

RARR: Researching and Revising What Language Models Say, Using Language Models

Luyu Gao¹*, Zhuyun Dai²*, Panupong Pasupat²*, Anthony Chen³*,
Arun Tejasvi Chaganty²*, Yicheng Fan²*, Vincent Y. Zhao², Ni Lao²,
Hongrae Lee², Da-Cheng Juan², Kelvin Guu²*

¹Carnegie Mellon University, ²Google Research, ³UC Irvine

luyug@cs.cmu.edu anthony.chen@uci.edu

{zhuyundai, ppasupat, arunchaganty, yichengfan, vzhao, nlao, hrlee, dacheng, kguu}@google.com

Abstract

Language models (LMs) now excel at many tasks such as question answering, reasoning, and dialog. However, they sometimes generate unsupported or misleading content. A user cannot easily determine whether their outputs are trustworthy or not, because most LMs do not have any built-in mechanism for *attribution* to external evidence. To enable attribution while still preserving all the powerful advantages of recent generation models, we propose RARR (*Retrofit Attribution using Research and Revision*), a system that 1) automatically finds attribution for the output of any text generation model, and 2) post-edits the output to fix unsupported content while preserving the original output as much as possible. When applied to the output of several state-of-the-art LMs on a diverse set of generation tasks, we find that RARR significantly improves attribution while otherwise preserving the original input to a much greater degree than previously explored edit models. Furthermore, the implementation of RARR requires only a handful of training examples, a large language model, and standard web search.¹

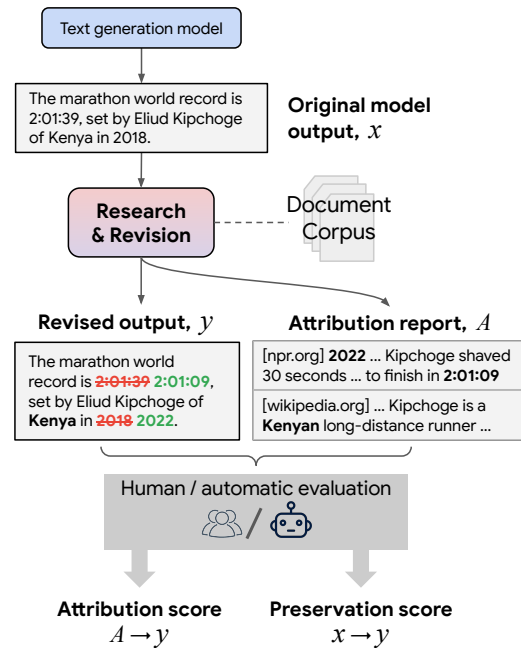


Figure 1: **The Editing for Attribution task.** The input x is a text passage produced by a generation model. Our *Research & Revision* model outputs an attribution report A containing retrieved evidence snippets, along with a revision y whose content can be *attributed* to the evidence in A while *preserving* other properties of x such as style or structure.

1 Introduction

Generative language models (LMs) and other text generation models are now the backbone of many AI systems. For example, large language models can perform multi-step reasoning (Nye et al., 2021; Wei et al., 2022), generate plans (Ahn et al., 2022), use tools and APIs (Shin et al., 2021; Thoppilan et al., 2022), and answer open-domain questions (Petroni et al., 2019; Roberts et al., 2020).

Despite these incredible advances, state-of-the-art LMs still frequently produce biased, misleading,

or unsupported content, colloquially called “hallucinations” (Maynez et al., 2020; Menick et al., 2022). To make LMs more trustworthy, we want to justify each generation by an *attribution report* (Rashkin et al., 2021; Bohnet et al., 2022) that contains supporting evidence from trusted sources (e.g., encyclopedia or articles) where appropriate.

Most existing LMs, such as those based on sequence-to-sequence architectures, lack a built-in mechanism for attribution. Even *retrieval-augmented models* (Guu et al., 2020; Lewis et al., 2020), which retrieve relevant documents and then condition on them to generate text, still do not guarantee attribution. Prior work has shown that

*Lead contributors. Please see Contributions section for details. °Work done during an internship at Google Research.

¹We release open-source implementations of RARR, the evaluation pipeline, and the evaluation sets at <https://github.com/anthonywchen/RARR>.

retrieval-augmented models generate text that either includes additional information outside the retrieved documents (Dziri et al., 2022), ignores the documents altogether (Krishna et al., 2021), or even contradicts the documents (Longpre et al., 2021). In fact, occasionally ignoring the retrievals can make the models more robust to bad retrievals (Khandelwal et al., 2020), illustrating that end-task performance and attribution are not always aligned.

Instead of constraining LMs to generate attributed text, we propose a model-agnostic approach to improve the attribution of any existing LM: *Retrofit Attribution using Research and Revision* (RARR). The approach is inspired by works on fact-checking² where simple research-and-revise workflows are effective at attributing or correcting unattributed claims made by humans (Thorne et al., 2018; Schuster et al., 2021; Thorne and Vlachos, 2021). As shown in Figure 1, after generating text with the LM, RARR does *research* to retrieve relevant evidence, and then *revises* the text to make it consistent with the evidence while preserving qualities like style or structure, enabling the revised text to be seamlessly used in place of the original. RARR can be viewed as a retrieval-augmented model where retrieval happens *after* generation rather than before. This allows RARR to stand on the shoulders of giant LMs without having to modify them to support attribution.

In our effort to expand the scope of Research & Revision models to handle the output of arbitrary LMs, we make the following contributions. First, we formalize the *Editing for Attribution* task and propose new **metrics** that evaluate revision models not just on their ability to produce well-attributed revisions, but also on their ability to otherwise *preserve* original properties of the text. Second, we use these metrics to **benchmark** how existing revision models perform on various types of LM outputs such as knowledge-intensive statements, reasoning chains, and dialog responses. Finally, we find that existing revision models do not always generalize across many tasks (and were not originally intended to), and therefore propose a new research-and-revise **model** that leverages the power of few-shot prompting in large language models to robustly generalize across domains.

²In this paper, we generally avoid the term “fact-checking” other than to reference relevant literature, because we only address attribution, and attribution does not entail correctness. Even if a claim is attributed to a particular source, it does not guarantee that the source is “correct” (Menick et al., 2022).

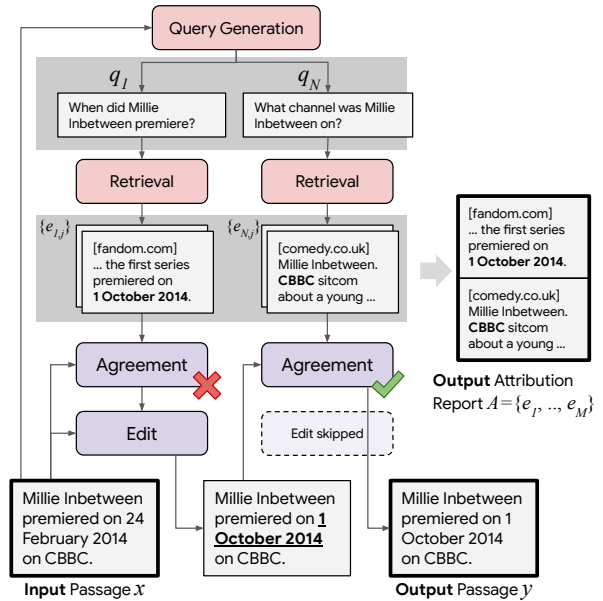


Figure 2: An overview of RARR, which improves attribution for a text passage via *Research & Revision*. Given the input text passage, the **research** stage uses a *query generator* to raise questions about different aspects of the text. The *retriever* then searches for evidence to investigate each query. The **revision** stage first runs an *agreement model* to detect disagreement between the text and the evidence, then runs an *edit model* to revise the text if needed. Finally, M evidence snippets are selected to form an attribution report.

2 Task formulation

We propose the task of *Editing for Attribution* as follows. As Figure 1 shows, the input to the system is a text passage x produced by a generation model. The output is a revised text passage y along with an *attribution report* A , which contains *evidence snippets* e_1, \dots, e_M that support the content in y . Optionally, the attribution report can contain additional information such as the alignment between evidence snippets and relevant parts in y .

We propose to measure the quality of the revised text y and attribution report A along two dimensions: (1) **attribution**: how much of the revised text y can be attributed to the evidence in A , and (2) **preservation**: how much the revised text y preserves aspects of the original text x .

2.1 Measuring attribution

Previously, Rashkin et al. (2021) proposed *Attributable to Identified Sources* (AIS), a human evaluation framework which considers a binary notion of attribution. Roughly speaking, a text passage y is attributable to a set A of evidence if a

generic hearer would affirm the statement “According to A , y ” under the context of y . A system either receives full credit (1.0) if *all* content in y can be attributed to A , and no credit (0.0) otherwise.

We propose a more fine-grained, sentence-level extension of AIS. We ask annotators to give an AIS score for each sentence s of y , and then report the average AIS score across all sentences:

$$\text{Attr}_{\text{AIS}}(y, A) = \text{avg}_{s \in y} \text{AIS}(s, A). \quad (1)$$

Since the AIS score is binary, this effectively measures the percentage of sentences in y that are fully attributed to A . When judging each sentence, we also give annotators access to the surrounding sentences and other necessary context, such as the question that the text passage responded to. We also impose the maximum number of evidence snippets in the attribution report A to make it concise enough for both the annotator and downstream users. By manually inspecting 30 examples from our benchmarks, we found $M = 5$ snippets to be sufficient for full attribution.

During model development, we define an automated metric, auto-AIS ($\text{Attr}_{\text{auto}}$), that approximates human AIS judgments. We utilize the natural language inference (NLI) model from [Honovich et al. \(2022\)](#), which correlates well with AIS scores. For each sentence s of y , and for each evidence snippet e in A , let $\text{NLI}(e, s)$ be the model probability of e entailing s . We then define

$$\text{Attr}_{\text{auto}}(y, A) = \text{avg}_{s \in y} \max_{e \in A} \text{NLI}(e, s). \quad (2)$$

To improve accuracy, we decontextualize ([Choi et al., 2021](#)) each sentence based on the entire context of y before computing the scores. See Appendix B for implementation details.

2.2 Measuring preservation

To measure preservation, we first ask annotators to decide if the revision preserves the text’s original intent (completely, somewhat, or not at all — see Appendix C for exact rubrics). Like AIS evaluation, we give annotators the necessary surrounding context. We define the binary metric $\text{Pres}_{\text{intent}}(x, y)$ to be 1.0 if the revision completely preserves the original intent, and 0.0 otherwise.

