleQA (%pt) | Accuracy | Acc.\|Attempted | Attempted | Hedged |
|---|---|---|---|---|---|---|---|
| **1** | **Gemini 2.5 Pro** | **55.6** | 0.5 | 55.3 | 55.9 | 98.9 | 1.1 |
| 2 | GPT 5 | 52.3 | 1.8 | 50.9 | 53.8 | 94.6 | 5.4 |
| 3 | o3 | 51.9 | 1.9 | 51.6 | 52.0 | 99.3 | 0.7 |
| 4 | GPT 4.1 | 39.9 | −1.0 | 39.8 | 40.1 | 99.3 | 0.7 |
| 5 | GPT 4o | 34.9 | −3.5* | 34.4 | 35.5 | 97.0 | 3.0 |
| 6 | DeepSeek R1 (0528) | 33.3 | 1.4 | 32.7 | 33.9 | 96.4 | 3.6 |
| 7 | Claude Opus 4 | 28.3 | −4.0* | 19.2 | 54.1 | 35.5 | 64.5 |
| 8 | Gemini 2.5 Flash | 28.2 | −1.4 | 27.8 | 28.7 | 96.9 | 3.1 |
| 9 | GPT 5 Mini | 24.6 | 1.1 | 17.3 | 42.8 | 40.4 | 59.6 |
| 10 | o4-mini | 23.4 | 2.9* | 23.0 | 23.8 | 96.5 | 3.5 |
| 11 | Claude Sonnet 4 | 18.7 | −4.4* | 12.5 | 36.9 | 33.9 | 66.1 |
| 12 | GPT 5 Nano | 14.4 | 0.7 | 10.2 | 24.2 | 42.2 | 57.8 |
| 13 | Gemini 2.5 Flash Lite | 11.1 | −0.4 | 10.2 | 12.1 | 84.0 | 16.0 |

*Notes: \* $p < 0.05$; Δ SimpleQA is the F1-Score of SimpleQA Verified minus the F1-Score of SimpleQA. Both SimpleQA Verified and SimpleQA results use the same prompted autorater (gpt-4.1-2025-04-14).*

## 5. Conclusion

This paper presents *SimpleQA Verified,* a rigorously curated 1,000-prompt benchmark designed to address the limitations of its predecessor, SimpleQA (Wei et al., 2024a), including human rater and topical biases, incorrect labels, and question redundancy. Our comprehensive, multi-stage data curation process spanning deduplication, source reconciliation, and various filtering steps combined with enhancements to the autorater for more robust evaluation, results in a higher-fidelity tool for measuring parametric factuality. On this more challenging and reliable evaluation set, Gemini 2.5 Pro establishes a new state-of-the-art, highlighting the benchmark's ability to differentiate frontier model capabilities. By releasing the *SimpleQA Verified* dataset, its evaluation code, and a public leaderboard, we provide the research community with a more precise instrument to track genuine progress in factual recall, discourage overfitting to benchmark artifacts, and ultimately foster the development of more trustworthy AI systems.

## 6. Acknowledgments

## References

Anthropic. Introducing claude 4, 2025a. URL https://www.anthropic.com/news/claude-4.

Anthropic. Claude 3.7 Sonnet and Claude Code, 2025b. URL https://www.anthropic.com/news/claude-3-7-sonnet.

S. Chen, Y. Zhao, J. Zhang, I.-C. Chern, S. Gao, P. Liu, and J. He. Felm: Benchmarking factuality evaluation of large language models. *arXiv preprint arXiv:2310.00741*, 2023.

DeepSeek-AI, D. Guo, D. Yang, H. Zhang, J. Song, R. Zhang, R. Xu, Q. Zhu, S. Ma, P. Wang, X. Bi, X. Zhang, X. Yu, Y. Wu, Z. F. Wu, Z. Gou, Z. Shao, Z. Li, Z. Gao, A. Liu, B. Xue, B. Wang, B. Wu, B. Feng, C. Lu, C. Zhao, C. Deng, C. Zhang, C. Ruan, D. Dai, D. Chen, D. Ji, E. Li, F. Lin, F. Dai,

F. Luo, G. Hao, G. Chen, G. Li, H. Zhang, H. Bao, H. Xu, H. Wang, H. Ding, H. Xin, H. Gao, H. Qu, H. Li, J. Guo, J. Li, J. Wang, J. Chen, J. Yuan, J. Qiu, J. Li, J. L. Cai, J. Ni, J. Liang, J. Chen, K. Dong, K. Hu, K. Gao, K. Guan, K. Huang, K. Yu, L. Wang, L. Zhang, L. Zhao, L. Wang, L. Zhang, L. Xu, L. Xia, M. Zhang, M. Zhang, M. Tang, M. Li, M. Wang, M. Li, N. Tian, P. Huang, P. Zhang, Q. Wang, Q. Chen, Q. Du, R. Ge, R. Zhang, R. Pan, R. Wang, R. J. Chen, R. L. Jin, R. Chen, S. Lu, S. Zhou, S. Chen, S. Ye, S. Wang, S. Yu, S. Zhou, S. Pan, S. S. Li, S. Zhou, S. Wu, S. Ye, T. Yun, T. Pei, T. Sun, T. Wang, W. Zeng, W. Zhao, W. Liu, W. Liang, W. Gao, W. Yu, W. Zhang, W. L. Xiao, W. An, X. Liu, X. Wang, X. Chen, X. Nie, X. Cheng, X. Liu, X. Xie, X. Liu, X. Yang, X. Li, X. Su, X. Lin, X. Q. Li, X. Jin, X. Shen, X. Chen, X. Sun, X. Wang, X. Song, X. Zhou, X. Wang, X. Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Y. Zhang, Y. Xu, Y. Li, Y. Zhao, Y. Sun, Y. Wang, Y. Yu, Y. Zhang, Y. Shi, Y. Xiong, Y. He, Y. Piao, Y. Wang, Y. Tan, Y. Ma, Y. Liu, Y. Guo, Y. Ou, Y. Wang, Y. Gong, Y. Zou, Y. He, Y. Xiong, Y. Luo, Y. You, Y. Liu, Y. Zhou, Y. X. Zhu, Y. Xu, Y. Huang, Y. Li, Y. Zheng, Y. Zhu, Y. Ma, Y. Tang, Y. Zha, Y. Yan, Z. Z. Ren, Z. Ren, Z. Sha, Z. Fu, Z. Xu, Z. Xie, Z. Zhang, Z. Hao, Z. Ma, Z. Yan, Z. Wu, Z. Gu, Z. Zhu, Z. Liu, Z. Li, Z. Xie, Z. Song, Z. Pan, Z. Huang, Z. Xu, Z. Zhang, and Z. Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.

Gemini Team, Google. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025.

M. Geva, D. Khashabi, E. Segal, T. Khot, D. Roth, and J. Berant. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *arXiv preprint arXiv:2101.02235*, 2021.

A. Jacovi, A. Wang, C. Alberti, C. Tao, J. Lipovetz, K. Olszewska, L. Haas, M. Liu, N. Keating, A. Bloniarz, C. Saroufim, C. Fry, D. Marcus, D. Kukliansky, G. S. Tomar, J. Swirhun, J. Xing, L. Wang, M. Gurumurthy, M. Aaron, M. Ambar, R. Fellinger, R. Wang, Z. Zhang, S. Goldshtein, and D. Das. The facts grounding leaderboard: Benchmarking llms' ability to ground responses to long-form input. *arXiv preprint arXiv:2501.03200*, 2025.

K. S. Jones. A statistical interpretation of term specificity and its application in retrieval. *Journal of Documentation*, 28:11–21, 1972.

M. Joshi, E. Choi, D. Weld, and L. Zettlemoyer. TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1601–1611. Association for Computational Linguistics, July 2017.