However, even if a revision preserves intent, it may still make superfluous modifications, such as reordering words, changing textual style, or including unnecessary additional information ([Thorne](#)

and [Vlachos, 2021](#)). Different tasks have different requirements for what should be preserved. Here, we desire a simple metric that can be readily computed for many tasks and that generally penalizes unnecessary changes. We thus define a metric based on the character-level Levenshtein edit distance ([Levenshtein, 1965](#)) between x and y :

$$\text{Pres}_{\text{Lev}}(x, y) = \max \left(1 - \frac{\text{Lev}(x, y)}{\text{length}(x)}, 0 \right) \quad (3)$$

This metric is 1.0 if x and y are the same, and 0.0 if y completely overwrites all parts of x . Pres_{Lev} is generally sensitive to any kind of change, but certainly does not capture all notions of preservation (e.g., preserving rhyme schemes or puns).

We want the revision to preserve the original intent while avoiding superfluous edits. To reflect this, we finally combine the two metrics as

$$\text{Pres}_{\text{comb}}(x, y) = \text{Pres}_{\text{intent}}(x, y) \cdot \text{Pres}_{\text{Lev}}(x, y). \quad (4)$$

which is 0.0 if the revision changes the intent and equal to $\text{Pres}_{\text{Lev}}(x, y)$ otherwise. Since $\text{Pres}_{\text{intent}}$ requires human annotation, we use Pres_{Lev} as an automated metric for model development.

2.3 Discussion

Optimizing for attribution alone cannot ensure a good revision: for example, an adversarial editor could ensure 100% attribution by simply replacing the input x with the text of any arbitrary retrieved document, which is trivially attributable to itself. Ideally, we want to maximize both attribution and preservation, while navigating any tradeoffs between the two. In our experiments, we report both metrics, as well as their harmonic mean (F1_{AP} , analogous to how recall and precision are combined in F1).

We emphasize that this evaluation scheme does not require any “gold” or “reference” edits (unlike many prior evaluations of text revision models), which are often only available for specialized domains. This enables us to broaden the scope to a much wider range of generation tasks.

3 Approach

We now present *Retrofit Attribution using Research and Revision* (RARR), a simple method for solving the *Editing for Attribution* task. As illustrated in Figure 2, given an input passage x , the research stage first generates a set of queries $\{q_1, \dots, q_N\}$, each investigating one aspect of x that potentially

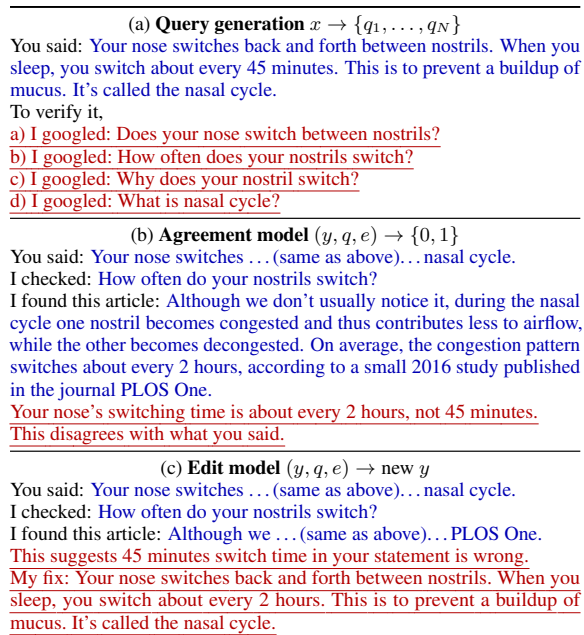


Figure 3: **Examples of few-shot examples** used to prompt the PaLM model (blue = input; red = output).

requires attribution. For each query q_i , it retrieves web documents and selects the best evidence snippets $\{e_{i1}, e_{i2}, \dots\}$. The revision stage then revises the original text x using the retrieval results $\{(q_1, e_{11}), \dots\}$, yielding a revised text y .

Most components for RARR are implemented using few-shot prompting (Brown et al., 2020). We use PaLM (Chowdhery et al., 2022) as our language model. Figure 3 shows some few-shot examples we use, while Appendix D lists the full prompts.

3.1 Research stage

Query generation We perform *comprehensive question generation* (CQGen) which produces a sequence of questions covering *all aspects* of the passage x that need to be verified and attributed. A similar strategy has been employed to train text-planning models (Narayan et al., 2022). A prompt with six human demonstrations was sufficient for PaLM to adequately learn the task. To increase diversity and coverage, we sample from our CQ-Gen model three times and take the union of the resulting queries.

Evidence retrieval For each query from CQGen, we use Google Search to retrieve $K = 5$ web pages. We extract candidate evidence snippets from each web page by running a sliding window of four sentences across the page, breaking at document headings. The evidence snippets for each query

are then ranked based on their relevance to the query. For this, we use an existing query-document relevance model trained following Ni et al. (2021), which computes a relevance score $S_{\text{relevance}}(q, e)$ between a query q and an evidence snippet e . We then keep the top $J = 1$ evidence for each query. The final retrieval result is $[(q_1, e_{11}), \dots, (q_1, e_{1J}), \dots, (q_N, e_{N1}), \dots, (q_N, e_{NJ})]$, where e_{ij} denotes the j^{th} evidence for the i^{th} query, and N denotes the total number of queries from CQGen (which can be different for each input x).

3.2 Revision stage

After retrieving evidence, certain parts of x may now be properly attributed, but other parts remain unattributed and should be revised. As illustrated in Figure 2, the revision stage initializes the output $y = x$. Then for each retrieved $(q, e) = (q_i, e_{ij})$, the *agreement model* checks if the evidence e disagrees with the current output y regarding the issue in query q . If a disagreement is detected, the *edit model* edits y to agree with e ; otherwise, it does nothing. The process continues until all retrievals are processed.

Agreement model The agreement model takes the partially edited passage y , a query q , and the evidence e as input. It then decides whether both y and e imply the same answer to the question in q . This form of question-guided agreement was previously explored by Honovich et al. (2021). We implement this by few-shot prompting PaLM using a chain-of-thought style prompt (Wei et al., 2022), where we ask the model to explicitly state the implied answers for both y and e before producing its judgment about their agreement.

Edit model The edit model is run only if a disagreement is detected. The model takes y , q and e as input, and outputs a new version of y that aims to agree with e while otherwise minimally altering y . We again use few-shot prompting and chain-of-thought, where we ask the model to first identify a particular span in y that needs to be edited before generating the revised y . This helps reduce the editor's deviation from the current y .³

³The editor occasionally produces large edits that bring the new revision close to e but far from the current y . Since this is rarely desirable, we reject edits with edit distance above 50 characters or 0.5 times the original text length.

3.3 Attribution report

Finally, we select at most $M = 5$ evidence snippets to form an attribution report A . Note that during evidence retrieval and revision, we may have encountered and used more than M snippets. Our goal is to find a subset of snippets that maximizes *coverage* over the potentially attributable points in the passage, as represented by the queries q_1, \dots, q_N . We use the relevance model from Section 3.1 as a proxy for measuring how much an evidence e covers the point raised by a query q . Then, we exhaustively search for $A \subseteq \{e_{11}, \dots, e_{NJ}\}$ of size at most M that maximizes

$$\text{Cover}(A, q_{1:N}) := \sum_{i=1}^N \max_{e \in A} S_{\text{relevance}}(q_i, e). \quad (5)$$

4 Related work

Fact-checking Our research builds upon works to identify whether a claim is supported or refuted by the given evidence (Thorne et al., 2018; Wang, 2017; Karadzhov et al., 2017; Augenstein et al., 2019; Wadden et al., 2020). In real-world scenarios such as the one which RARR operates in, relevant evidence may not be provided, necessitating retrieval (Fan et al., 2020; Piktus et al., 2021).

Post-hoc editing for factuality Recent work has gone beyond checking the validity of a claim to correcting a piece of text to be factually consistent with a set of evidence via post-hoc editing (Shah et al., 2020; Thorne and Vlachos, 2021; Schuster et al., 2021; Balachandran et al., 2022; Cao et al., 2020; Iso et al., 2020). FRUIT (Logan IV et al., 2022) and PEER (Schick et al., 2022) both implement an editor that is fine-tuned on Wikipedia edit history with the goal of updating outdated information and collaborative writing respectively. Evidence-based Factual Error Correction (EFEC; Thorne and Vlachos, 2021) also implements a full research-and-revise workflow trained on Wikipedia passages (Thorne et al., 2018). A key differentiator of RARR is its ability to edit the output of any generation model without being restricted by the domain, task, or the need for training data.

Measuring attribution A key part of improving attribution is being able to quantify it. Apart from human evaluation (Rashkin et al., 2021), several automated evaluation methods have been proposed. Our work uses an entailment-based metric, which measures whether the referenced evidence entails

PaLM outputs on NQ (factoid statements)
Millie Inbetween is a British comedy television series. It premiered on 24 February 2014 on BBC One. The first series was produced by John Yorke and Phil Clymer.
PaLM outputs on SQA (reasoning chains)
The highest point of Mount Wycheproof is 70 metres. Edmund Hillary climbed Mount Everest, which is 8,848 metres. So Mount Wycheproof would be a breeze for Edmund Hillary.
LaMDA outputs on QReCC (knowledge-intensive dialogs)
When was Welsh social reformer Robert Owen born? Robert Owen was born on 14 May 1771 ... Did he have another job? In 1810 he moved to Manchester and established a draper's shop.

Figure 4: **Examples of input passages.** For QReCC, prior dialog turns are also given as the context.

the output text (Bohnet et al., 2022; Kryscinski et al., 2020; Goyal and Durrett, 2021). A common alternative is to evaluate whether the output text contains the same factual information as the evidence; e.g., by checking if both yield the same answer to the same question (Wang et al., 2020). We use this notion of attribution in RARR’s agreement model rather than for evaluation.

Retrieval-augmented models Models with a retrieval component have seen successes in question answering (Chen et al., 2017; Lee et al., 2019; Nakano et al., 2021), machine translation (Zhang et al., 2018), code generation (Hayati et al., 2018), language modeling (Khandelwal et al., 2020), and other knowledge-intensive tasks (Lewis et al., 2020). Their retrievals are not necessarily attributions (Dziri et al., 2022; Longpre et al., 2021) and typically are not used to revise an existing output. An exception is LaMDA (Thoppilan et al., 2022), a language model for dialog that performs revision by training on human annotations.

5 Experiments

5.1 Evaluation setups

RARR aspires to be a general-purpose method for improving the attribution of any text generation model in any text domain. We thus construct evaluation benchmarks by taking the task input from three diverse datasets, and prompting different generation models to produce *long-form outputs* which may contain “hallucinations,” as demonstrated in Figure 4. These long-form outputs serve as input text passages to RARR. We generate 150 development and 150 test passages for each combination of generation model and source dataset.

Factoid statements We prompt PaLM 540B and GPT-3 text-davinci-002 to generate long-form answers to questions from the Natural Questions dev set (NQ; Kwiatkowski et al., 2019). The resulting passages are mostly coherent but often contain factual errors. This setup examines the ability to attribute a diverse range of factoid knowledge.

Reasoning chains Language models can generate reasoning chains to answer complex questions (Wei et al., 2022). We use PaLM and GPT-3 to generate reasoning chains for the StrategyQA train set (SQA; Geva et al., 2021). This setup tests whether the revision model can provide better attribution for intermediate steps of reasoning, while preserving the overall reasoning process.

Knowledge-intensive dialogs We consider the conversational QA task from the QReCC dev set (Anantha et al., 2021). Given the previous dialog turns, which are rounds of questions and answers $(Q_1, A_1, Q_2, A_2, \dots, Q_k)$, we use LaMDA and GPT-3 to answer to the final question Q_k conditioned on the dialog history. The answer tends to be context-dependent, featuring pronouns and implicit references. All dialog turns are given alongside the answer as inputs to the revision model.

5.2 Models

We compare RARR to several systems that have a research-and-revise workflow.

EFEC We consider EFEC (Thorne and Vlachos, 2021) as a representative fine-tuned editor. EFEC fine-tunes a T5-based model to revise text conditioned on multiple evidence snippets using both semi-supervised and fully-supervised approaches. We compare against their fully-supervised approach, which performed best in their experiments. EFEC uses a neural retrieval model (Karpukhin et al., 2020) to retrieve from Wikipedia; however, not all passages in our experiments are supported by Wikipedia articles. To more fairly compare the editing capabilities of EFEC, we instead use the evidence retrieved by our research stages (CQGen and web search). Note that the EFEC editor conditions on multiple pieces of evidence at once, while our editor iteratively conditions on one at a time.

LaMDA LaMDA (Thoppilan et al., 2022) generates responses in three steps: 1) generate a “base response”; 2) generate search queries from the base response; 3) generate a “revised response” conditioned on the base response and retrieved evidence.

Model	Attribution		Preservation			
	auto-AIS	AIS	intent	Lev	comb	F1 _{AP}
PaLM outputs on NQ						
EFEC	45.6 → 64.3	35.4 → 48.3	16.0	39.1	10.4	17.1
LaMDA	39.5 → 49.9	18.3 → 30.4	26.0	39.6	21.1	24.9
RARR	45.6 → 54.9	35.4 → 43.4	90.0	89.6	83.1	57.0
PaLM outputs on SQA						
EFEC	37.8 → 58.6	24.5 → 51.7	6.0	31.0	3.8	7.1
LaMDA	32.7 → 43.2	15.8 → 27.0	40.0	46.4	33.7	30.0
RARR	37.6 → 45.1	24.5 → 31.5	92.6	89.9	84.6	45.9
LaMDA outputs on QReCC						
EFEC	19.1 → 47.4	13.2 → 48.7	39.7	39.4	23.7	31.9
LaMDA	16.4 → 36.2	16.0 → 27.1	21.3	24.8	12.0	16.6
RARR	18.8 → 29.4	13.2 → 28.3	95.6	80.2	78.1	41.5

Table 1: **Evaluation results.** For attribution, we report the AIS scores of the texts both before and after editing (before → after). For preservation, we report intent preservation $\text{Pres}_{\text{intent}}$, Levenshtein similarity Pres_{Lev} , and the combined $\text{Pres}_{\text{comb}}$. We summarize Attr_{AIS} and $\text{Pres}_{\text{comb}}$ using their harmonic mean (F1_{AP}).

To apply LaMDA on a given text x , we simply set the base response in step 1 to x , and then run steps 2 and 3 (we call these latter two stages “LaMDA Research”). LaMDA was trained as a dialog system, and always expects a dialog context where the user speaks first. So, for non-dialog tasks, we insert an artificial user utterance as dialog history: “*Tell me something interesting.*” For the attribution report, we take all evidence documents retrieved by LaMDA during its research process.

RARR Our model uses few-shot prompting on PaLM 540B for query generation, the agreement model, and the edit model. We use the same prompts for all tasks except when the context comes from a dialog, where we slightly modify the prompts to use the dialog context (e.g., CQGen now maps dialog context + x to queries). The query-evidence relevance model $S_{\text{relevance}}$ is a pretrained T5-large model (Raffel et al., 2020) fine-tuned following Ni et al. (2021) on MS MARCO (Nguyen et al., 2016). See Appendix D for the few-shot prompting strategies and more modeling details.

5.3 Results

For the main experiments, we report results on passages generated by PaLM and LaMDA. Results on GPT-3 passages show similar trends (Appendix A). Table 1 and Figure 5 show attribution and preservation results for each model and dataset. We also report F1_{AP} , the harmonic mean of the two metrics, which is shown as level curves in Figure 5.

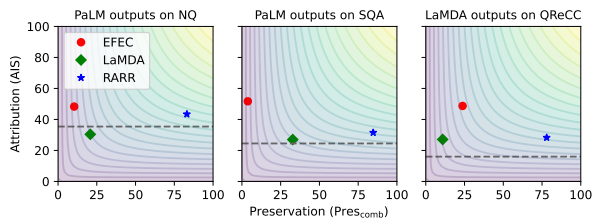


Figure 5: **Attribution and preservation scores.** Dashed lines indicate the highest attribution score obtained by any of the models *before* editing: points above the line have better attribution after revision. The contours are $F1_{AP}$ level curves: points along a contour have equivalent $F1_{AP}$. Different models make very different trade-offs between attribution and preservation. Only RARR has a robust $F1_{AP}$ across all tasks.

RARR significantly improves attribution while preserving most of the original text. In terms of $F1_{AP}$, RARR is the only method that performs robustly across all three datasets, and significantly outperforms prior methods on NQ and SQA.

We found that RARR is the only method that preserves the original intent of x over 90% of the time — EFEC and LaMDA only manage to preserve the original intent 6–40% of the time. We also see that editing is crucial to improve attribution: if we only retrieve evidence to support the original response x without editing, attribution ranges from the low 10s to mid 30s. After editing, RARR can increase attribution by up to 13% absolute, while changing only 10–20% of the text.

As noted in Section 2, one can sacrifice preservation for higher attribution. EFEC is able to obtain strong $F1_{AP}$ on QReCC by making larger changes to the text in exchange for a higher attribution score. However, it occupies a very different point from RARR on the attribution-preservation trade-off curve, as visualized in Figure 5.

6 Analysis

6.1 Qualitative analysis

Human oracle To understand the remaining headroom in our task, we ask: *what is the minimal amount of editing needed to make a text passage fully attributed?* The answer would depend on the quality of the LM that generated the text as well as the task difficulty. As an approximation, we manually edited 30 examples in our NQ benchmark until we judged them to be 100% attributable. We achieved a preservation score of 88%, which (when combined with 100% attribution) translates to 93.6 $F1_{AP}$, indicating a significant headroom.

x : Justice Ashok Kumar Mathur headed the 7th central pay commission in India. It was created in 2014 and submitted its report in 2016 .	Attribution: 50%	Preservation: 100%
EFEC: The 7th central pay commission in India was created in 2014.	Attribution: 100%	Preservation: 0%
LaMDA: I heard the 7th CPC made recommendations for increasing the minimum salary pay from Rs 7066 to 18k per month for new central government employees.	Attribution: 0%	Preservation: 0%
RARR: Justice Ashok Kumar Mathur headed the 7th central pay commission in India. It was created in 2014 and submitted its report in 2015 .	Attribution: 100%	Preservation: 100%
evidence: The 7th Central Pay Commission (Chair: Justice A. K. Mathur) submitted its report on November 19, 2015 . The Commission had been appointed in February 2014 , to look at remuneration for central government employees. ...		

Figure 6: **Example model outputs and human judgment of their attribution and preservation scores.** EFEC reduces the passage x into a single sentence. LaMDA changes the writing style. RARR preserves the structure of the input passage. We show one evidence retrieved by RARR to help explain the example.

Analyzing the baselines As exemplified in Figure 6, EFEC frequently attempts to summarize the entire passage into one sentence, or drops later sentences. This is likely due to EFEC’s training data, which was limited to single sentences. This behavior generally increases the attribution score, because it is usually easier to make one sentence fully attributable than many sentences. However, in datasets where the claim contains multiple sentences (NQ and SQA), such a behavior yields low preservation scores, and also results in outputs that are less informative. We expect that EFEC could perform much better if its training data were augmented to include multiple sentences. LaMDA Research achieves similar attribution scores to RARR. But as mentioned in Section 5.2, the intent and linguistic style of the output tend to deviate from the input, resulting in lower preservation scores (Figure 6). We emphasize that this is not a purely apples-to-apples comparison since LaMDA was not optimized for preservation. Overall, these experiments are mainly meant to illustrate that prior models were simply not designed for the task of *Editing for Attribution*, rather than to mark RARR as the best method.

Analyzing RARR For the research stage, the question generation model had comprehensive coverage: a manual inspection of 40 examples shows $> 80\%$ with questions that fully cover all aspects of the input text. The retriever was strongest at re-

(a) Correctly revising an entity	
<i>y</i> :	If She Knew What She Wants was written by Henry Roth .
<i>e</i> :	[en.wikipedia.org] "If She Knew What She Wants" is a song written by American singer-songwriter Jules Shear and introduced on ...
<i>y'</i> :	If She Knew What She Wants was written by Jules Shear .
(b) Correctly revising a number	
<i>y</i> :	God Save the Queen became the British national anthem in 1745
<i>e</i> :	[www.britannica.com] The oldest national anthem is Great Britain's "God Save the Queen," which was described as a national anthem in 1825
<i>y'</i> :	God Save the Queen became the British national anthem in 1825 (The year 1745 was when the song was first performed.)
(c) Performing a necessary larger revision	
<i>y</i> :	"It's My Party" is a song written and composed by American singer-songwriter and producer Walter Gold . The song was published in 1962 and first recorded by American singer-songwriter, producer and arranger Lesley Gore.
<i>e</i> :	[secondhandsongs.com] Song: It's My Party written by John Gluck, Wally Gold, Seymour Gottlieb, Herbert Weiner
<i>y'</i> :	"It's My Party" is a song written and composed by John Gluck, Wally Gold, Seymour Gottlieb and Herb Weiner . The song was published in 1963 and first recorded by American singer Lesley Gore.
(d) Confused by misleading evidence	
<i>y</i> :	In the Heat of the Night is a police procedural crime drama TV series produced by MGM Television. It aired on NBC from September 8, 1988 to May 3, 1992. ...
<i>e</i> :	[metv.com] Watch In the Heat of the Night on MeTV .
<i>y'</i> :	In the Heat of the Night is a police procedural crime drama TV series produced by MGM Television. It aired on MeTV from September 8, 1988 to May 3, 1992. ... (The show did actually air on NBC in the 1980s, but the model fails to realize that MeTV is a re-run much later in time.)
(e) Improved attribution, but reasoning doesn't adapt to new info	
	(<i>Q</i> : Does Homer Simpson need two hands worth of fingers to count to 5?)
<i>y</i> :	Homer Simpson has 5 fingers and 2 hands. Therefore, he does not need two hands worth of fingers to count to 5. He only need one hand to count to 5.
<i>e</i> :	[refinery29.com] Other iconic animated characters, such as Homer Simpson and his brood, ... only have four fingers , too.
<i>y'</i> :	Homer Simpson has 4 fingers and 2 hands. Therefore, he does not need two hands worth of fingers to count to 5. He only need one hand to count to 5. (Having only 4 fingers per hand, he does need two hands to count to 5.)