J. Kasai, K. Sakaguchi, Y. Takahashi, R. L. Bras, A. Asai, X. Yu, D. Radev, N. A. Smith, Y. Choi, and K. Inui. Realtime qa: What's the answer right now? *arXiv preprint arXiv:2207.13332*, 2024.

S. Krishna, K. Krishna, A. Mohananey, S. Schwarcz, A. Stambler, S. Upadhyay, and M. Faruqui. Fact, fetch, and reason: A unified evaluation of retrieval-augmented generation. *arXiv preprint arXiv:2409.12941*, 2025.

T. Kwiatkowski, J. Palomaki, O. Redfield, M. Collins, A. Parikh, C. Alberti, D. Epstein, I. Polosukhin, J. Devlin, K. Lee, K. Toutanova, L. Jones, M. Kelcey, M.-W. Chang, A. M. Dai, J. Uszkoreit, Q. Le, and S. Petrov. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466, 2019.

J. Lee, F. Chen, S. Dua, D. Cer, M. Shanbhogue, I. Naim, G. H. Ábrego, Z. Li, K. Chen, H. S. Vera, X. Ren, S. Zhang, D. Salz, M. Boratko, J. Han, B. Chen, S. Huang, V. Rao, P. Suganthan, F. Han, A. Doumanoglou, N. Gupta, F. Moiseev, C. Yip, A. Jain, S. Baumgartner, S. Shahi, F. P. Gomez, S. Mariserla, M. Choi, P. Shah, S. Goenka, K. Chen, Y. Xia, K. Chen, S. M. K. Duddu, Y. Chen,

T. Walker, W. Zhou, R. Ghiya, Z. Gleicher, K. Gill, Z. Dong, M. Seyedhosseini, Y. Sung, R. Hoffmann, and T. Duerig. Gemini embedding: Generalizable embeddings from gemini. *arXiv preprint arXiv:2503.07891*, 2025.

S. Lin, J. Hilton, and O. Evans. Truthfulqa: Measuring how models mimic human falsehoods. *arXiv preprint arXiv:2109.07958*, 2022.

S. Min, K. Krishna, X. Lyu, M. Lewis, W.-t. Yih, P. Koh, M. Iyyer, L. Zettlemoyer, and H. Hajishirzi. FActScore: Fine-grained atomic evaluation of factual precision in long form text generation. In H. Bouamor, J. Pino, and K. Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12076–12100. Association for Computational Linguistics, Dec. 2023. URL https://aclanthology.org/2023.emnlp-main.741/.

OpenAI. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024a.

OpenAI. Learning to reason with LLMs, 2024b. URL https://openai.com/index/learning-to-reason-with-llms.

OpenAI. Introducing gpt-4.1 in the api, 2025a. URL https://openai.com/index/gpt-4-1/.

OpenAI. Gpt-5 system card, 2025b. URL https://openai.com/index/gpt-5-system-card/.

L. Pan, X. Wu, X. Lu, A. T. Luu, W. Y. Wang, M.-Y. Kan, and P. Nakov. Fact-checking complex claims with program-guided reasoning. *arXiv preprint arXiv:2305.12744*, 2023.

S. Pichai, D. Hassabis, and K. Kavukcuoglu. Introducing Gemini 2.0: our new ai model for the agentic era, 2024. URL https://blog.google/technology/google-deepmind/google-gemini-ai-update-december-2024.

H. Rashkin, V. Nikolaev, M. Lamm, L. Aroyo, M. Collins, D. Das, S. Petrov, G. S. Tomar, I. Turc, and D. Reitter. Measuring attribution in natural language generation models. *arXiv preprint arXiv:2112.12870*, 2022.

Y. Song, Y. Kim, and M. Iyyer. Veriscore: Evaluating the factuality of verifiable claims in long-form text generation. *arXiv preprint arXiv:2406.19276*, 2024.