Figure 7: **Example revisions from RARR, both good and bad.** *y* = partially edited passage; *e* = evidence; *y'* = passage after editing with *e*.

searching content involving distinct entities (e.g., a movie, a major event, or a person). In contrast, we found significant headroom for better attribution of statements involving generic objects and more abstract claims (e.g. "Video games require electricity."— since this is obvious to most humans, retrieved articles from the web tend to address related but different topics). We suspect that a significant amount of attribution headroom on our benchmarks would benefit from a better research stage.

For the revision stage, RARR was able to revise many unattributed claims, especially those involving entities and numbers (Figures 7a and 7b). It can also perform larger revisions when necessary (Figure 7c). Moreover, RARR abstains from editing when the claim is already well-attributed: on NQ, among the inputs with near-perfect attribution (pre-edit $\text{Attr}_{\text{AIS}} > 0.9$), RARR does not make an edit in 90% of the cases. However, the system also has several shortcomings. Some erroneous edits arise from misleading irrelevant evidence (Figure 7d). We also observed an interesting challenge when revising reasoning chains, where the model suc-

cessfully revised an incorrect claim, but did not revise subsequent reasoning steps that depend on the earlier claim (Figure 7e). In this case, further editing to improve logical coherence could help.

6.2 Ablations

Ablating query generation RARR uses generated questions as search queries for evidence retrieval. We consider two natural alternatives: using the entire input passage as a single search query, or using each sentence as a search query. For the former, we retrieve $J = 3$ evidence snippets to make the amount a closer match to other methods.

The results are in Table 2. Using the entire input passage as the query gives poor results, as the retrieved evidence tends to not focus on potentially unattributed parts in the passage. Using sentences as queries gives results closer to the full CQGen, but a closer analysis reveals two caveats.

First, sentences-as-queries are more effective when such sentences "mimic" content on the Web, and are less effective otherwise. In Table 3, we test this by excluding all of Wikipedia from web search results (since many PaLM outputs for NQ have a Wikipedia style). The attribution performance of sentences-as-queries drops significantly, while CQGen is more robust.

Second, sentence-as-queries tends to retrieve passages that may encourage confirmation bias. Consider the example "Georgia is called the Peach State, but California actually produces the most peaches." Retrieval using sentences-as-queries found an article echoing that California produces the most peaches, while CQGen generated the more impartial query "Which state produces the most peaches?" and found a newer article saying that South Carolina replaced California as the top peach producer. In this case, RARR using CQGen needs to sacrifice more preservation score to edit the text, leading to a lower F1_{AP} score. This underscores that attribution alone cannot measure "correctness" since not all evidence is up-to-date or reliable.

Ablating agreement model We try removing the agreement model, which effectively forces the model to revise the passage based on every retrieved evidence. The results are shown in Table 2. As expected, more revision leads to less preservation score and spurious changes to the text passage, as demonstrated in Figure 8.

Impact on downstream task performance We have measured preservation using the metric de-

Model	PaLM outputs on NQ			PaLM outputs on SQA			LaMDA outputs on QReCC		
	Attr _{auto}	Pres _{Lev}	F1 _{AP}	Attr _{auto}	Pres _{Lev}	F1 _{AP}	Attr _{auto}	Pres _{Lev}	F1 _{AP}
Full RARR	45.6 → 54.9	89.6	68.1	37.6 → 45.1	89.9	60.0	18.8 → 29.4	80.2	43.1
no agreement model	45.6 → 50.6	82.6	62.8	37.8 → 46.9	83.4	60.0	18.8 → 28.8	72.0	41.2
query = input	45.4 → 47.2	98.4	63.8	39.4 → 30.3	98.8	46.4	19.7 → 20.6	96.3	34.0
query = sentence	49.1 → 52.1	97.0	67.8	43.7 → 44.3	98.8	61.2	19.0 → 19.6	97.0	32.6

Table 2: **Ablation results.** We report the automatic metrics: Attr_{auto}, Pres_{Lev}, and harmonic mean between the two (F1_{AP}). We show auto-AIS scores both before and after editing (before → edit), with respect to the attribution report *A* produced by the model. Even though sentence-as-queries may achieve similar F1_{AP} as RARR, it is less robust to corpus shifts and tends to retrieve passages that may encourage confirmation bias.

Model	NQ F1 _{AP}		SQA F1 _{AP}	
	orig	no wiki	orig	no wiki
Full RARR	68.1	64.3	60.0	57.6
query = sentence	67.8	60.3	61.2	56.7

Table 3: **The impact of excluding Wikipedia from the retrieval corpus.** CQGen (full RARR) is more robust to Wikipedia’s absence, while using sentences-as-queries suffers a bigger drop in performance.

x: The Crown-of-thorns starfish is native to the Great Barrier Reef... The starfish was introduced to the Great-Barrier-Reef by **ocean currents**.
e: [invasivespeciesinfo.gov] **Ballast water** is one of the major pathways for the introduction of nonindigenous marine species...
y: The Crown-of-thorns starfish is native to the Great Barrier Reef... The starfish was introduced to the Great-Barrier-Reef by **ballast water**.

Figure 8: **Disabling the agreement model leads to over-edits.** Here, the evidence *e* does not explicitly disagree with *x*, but without an agreement model to detect this, the edit model makes an unsupported change.

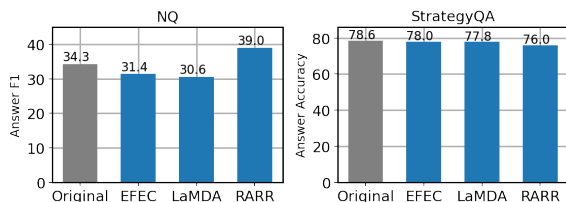


Figure 9: **Downstream task performance** on NQ and SQA. RARR’s revisions lead to better answer accuracy on NQ. No models improved answer accuracy on SQA.

finned in Section 2.2. However, another measure of preservation is whether the revised text can still be used to perform the task that it was originally generated for. Following EFEC, we quantitatively evaluate this on short answer tasks NQ and SQA, and we summarize the result in Figure 9.

For NQ, each original text *x* is a long-form response to a factoid question. To determine whether the revised text *y* still serves this purpose, we feed the factoid question and *y* back into PaLM and prompt it to extract a *short* answer from *y*. We find

that RARR not only preserves the short answer accuracy but actually improves it by roughly 5%.

For SQA, each original text is a reasoning chain that helps to answer a yes/no question. We feed the SQA question and *y* back into PaLM and prompt it to output a yes/no answer, and evaluate answer accuracy. Here, we find that increasing attribution comes at a slight cost in downstream task performance: answer accuracy drops modestly for *all* revision models (up to 2.6%). We suspect that this may be due to noisy retrievals, which sometimes provide misleading evidence (exemplified in Figure 7d). Furthermore, even though revisions can address factoid errors in the passage (e.g., “*Homer Simpson has 5 fingers*” from Figure 7e), RARR currently does not try to modify subsequent reasoning steps which may no longer be logically entailed (e.g., “*He only needs one hand to count to 5*”).

7 Conclusion

Language models have developed increasingly good “procedural” knowledge of what should be discussed and how it should be presented, but often struggle to memorize “factoid” knowledge and produce unsubstantiated claims. We proposed RARR, a framework for revising such claims to make them attributable to the researched evidence. From experiments on text passages generated by different models on various domains, we showed that RARR can revise the passages to improve attribution while preserving other desirable properties such as writing style or structure. Furthermore, RARR sits on top of existing generation models without needing to re-design or re-train LMs.

Major headroom still remains, as discussed in Section 6 and the Limitations section. We hope our analysis of RARR would help with developing new approaches for integrating attribution to LMs.

8 Limitations

Limitations of our task definition Depending on the application, attribution and preservation may not deserve equal weight. For instance, if there are multiple acceptable options for the output, such as in a dialog system, we might trade-off preservation for attribution, similar to how LaMDA behaves in our experiments.

Our evaluation metrics also do not measure all aspects of attribution. For instance, some sentences are self-evident and do not require attribution (e.g., “*I agree.*”) but would be penalized in our evaluation. It is also necessary to note that linguistic assertions have varying scope: for example, there is a difference between “*Frozen is a scary movie*” and “*I got scared watching Frozen*” — while expressing a similar sentiment, the former makes a more general statement that many would disagree with, while the latter is scoped to the speaker’s own experience. In some applications, one could even argue that the latter case does not require attribution, since the speaker is their own source-of-truth. In addition to varying scope, utterances can also make assertions with varying levels of directness. For example, according to standard linguistics, “*John ate some of the cookies*” yields the implicature that John did not eat *all* of the cookies, even though it is not logically entailed. This raises the question of which implicatures or implied assertions should be detected and attributed, which should be explored in future work. For more nuances, we refer to [Rashkin et al. \(2021\)](#).

For preservation, we wish to explore other properties that should be preserved, such as discourse or logical coherence. Additionally, if the input text passage is completely misguided or flawed, it can be difficult to revise the text without significant changes, which would be heavily penalized by the current metrics.

Limitations of our model While we aspire to improve attribution for arbitrary text, it is clear that RARR is not yet fully general. For example, the current implementation of RARR would not be well-prepared to edit poetry (where preserving rhyme matters) or long documents, primarily because we do not provide examples of such inputs in our few-shot LLM prompts. However, we do believe that future developers may be able to quickly adapt RARR to such tasks by simply changing the prompts. Second, RARR tends to preserve rather

than delete claims that it cannot attribute. Some of these claims genuinely do not require attribution, but others are hallucination and should be removed. Judging whether a claim requires attribution can be subjective and challenging. Finally, our model is computationally costly, since it is based on prompting a large language model. One potential solution is to leverage recent synthetic data generation recipes to train a smaller model ([Lee et al., 2021](#); [Schick et al., 2022](#)).

9 Ethical considerations

Partial attribution When RARR is not 100% successful in making text consistent with retrieved evidence, the revised text will be partially attributed. One could identify unattributed parts using either the automated attribution score (Attr_{AIS}) or the relevance scores used to generate the attribution report (Section 3.3). Such information should be presented to avoid misleading readers into thinking that the entire revision is attributed.

Evidence trustworthiness RARR seeks to improve attribution for the output of any generative model. However, even if RARR can attribute content to a particular source, the user must still consider whether the source itself is trustworthy. Even for sources that are traditionally considered “authoritative” (such as an encyclopedia), there may still be factual inaccuracies or biases. This work does not address the question of whether a source is trustworthy, or the related topic of misinformation. While we do not provide a means for judging trustworthiness, the design of RARR does allow for the research stage to restrict its search over a user-specified corpus, based on what the user deems trustworthy.

Conflicting evidence There is also the possibility that some content may be simultaneously supported by certain sources, while contradicted by others. This can easily occur for content involving subjective or imprecise claims. The current implementation and evaluation for RARR does not explicitly address this issue — we adopted a “permissive” definition of attribution, where we consider content to be attributed if there exists any source that supports it. For some applications, a more restrictive definition that requires both existence of supporting sources and absence of contradicting sources would be needed.

Acknowledgments

We wish to thank Raphael Hoffmann, Slav Petrov, Dipanjan Das, Michael Collins, Iftekhar Naim, Kristina Toutanova, William Cohen, Sundeep Tirumalareddy, Samer Hassan, Quoc Le and Heng-Tze Cheng for their research mentorship, feedback and support. We are grateful to Hao Zhou and Petr Pilar for helping us experiment with LaMDA and motivating our dialog experiments. We also wish to thank Tal Schuster for pointing us to relevant work in the fact checking literature, and helping us reproduce it. We thank Vitaly Nikolaev, David Reitter and Roe Aharoni for helping us use AIS and auto-AIS. We also wish to thank Jianmo Ni and Honglei Zhuang for developing the query-evidence relevance model we use, Daniel Andor for developing the sentence decontextualization model we use, and Ran Tian for the initial prototype of CQGen. Finally, we thank Kathy Meier-Hellstern, Philip Parham and Diane Korngiebel for their thoughtful feedback on ethical considerations.

Contributions

Luyu Gao: Designed RARR’s few-shot prompting strategies and implemented the first PaLM-based prototype. Analyzed results, and advised on the design of human and automatic evaluation.

Zhuyun Dai: Proposed the evaluation setup of editing long-form generations from PaLM/LaMDA on various QA datasets. Hosted and mentored Luyu Gao (student researcher) in prototyping RARR. Implemented the final models, designed overall experiments, and obtained main results and ablations (together with Ice Pasupat). Contributed many parts of the writing.

Ice Pasupat: Implemented the final models, designed overall experiments, and obtained main results and ablations (together with Zhuyun Dai). Automated experimental infrastructure, conducted error analyses, and oversaw many parts of the paper writing.

Anthony Chen: Developed the automatic evaluation for attribution and preservation and worked with Arun Chaganty to design human evaluation. Developed the open-source implementation (GPT-3 RARR), made improvements to prompts, and helped with writing.

Arun Chaganty: Led and implemented all human evaluation. Proposed the two-dimensional attribution + preservation metric (together with

Kelvin Guu). Advised on model design and contributed many parts of the writing.

Yicheng Fan: Worked with Kelvin Guu to develop the first prototype of RARR. Proposed multiple retrieval strategies and implemented the EFEC baseline.

Vincent Zhao: Co-hosted and mentored Luyu Gao (student researcher) in prototyping RARR. Enabled bulk inference for PaLM. Proposed the downstream task evaluation.

Ni Lao: Research mentorship, advising and contributed many parts of the writing.

Hongrae Lee: Research mentorship and advising. Helped integrate RARR with Google Search and evaluate LaMDA.

Da-Cheng Juan: Research mentorship and early design discussions.

Kelvin Guu: Proposed the original research-and-revise concept, implemented the first prototype, initiated the project and involved all collaborators. Implemented baselines (together with Yicheng Fan). Research mentorship, oversaw project coordination and paper writing.

References

- Michael Ahn et al. 2022. [Do as I can, not as I say: Grounding language in robotic affordances](#). *ArXiv*, abs/2204.01691.
- Raviteja Anantha, Svitlana Vakulenko, Zhucheng Tu, Shayne Longpre, Stephen G. Pulman, and Srinivas Chappidi. 2021. [Open-domain question answering goes conversational via question rewriting](#). In *NAACL*.
- Isabelle Augenstein, Christina Lioma, Dongsheng Wang, Lucas Chaves Lima, Casper Hansen, Christian Hansen, and Jakob Grue Simonsen. 2019. [MultiFC: A real-world multi-domain dataset for evidence-based fact checking of claims](#). In *EMNLP*.
- Vidhisha Balachandran, Hannaneh Hajishirzi, William Cohen, and Yulia Tsvetkov. 2022. [Correcting diverse factual errors in abstractive summarization via post-editing and language model infilling](#). In *EMNLP*.
- Bernd Bohnet, Vinh Quang Tran, Pat Verga, Roe Aharoni, Daniel Andor, Livio Baldini Soares, Jacob Eisenstein, Kuzman Ganchev, Jonathan Herzig, Kai Hui, Tom Kwiatkowski, Ji Ma, Jianmo Ni, Tal Schuster, William W. Cohen, Michael Collins, Dipanjan Das, Donald Metzler, Slav Petrov, and Kellie Webster. 2022. [Attributed question answering: Evaluation and modeling for attributed large language models](#). *ArXiv*.

- Samuel R. Bowman, Gabor Angeli, Christopher Potts, and Christopher D. Manning. 2015. [A large annotated corpus for learning natural language inference](#). In *EMNLP*.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, T. J. Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeff Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *NeurIPS*.
- Mengyao Cao, Yue Dong, Jiapeng Wu, and Jackie Chi Kit Cheung. 2020. [Factual error correction for abstractive summarization models](#). In *EMNLP*.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. [Reading Wikipedia to answer open-domain questions](#). In *ACL*.
- Eunsol Choi, Jennimaria Palomaki, Matthew Lamm, Tom Kwiatkowski, Dipanjan Das, and Michael Collins. 2021. [Decontextualization: Making sentences stand-alone](#). *TACL*, 9:447–461.
- Aakanksha Chowdhery et al. 2022. [PaLM: Scaling language modeling with pathways](#). *ArXiv*, abs/2204.02311.
- Nouha Dziri, Sivan Milton, Mo Yu, Osmar Zaiane, and Siva Reddy. 2022. [On the origin of hallucinations in conversational models: Is it the datasets or the models?](#) In *NAACL*.
- Alexander R. Fabbri, Wojciech Kryscinski, Bryan McCann, Richard Socher, and Dragomir Radev. 2021. [SummEval: Re-evaluating summarization evaluation](#). *TACL*, 9:391–409.
- Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Auli. 2019. [ELI5: Long form question answering](#). In *ACL*.
- Angela Fan, Aleksandra Piktus, Fabio Petroni, Guillaume Wenzek, Marzieh Saeidi, Andreas Vlachos, Antoine Bordes, and Sebastian Riedel. 2020. [Generating fact checking briefs](#). In *EMNLP*.
- Mor Geva, Daniel Khashabi, Elad Segal, Tushar Khot, Dan Roth, and Jonathan Berant. 2021. [Did Aristotle use a laptop? a question answering benchmark with implicit reasoning strategies](#). *TACL*, 9:346–361.
- Tanya Goyal and Greg Durrett. 2021. [Annotating and modeling fine-grained factuality in summarization](#). In *NAACL*.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020. [REALM: Retrieval-augmented language model pre-training](#). In *ICML*.
- Shirley Anugrah Hayati, Raphaël Olivier, Pravalika Avvaru, Pengcheng Yin, Anthony Tomasic, and Graham Neubig. 2018. [Retrieval-based neural code generation](#). In *EMNLP*.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. [Measuring massive multitask language understanding](#). In *ICLR*.
- Or Honovich, Roei Aharoni, Jonathan Herzig, Hagai Taitelbaum, Doron Kukliansy, Vered Cohen, Thomas Scialom, Idan Szpektor, Avinatan Hassidim, and Y. Matias. 2022. [TRUE: Re-evaluating factual consistency evaluation](#). In *Workshop on Document-grounded Dialogue and Conversational Question Answering*.
- Or Honovich, Leshem Choshen, Roei Aharoni, Ella Neeman, Idan Szpektor, and Omri Abend. 2021. [Q²: Evaluating factual consistency in knowledge-grounded dialogues via question generation and question answering](#). In *EMNLP*.
- Hayate Iso, Chao Qiao, and Hang Li. 2020. [Fact-based text editing](#). In *ACL*.
- Georgi Karadzhov, Preslav Nakov, Lluís Màrquez, Alberto Barrón-Cedeño, and Ivan Koychev. 2017. [Fully automated fact checking using external sources](#). In *RANLP*.
- Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Yu Wu, Sergey Edunov, Danqi Chen, and Wen tau Yih. 2020. [Dense passage retrieval for open-domain question answering](#). In *EMNLP*.
- Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke Zettlemoyer, and Mike Lewis. 2020. [Generalization through memorization: Nearest neighbor language models](#). In *ICLR*.
- Tushar Khot, Ashish Sabharwal, and Peter Clark. 2018. [SciTail: A textual entailment dataset from science question answering](#). In *AAAI*.
- Kalpesh Krishna, Aurko Roy, and Mohit Iyyer. 2021. [Hurdles to progress in long-form question answering](#). In *NAACL*.
- Wojciech Kryscinski, Bryan McCann, Caiming Xiong, and Richard Socher. 2020. [Evaluating the factual consistency of abstractive text summarization](#). In *EMNLP*.
- Tom Kwiatkowski et al. 2019. [Natural Questions: A benchmark for question answering research](#). *TACL*, 7:453–466.
- Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. 2019. [Latent retrieval for weakly supervised open domain question answering](#). In *ACL*.
- Kenton Lee, Kelvin Guu, Luheng He, Timothy Dozat, and Hyung Won Chung. 2021. [Neural data augmentation via example extrapolation](#). *ArXiv*, abs/2102.01335.