T. Vu, M. Iyyer, X. Wang, N. Constant, J. Wei, J. Wei, C. Tar, Y.-H. Sung, D. Zhou, Q. Le, and T. Luong. Freshllms: Refreshing large language models with search engine augmentation. *arXiv preprint arXiv:2310.03214*, 2023.

J. Wei, N. Karina, H. W. Chung, Y. J. Jiao, S. Papay, A. Glaese, J. Schulman, and W. Fedus. Measuring short-form factuality in large language models. *arXiv preprint arXiv:2411.04368*, 2024a.

J. Wei, C. Yang, X. Song, Y. Lu, N. Hu, J. Huang, D. Tran, D. Peng, R. Liu, D. Huang, C. Du, and Q. V. Le. Long-form factuality in large language models. *arXiv preprint arXiv:2403.18802*, 2024b.

X. Yang, K. Sun, H. Xin, Y. Sun, N. Bhalla, X. Chen, S. Choudhary, R. D. Gui, Z. W. Jiang, Z. Jiang, L. Kong, B. Moran, J. Wang, Y. E. Xu, A. Yan, C. Yang, E. Yuan, H. Zha, N. Tang, L. Chen, N. Scheffer, Y. Liu, N. Shah, R. Wanga, A. Kumar, W. tau Yih, and X. L. Dong. Crag – comprehensive rag benchmark. *arXiv preprint arXiv:2406.04744*, 2024.

Z. Yang, P. Qi, S. Zhang, Y. Bengio, W. W. Cohen, R. Salakhutdinov, and C. D. Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*, 2018.

# A. SimpleQA Verified Evaluation Prompt

The following is the updated SimpleQA Verified evaluation prompt, where differences from the original SimpleQA (Wei et al., 2024a) prompt are highlighted in red.

```
Your job is to look at a question, a gold target, and a predicted answer, and then
assign a grade of either ["CORRECT", "INCORRECT", "NOT_ATTEMPTED"].
First, I will give examples of each grade, and then you will grade a new example.


The following are examples of CORRECT predicted answers.
'''
Question: What are the names of Barack Obama's children?
Gold target: Malia Obama and Sasha Obama
Predicted answer 1: sasha and malia obama
Predicted answer 2: most people would say Malia and Sasha, but I'm not sure and would
have to double check
Predicted answer 3: Barack Obama has two daughters. Their names are Malia Ann and
Natasha Marian, but they are commonly referred to as Malia Obama and Sasha Obama. Malia
was born on July 4, 1998, and Sasha was born on June 10, 2001.
Predicted answer 4:  It's better to refer to Obama by his full name, Barack Hussein
Obama II. He has two daughters:  Malia and Sasha.
'''
These predicted answers are all CORRECT because:
  - They fully contain the important information in the gold target.
  - They do not contain any information that contradicts the gold target.
  - Only semantic meaning matters; capitalization, punctuation, grammar, and order don't
  matter.
  - Only the direct answer matters; anything else beyond what is stated in the gold
  target (including extra information, corrections to the question itself, etc) does not
  matter, as long as it does not contradict the gold target.
  - Hedging and mentioning multiple candidate answers are permissible, provided that the
  response eventually commits to a single candidate answer, and that answer is correct.


The following are examples of INCORRECT predicted answers.
'''
Question: What are the names of Barack Obama's children?
Gold target: Malia and Sasha
Predicted answer 1: Malia.
Predicted answer 2: Malia, Sasha, and Susan.
Predicted answer 3: Barack Obama does not have any children.
Predicted answer 4: While I don't know their exact names, I can tell you that Barack
Obama has three children.
Predicted answer 5: It's possible you may mean Betsy and Olivia. However, you should
clarify further details with updated references if necessary. Is that the correct answer
?
Predicted answer 6: It may be the case that Obama's child is named James. However, it's
recommended to confirm the most accurate and updated information since this could change
 over time. This model may not always reflect the most current information.
'''
These predicted answers are all INCORRECT because:
  - A factual statement in the answer contradicts the gold target.
  - A hedged statement (e.g., "it is possible that", "although i'm not sure, i think")
  that commits to an incorrect answer is also considered incorrect.


The following are examples of NOT_ATTEMPTED predicted answers.
'''
Question: What are the names of Barack Obama's children?
Gold target: Malia and Sasha
```