- Vladimir I. Levenshtein. 1965. Binary codes capable of correcting deletions, insertions, and reversals. *Soviet physics. Doklady*, 10:707–710.
- Patrick Lewis, Ethan Perez, Aleksandara Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Kuttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks. In *NeurIPS*.
- Robert L Logan IV, Alexandre Passos, Sameer Singh, and Ming-Wei Chang. 2022. **FRUIT: Faithfully reflecting updated information in text**. In *NAACL*.
- Shayne Longpre, Kartik Kumar Perisetla, Anthony Chen, Nikhil Ramesh, Chris DuBois, and Sameer Singh. 2021. **Entity-based knowledge conflicts in question answering**. In *EMNLP*.
- Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan T. McDonald. 2020. **On faithfulness and factuality in abstractive summarization**. In *ACL*.
- Jacob Menick, Maja Trebacz, Vladimir Mikulik, John Aslanides, Francis Song, Martin Chadwick, Mia Glaese, Susannah Young, Lucy Campbell-Gillingham, Geoffrey Irving, and Nathan McAleese. 2022. **Teaching language models to support answers with verified quotes**. *ArXiv*, abs/2203.11147.
- Reiichiro Nakano, Jacob Hilton, S. Arun Balaji, Jeff Wu, Long Ouyang, Christina Kim, Christopher Hesse, Shantanu Jain, Vineet Kosaraju, William Saunders, Xu Jiang, Karl Cobbe, Tyna Eloundou, Gretchen Krueger, Kevin Button, Matthew Knight, Benjamin Chess, and John Schulman. 2021. **WebGPT: Browser-assisted question-answering with human feedback**. *ArXiv*, abs/2112.09332.
- Shashi Narayan, Joshua Maynez, Reinald Kim Amplayo, Kuzman Ganchev, Annie Louis, Fantine Huot, Dipanjan Das, and Mirella Lapata. 2022. **Conditional generation with a question-answering blueprint**. *ArXiv*, abs/2207.00397.
- Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. **MS MARCO: A human generated machine reading comprehension dataset**. In *CoCo@NIPS*.
- Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gustavo Hernández Ábrego, Ji Ma, Vincent Zhao, Yi Luan, Keith B. Hall, Ming-Wei Chang, and Yinfei Yang. 2021. **Large dual encoders are generalizable retrievers**. *ArXiv*, abs/2112.07899.
- Maxwell Nye, Anders Johan Andreassen, Guy Gur-Ari, Henryk Michalewski, Jacob Austin, David Bieber, David Dohan, Aitor Lewkowycz, Maarten Bosma, David Luan, Charles Sutton, and Augustus Odena. 2021. **Show your work: Scratchpads for intermediate computation with language models**. *ArXiv*, abs/2112.00114.
- Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H. Miller, and Sebastian Riedel. 2019. **Language models as knowledge bases?** In *EMNLP*.
- Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Dmytro Okhonko, Samuel Broscheit, Gautier Izacard, Patrick Lewis, Barlas Ouguz, Edouard Grave, Wen tau Yih, and Sebastian Riedel. 2021. **The web is your oyster - knowledge-intensive nlp against a very large web corpus**. *ArXiv*, abs/2112.09924.
- Colin Raffel, Noam M. Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. **Exploring the limits of transfer learning with a unified text-to-text transformer**. *JMLR*, 21.
- Hannah Rashkin, Vitaly Nikolaev, Matthew Lamm, Lora Aroyo, Michael Collins, Dipanjan Das, Slav Petrov, Gaurav Singh Tomar, Iulia Turc, and David Reitter. 2021. **Measuring attribution in natural language generation models**. *ArXiv*, abs/2112.12870.
- Adam Roberts, Colin Raffel, and Noam M. Shazeer. 2020. **How much knowledge can you pack into the parameters of a language model?** In *EMNLP*.
- Timo Schick, Jane Dwivedi-Yu, Zhengbao Jiang, Fabio Petroni, Patrick Lewis, Gautier Izacard, Qingfei You, Christoforos Nalmpantis, Edouard Grave, and Sebastian Riedel. 2022. **PEER: A collaborative language model**. *ArXiv*, abs/2208.11663.
- Tal Schuster, Adam Fisch, and Regina Barzilay. 2021. **Get your vitamin C! robust fact verification with contrastive evidence**. In *NAACL*.
- Darsh J. Shah, Tal Schuster, and Regina Barzilay. 2020. **Automatic fact-guided sentence modification**. In *AAAI*.
- Richard Shin, Christopher H Lin, Sam Thomson, Charles Chen, Subhro Roy, Emmanouil Antonios Platanios, Adam Pauls, Dan Klein, Jason Eisner, and Benjamin Van Durme. 2021. **Constrained language models yield few-shot semantic parsers**. In *EMNLP*.
- Romal Thoppilan et al. 2022. **LaMDA: Language models for dialog applications**. *ArXiv*, abs/2201.08239.
- James Thorne and Andreas Vlachos. 2021. **Evidence-based factual error correction**. In *ACL*.
- James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. 2018. **FEVER: a large-scale dataset for fact extraction and verification**. In *NAACL*.
- David Wadden, Kyle Lo, Lucy Lu Wang, Shanchuan Lin, Madeleine van Zuylen, Arman Cohan, and Hananeh Hajishirzi. 2020. **Fact or fiction: Verifying scientific claims**. In *EMNLP*.
- Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020. **Asking and answering questions to evaluate the factual consistency of summaries**. In *ACL*.

William Yang Wang. 2017. “Liar, liar pants on fire”: A new benchmark dataset for fake news detection. In *ACL*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Chi, Quoc Le, and Denny Zhou. 2022. Chain of thought prompting elicits reasoning in large language models. *ArXiv*, abs/2201.11903.

Adina Williams, Nikita Nangia, and Samuel R. Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In *NAACL*.

Jingyi Zhang, Masao Utiyama, Eiichiro Sumita, Graham Neubig, and Satoshi Nakamura. 2018. Guiding neural machine translation with retrieved translation pieces. In *NAACL*.

Yuan Zhang, Jason Baldridge, and Luheng He. 2019. PAWS: Paraphrase adversaries from word scrambling. In *NAACL*.

A Additional experiments and analysis

Model variance The main experiments in Section 5 are based on a single run. We ran automated evaluation on 3 random runs of RARR, using PaLM outputs on NQ as input passages. The standard deviations of $\text{Attr}_{\text{auto}}$, Pres_{Lev} , and F1_{AP} are 1.2, 0.5, and 1.0 respectively.

Impact of the retriever choice We tried using Microsoft Bing in place of Google Search, with near identical results (< 1% difference).

Impact of model scale Many components in RARR work by few-shot prompting PaLM, a large 540B parameter LM. To assess the benefit of LM scaling, we replaced PaLM 540B with a smaller 62B parameter PaLM. As shown in Table 4, we found that 540B outperforms 62B by a large margin, suggesting that RARR could potentially further improve with even more scaling. We also experimented with keeping the editor stage at 540B while shrinking the query generation stage to 64B — this yielded a relatively small performance drop, suggesting that model scaling is more important for the editor.

Impact of model type Few-shot prompting has proven to be effective for many recent large language models. We try replacing the query generation model, agreement model, and edit model with GPT-3 text-davinci-003. The few-shot prompts were slightly tuned to fit the GPT-3 model. Table 4 shows the results, which are slightly better than RARR implemented with PaLM 540B on all three

datasets. We will release this open-source version of RARR that uses GPT-3 as the backbone.

Results on GPT-3 passages Table 5 shows automated evaluation results on passages generated by GPT-3. The results follow the same trend as the results on PaLM and LaMDA passages.

Challenging domains We report results on tasks where attribution was particularly hard, and significant future work is needed.

We considered news article summaries produced by summarization models from SummEval (Fabbri et al., 2021) (e.g., “*John Doe was left homeless when the storms hit Staten Island, New York . . .*”). Results are shown in Table 6. First, we note that the before-edit auto-AIS scores for all models are low. These news article summaries are often about less widely known people and events, which is challenging for retrievers, leading to low attribution. For example, our query generator may ask “*where does John Doe live*” but get results for a different John Doe. EFEC and LaMDA also face this issue, but instead trade preservation for attribution and rewrite the text to a different topic. This result suggests that using web search with standard question generation methods may fail to capture important context from the input, and is not sufficient for the attribution task.

We also considered long-form explanations generated by PaLM for the ELI5 dataset (Fan et al., 2019) (Table 6). ELI5 was collected from online forums, so many answers tend to have subjective opinions instead of specific entities and facts (e.g., “*How do our brains interpret scary music? To me, scary music often sounds a little bit like a person . . .*”), and are thus difficult to attribute. Sometimes the whole output is based on a false premise and needs to be completely rewritten, in which case RARR cannot satisfactorily edit due to our revision threshold (Section 3.2).

Finally, we considered technical explanations to questions from the MMLU dataset (Hendrycks et al., 2021) which covers diverse subjects from social science, humanities, STEM, and others.⁴ An example input looks like “*Every time you remove an edge from a complete graph, you divide it into two connected components. So, a complete graph with 13 vertices must have 12 connected components.*” Results are shown in Table 7. RARR im-

⁴MMLU has questions from 57 subjects; we took 10 random question from each topic and generated answer explanations by prompting PALM 540B.

Model	PaLM outputs on NQ			PaLM outputs on SQA			LaMDA outputs on QReCC		
	Attr _{auto}	Pres _{Lev}	F1 _{AP}	Attr _{auto}	Pres _{Lev}	F1 _{AP}	Attr _{auto}	Pres _{Lev}	F1 _{AP}
Full RARR	45.6 → 54.9	89.6	68.1	37.6 → 45.1	89.9	60.0	18.8 → 29.4	80.2	43.1
qgen 62B, editor 540B	45.9 → 54.6	87.8	67.4	37.0 → 40.5	90.0	55.9	15.8 → 28.4	76.1	41.4
qgen 62B, editor 62B	45.9 → 49.9	91.0	64.4	37.0 → 38.3	93.0	54.2	15.8 → 21.9	71.6	33.5
GPT-3	44.3 → 55.0	90.6	68.5	38.6 → 46.6	89.3	61.2	18.3 → 28.6	89.8	43.4

Table 4: **Additional ablation results.** We report the automatic metrics: Attr_{auto}, Pres_{Lev}, and harmonic mean between the two (F1_{AP}). We show auto-AIS scores both before and after editing (before → edit), with respect to the attribution report A produced by the model.