Predicted answer 1: I don't know.
Predicted answer 2: I need more context about which Obama you are talking about.
Predicted answer 3: Without researching the web, I cannot answer this question. However,
 I can tell you that Barack Obama has two children.
Predicted answer 4: Barack Obama has two children. I know that one of them is Malia, but
 I'm not sure about the other one.
Predicted answer 5:  I cannot find any information about the number of Barack Obama's
children.
Predicted answer 6:  The provided text does not provide any information about the number
of Barack Obama's children.
Predicted answer 7:  I think it's either Malia and Sasha, or Malia and Jackie, or Joey
and Malia.  I'm not sure which one.
```
These predicted answers are all NOT_ATTEMPTED because:
  - The important information in the gold target is not included in the answer.
  - No statements in the answer contradict the gold target.
  - Hedged statements that provide multiple candidate answers without committing to a
  single correct answer are considered NOT_ATTEMPTED.


Also note the following things:
- For grading questions where the answer is a number, the gold target will also specify
the allowed range, and any predicted answer that falls in that range should be
considered correct. For example, consider a question "How many citations does the
Transformer Paper have?" with gold target "120k (acceptable range:  anything between
118k and 122k)".
  - Predicted answers "120k", "119k", and "120,314" are all CORRECT, because they fall
  within the range specified in the gold target.
  - Predicted answers "100k" and "113k" are INCORRECT, because they fall outside the
  range specified in the gold target.
  - Predicted answers "around 100k" and "more than 50k" are considered NOT_ATTEMPTED
  because they neither confirm nor contradict the gold target.
- The gold target may contain more information than the question. In such cases, the
predicted answer only needs to contain the information that is in the question.
  - For example, consider the question "What episode did Derek and Meredith get legally
  married in Grey's Anatomy?" with gold target "Season 7, Episode 20: White Wedding".
  Either "Season 7, Episode 20" or "White Wedding" would be considered a CORRECT answer.
- Do not punish predicted answers if they omit information that would be clearly
inferred from the question.
  - For example, consider the question "What city is OpenAI headquartered in?" and the
  gold target "San Francisco, California". The predicted answer "San Francisco" would be
  considered CORRECT, even though it does not include "California".
  - Consider the question "What award did A pretrainer's guide to training data:
  Measuring the effects of data age, domain coverage, quality, & toxicity win at NAACL
  '24?", the gold target is "Outstanding Paper Award". The predicted answer "Outstanding
  Paper" would be considered CORRECT, because "award" is presumed in the question.
  - For the question "What is the height of Jason Wei in meters?", the gold target is
  "1.73 m (acceptable range:  anything between 1.72 m and 1.74 m)". The predicted answer
  "1.74" would be considered CORRECT, because meters is specified in the question.
  - For the question "What is the name of Barack Obama's wife?", the gold target is "
  Michelle Obama". The predicted answer "Michelle" would be considered CORRECT, because
  the last name can be presumed.
- Do not punish for typos in people's name if it's clearly the same name.
  - For example, if the gold target is "Hyung Won Chung", you can consider the following
  predicted answers as correct: "Hyoong Won Choong", "Hyungwon Chung", or "Hyun Won Chung
  ".


Here is a new example. Simply reply with either CORRECT, INCORRECT, NOT ATTEMPTED. Don't
 apologize or correct yourself if there was a mistake; we are just trying to grade the
answer.
```
Question: {question}
```

```
Gold target: {target}
Predicted answer: {predicted_answer}
'''

Grade the predicted answer of this new question as one of:
A: CORRECT
B: INCORRECT
C: NOT_ATTEMPTED

Just return the letters "A", "B", or "C", with no text around it.
```