Model	GPT-3 outputs on NQ			GPT-3 outputs on SQA			GPT-3 outputs on QReCC		
	Attr _{auto}	Pres _{Lev}	F1 _{AP}	Attr _{auto}	Pres _{Lev}	F1 _{AP}	Attr _{auto}	Pres _{Lev}	F1 _{AP}
EFEC	48.3 → 66.8	41.5	51.2	32.6 → 50.6	29.4	37.2	26.4 → 53.1	39.0	44.9
LaMDA	36.2 → 61.1	45.9	52.4	22.3 → 27.3	43.3	33.5	19.0 → 33.9	28.3	30.8
PaLM RARR	48.3 → 57.2	89.6	69.8	32.6 → 36.3	91.6	52.0	26.4 → 31.1	87.7	45.9
GPT-3 RARR	48.0 → 59.3	91.8	72.0	34.7 → 37.0	91.8	52.8	23.2 → 25.3	89.7	39.5

Table 5: **Results on passages from GPT-3.** We report the automatic metrics: Attr_{auto}, Pres_{Lev}, and harmonic mean between the two (F1_{AP}). We show auto-AIS scores both before and after editing (before → edit), with respect to the attribution report A produced by the model. The results show similar trends as the results on passages from PaLM and LaMDA in Table 1.

Model	Attr _{auto}	Pres _{Lev}	F1 _{AP}
SummEval			
EFEC	17.9 → 34.6	20.9	26.0
LaMDA	10.3 → 28.8	28.1	28.4
RARR	18.3 → 16.9	92.9	28.6
ELI5			
EFEC	18.2 → 41.2	17.2	24.2
LaMDA	19.9 → 40.1	31.2	35.1
RARR	18.5 → 18.9	97.2	31.7

Table 6: **Results on ELI5 and SummEval.**

MMLU Category	RARR		
	Attr _{auto}	Pres _{Lev}	F1 _{AP}
Humanities	26.6 → 29.6	6.6	45.0
Social Sciences	35.5 → 40.7	7.6	56.5
STEM	37.8 → 41.5	7.2	57.4
Other	36.9 → 41.7	7.1	57.6

Table 7: **RARR results on MMLU.**

proves attribution of the explanations on all four categories of MMLU, although the increases are relatively small. We also found that RARR’s performance is low on examples with mathematical reasoning, as these are beyond the capability of the edit model with our current prompt.

B Details on automated evaluation

Sentence splitting When computing the attribution score, we use spaCy en_core_web_sm v3.0.0a1 to segment the text passage into sentences.

(More recent models gave similar results.) While each sentence may contain multiple claims that could be attributed independently, there is currently no linguistic consensus on what constitutes a claim. Instead of depending on a particular definition of claims, we use sentences as claims for simplicity and reproducibility. The same segmentation is also used for human evaluation.

Decontextualization We decontextualize each sentence in the text passage before computing the attribution score. We use the model from Choi et al. (2021), which is a T5 model fine-tuned to map the input “[HEAD] [SEP] *context and passage* [start] *sentence* [end]” to the output “[OPCODE] *decontextualized sentence*”, where the OPCODE can be “done” (success), “un” (unnecessary), or “imp” (impossible). We feed the passage’s context (questions for NQ and SQA; dialog context for QReCC) along with the passage itself to the input. We use beam search with beam size 8 and discard any result whose number of tokens differ by more than 4.

NLI model We obtained a newer version of the end-to-end NLI model from the authors of Honovich et al. (2022), which was trained on MNLI, SNLI, FEVER, PAWS, SciTail and VitaminC (Williams et al., 2018; Bowman et al., 2015; Thorne et al., 2018; Zhang et al., 2019; Khot et al., 2018; Schuster et al., 2021). The model is a T5

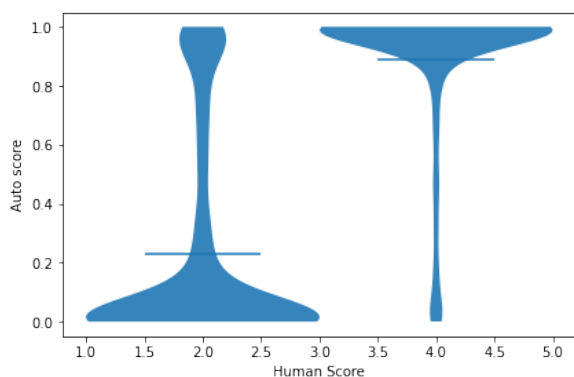


Figure 10: Violin plot illustrating the strong correlation between human AIS and auto-AIS labels on our NQ benchmark. Pearson correlation is 0.74 (N=450). y -axis is auto-AIS score, the two violins correspond to a human label of 0 or 1.

model fine-tuned to map the input “premise: *evidence* hypothesis: *claim sentence*” to either “1” (entailed) or “0” (not entailed). As suggested by the authors, we use the probability of producing “1” as the entailment score.

Comparing human and automated evaluation

We conducted correlation studies between human and automatic metrics and found strong Pearson correlation (attribution = 0.74; preservation = 0.62). We visualize the correlation between human and automated attribution scores on NQ and SQA in Figure 10. We found that the AIS scores from human correlate well with auto-AIS scores, with some bias for non-attributed sentences to be judged as attributed by auto-AIS.

C Details on human evaluation

To end-goal of RARR is to improve the *attribution* of generation models through post-editing while *preserving* the original intent. Attribution and preservation are both subjective properties that may change with even small edits. In the main paper, we present two automatic metrics to conveniently gauge these properties, but rely on a human evaluation as the gold standard. In this section, we describe how we conducted the human evaluation and what instructions and examples annotators were provided.

Rater recruitment and training We engaged with a vendor supplier of full-time crowd workers to recruit human annotators for our task. Annotators were asked to review the instructions below and were provided direct feedback on their

responses during the pilot annotation runs. We had 3 annotators rate each example in the pilot phase to measure inter-annotator agreement, and had a single rater annotate each example afterwards.

C.1 Instructions: Overview

In this task you will evaluate the quality of text generated by a system (the “**passage**”) based on how well it represents information from multiple pieces of “**evidence**”.

We will be using two categories to evaluate the quality of the passage: **Attribution** and **Intent Similarity**. You will evaluate these categories in succession. In some tasks, you will only evaluate Attribution. The task interface will guide you through the flow; you can also see the overall task flow in the diagram below.

Note: The passage may appear very fluent and well-formed, but still contain slight inaccuracies that are not easy to discern at first glance. Pay close attention to the text. Read it carefully as you would when proofreading.

C.2 Instructions: Attribution

In this step, you will evaluate how much of the passage is attributable to one or more pieces of evidence (Figure 11).

In the interface, the passage of text and the context in which it was generated is shown on the left, and each piece of evidence is shown on the right. You will use all three (context, passage, evidence) to answer the following question for each sentence in the passage: *Is all of the information provided by this sentence fully supported by at least one piece of evidence?*

Determining the information provided by the sentence. Three points are key when determining information provided by the sentence:

1. The context and the other sentences of the passage are often critical in understanding the information provided by the sentence.
2. The context should only be used to understand the information provided by the sentence.
3. The evidence should be completely ignored for this step.

Consider the following example:

Context: who plays doug williams in days of our lives

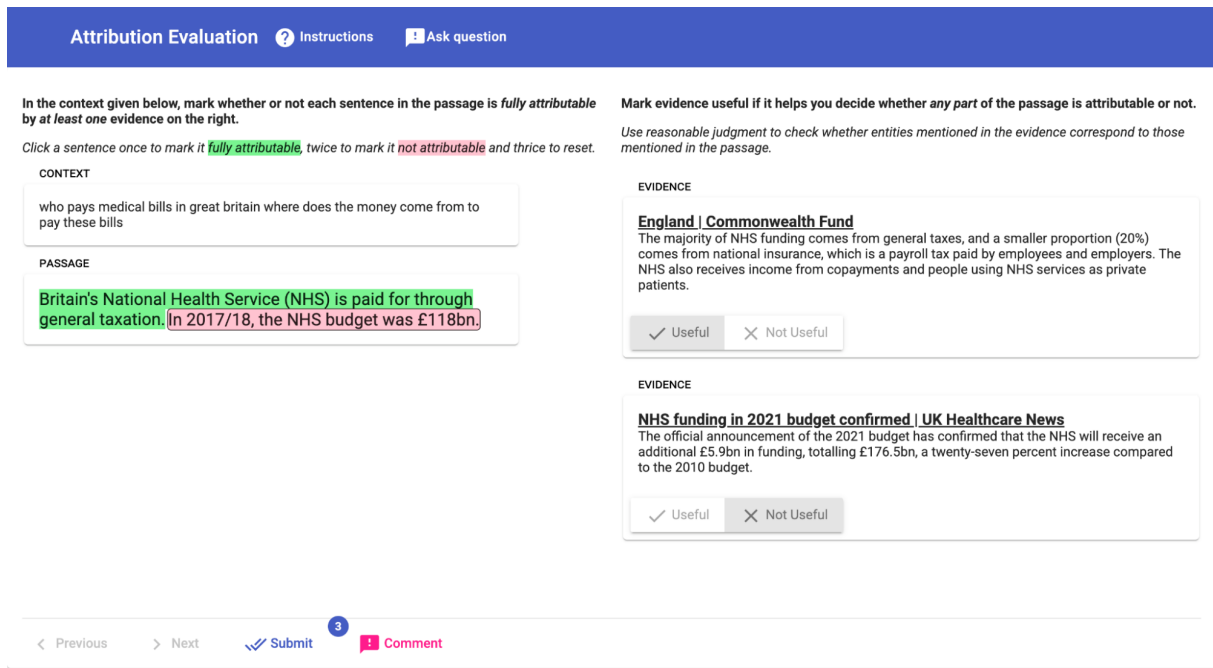


Figure 11: Screenshot of interface to annotate attribution at the sentence level. annotators were asked to mark sentences as being **fully attributable** or **not fully attributable** by clicking each sentence, and rating each piece of evidence as being useful or not in helping determine attribution of the passage. Annotators were also presented with the context of the generation.

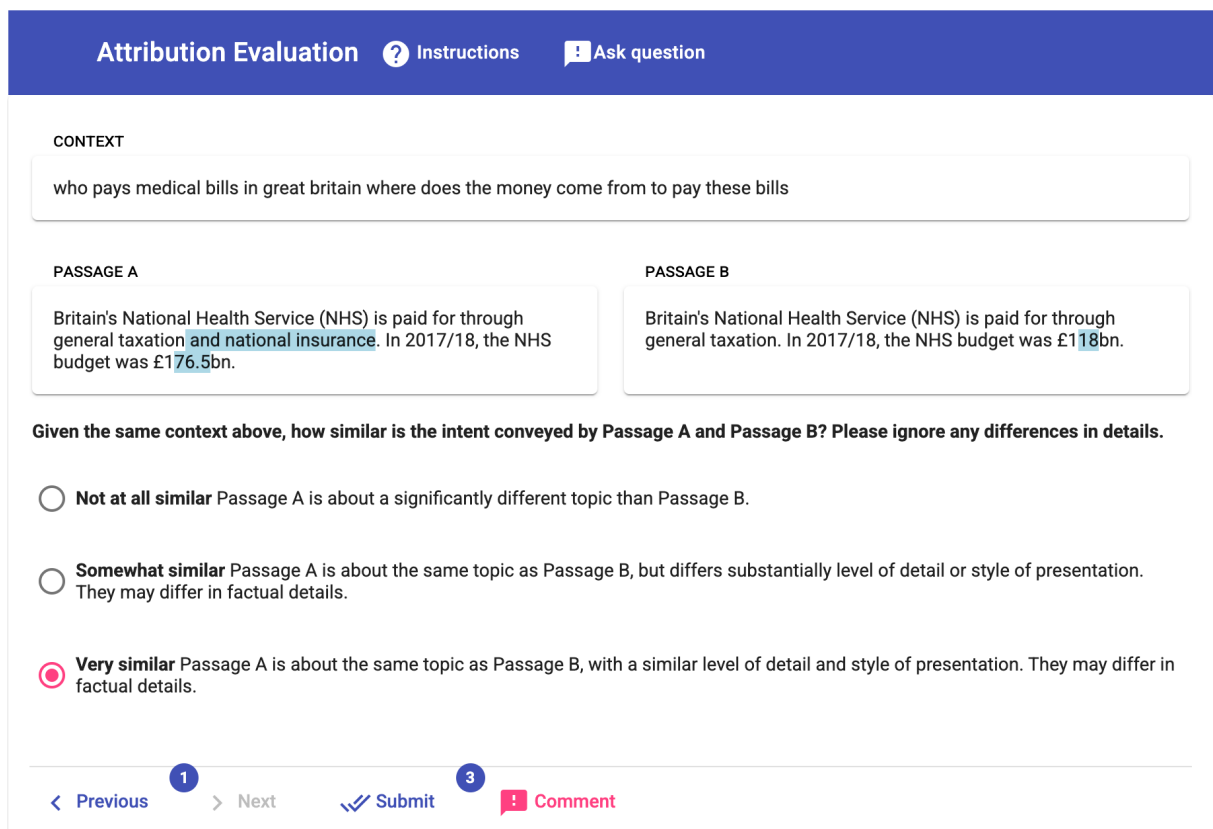


Figure 12: Screenshot of the preservation interface. Annotators are asked to read compare two passages and rate how similar the intent conveyed by the two passages is.

Passage: In the American daytime drama *Days of Our Lives*, Doug Williams and Julie Williams are portrayed by Bill Hayes and Susan Seaforth Hayes.

In the above example, the meaning of the passage is clear even without seeing the query. But consider another example:

Context: who plays doug williams in days of our lives

Passage: he is played by Bill Hayes

Passage (interpreted): Doug Williams is played by Bill Hayes in days of our lives

In this case the pronoun “he” depends on the context, but it is clear that the intended meaning of the passage can be reasonably interpreted as “Doug Williams is played by Bill Hayes in days of our lives”. This interpretation is the “information provided by the passage”.

Pronouns such as he/she/it/they etc. are one case where context is needed to figure out the intended meaning of the system response. Here’s another example (given with paraphrases of the information highlighted below):

Context: when is the last time the us lost basketball at the olympics

Passage: The last time they lost was in 2004, when Argentina defeated the US 89–79. Most recently, they won gold in 2016.

Passage (interpreted): The last time the United States lost basketball at the Olympics was in 2004.

The context should only be used to determine the information provided by the passage; at times, the passage may be about a slightly different topic than the context, for example:

Context: the south west wind blows across nigeria between

Passage: The Harmattan is a dry and dusty northeasterly trade wind that blows across West Africa from December to March. It is very dusty because it blows across the Sahara.

Here, the passage talks about a *northeasterly wind*, while the context asks about a *south-west wind*, but the passage can be fully understood.

In general, use your best judgment to determine the information provided by the passage. If the passage is hard to understand and you are unsure what the intended meaning of the passage is, *mark the sentences as not attributed* and enter a comment with an explanation. As one example, take the following:

Context: how many NBA championships did Michael Jordan win?

Passage: it is the best team in the NBA

Determining if the information accurately represents the evidence. Two points are key when determining whether the information accurately represents the evidence: When interpreting a piece of evidence, use only the title and text of that specific evidence. Completely ignore the context, passage and all other evidence. Check all the information in a sentence. If only some information is supported by the evidence, mark the sentence as not fully attributable.

Consider the following example:

Context: when did reba mcentire record back to god

Passage: Back to God was released by McEntire in 2017.

Evidence: “Back to God” is a song performed by American singer, Reba McEntire. It was released as the second single from her 2017 album, Sing it Now: Songs of Faith & Hope, on January 20, 2017.

In the above example, it is reasonable to conclude that the evidence supports all the information in the passage, and we can mark the passage as being fully attributable. But consider another example:

Context: who won the womens 2017 ncaa basketball tournament

Passage: South Carolina Gamecocks won the 2017 NCAA Women’s Division I Basketball Tournament.

Evidence: The South Carolina Gamecocks defeated the Mississippi State Bulldogs, 67–55, to claim their first-ever national championship.

In this case, while the evidence also mentions the “South Carolina Gamecocks”, it isn’t clear that the national championship being mentioned is indeed the 2017 NCAA Women’s Division I Basketball

Tournament. The passage should be marked as not attributable.

Finally, when the passage contains multiple sentences, evaluate whether each sentence can be fully attributed to one or more pieces of evidence—it is possible for one sentence to be attributed while another is not. For example:

Context: who won the womens 2017 ncaa basketball tournament

Passage: South Carolina Gamecocks won the 2017 NCAA Women’s Division I Basketball Tournament. The final score is 67-55. The championship game was held in Dallas, Texas.

Evidence 1: The South Carolina Gamecocks defeated the Mississippi State Bulldogs, 67–55, to claim their first-ever national championship.

Evidence 2: The 2017 NCAA Women’s Division I Basketball Tournament was played from Friday, March 17 to Sunday, April 2, 2017, with the Final Four played at the American Airlines Center in Dallas, Texas on March 31 and April 2.

The first two sentences cannot be attributed to either evidence for the same reason as the previous example, but the last sentence is fully supported by Evidence 2 and should be marked as attributed.

In general, you should use your best judgment in determining whether all of the information provided by the passage is “an accurate representation of information in at least one evidence”. See Table 8 for additional examples.

We give the following final notes of guidance:

- **Marking evidence as useful.** When reviewing each piece of evidence, mark it as useful if it helps you judge the attributability of any sentence, and mark it not useful if not. In the above example Evidence 1 is not useful because it didn’t contain enough context to actually help you assess if the passage was attributable, but Evidence 2 was useful.
- **Contradicting evidence.** Mark a sentence as being attributed if any piece of evidence supports it: if two pieces of evidence contradict each other, but one of them supports the passage, mark the sentence as fully attributable.

- **More on the concept of “accurate representation”.** We take as inspiration the journalist’s conception of “accurate representation”. For example, take [this excerpt on Accuracy in the NPR Ethics Handbook](#): “When quoting or paraphrasing anyone . . . consider whether the source would agree with the interpretation...” In other words, if you had written the source document, consider whether you would view the system response as an accurate representation of information in that source document.

C.3 Instructions: Intent Similarity

In this step, you will evaluate how much similar the passage is to another passage (Figure 12).

In the interface, the passage A and passage B are both text generated by a system—given the same context in which it was generated. You will use all three (context, passage A, passage B) to answer the following question: *How similar is the intent expressed by Passage A and Passage B? Please ignore any differences in details.*

Two points are key when determining whether the two passages convey the same intent:

1. Judge the similarity solely based on the similarity in the type and quantity of information provided by each passage.
2. Ignore any differences in factual details between the two passages.

Consider the following examples:

Context: who pays medical bills in great britain where does the money come from to pay these bills

Passage A: Britain’s National Health Service (NHS) is paid for through general taxation and national insurance. In 2017/18, the NHS budget was £176.5bn.

Passage B: Britain’s National Health Service (NHS) is paid for through general taxation. In 2017/18, the NHS budget was £118bn.

Rating: Very similar. Passage A is about the same topic as Passage B, with a similar level of detail and style of presentation. They may differ in factual details.

The above example should be rated “very similar” because both passages include information about (1) how the NHS is paid for, and (2) what its budget

Context + Passage	Evidences	Notes
<p><i>Context: who played morticia in the addams family tv show</i> <u>The Addams Family is an American animated sitcom TV series.</u> <u>It was first aired on NBC on September 24, 1973. Carolyn Jones played the role of Morticia.</u></p>	<p>1/ The Addams Family (1973 TV series): The Addams Family is an American animated sitcom adaptation of the Charles Addams comic. The series was produced in 1973 and was rebroadcast the following season. 2/ The Addams Family (TV Series 1964–1966): When The Addams Family went off the air in 1966, network executives in charge of children’s programming for NBC brought them back in 1973 for their own Saturday Morning cartoon show featuring the voices of Carolyn Jones from the original series.</p>	<p>While the evidence supports the show being aired in 1973, it doesn’t specify the exact date (September 24, 1973). Similarly, while the evidence mentions Carolyn Jones as being a voice actor, it doesn’t say she played the role of Morticia.</p>
<p><i>Context: when will the la sagrada familia be finished</i> <u>The La Sagrada Familia is a large Roman Catholic church in Barcelona. It is designed by Antoni Gaudi.</u> <u>It started construction in 1882, and the construction is still going on. The estimated date to finish is 2026.</u></p>	<p>1/ Sagrada Família - Wikipedia: The Basílica i Temple Expiatori de la Sagrada Família is a church in the Eixample district of Barcelona, Catalonia, Spain, and is currently the largest unfinished Roman Catholic church. 2/ Find Out Sagrada Familia’s Expected Finish Date: Visiting the breathtaking Sagrada Familia today also means witnessing the slow progress towards the completion of the project. Sagrada Familia is now expected to be completed in 2026, the centenary of Gaudi’s death. It’s a reasonable inference that La Sagrada Familia is the same as Sagrada Familia, even though the names differ slightly.</p>	<p>While Evidence 2 mentions Gaudi, isn’t clear this is a reference to Antoni Gaudi and further doesn’t say that he designed the church.</p>

Table 8: Additional examples for annotating attribution.

in 2017/18 was, though they differ in their actual answers to these questions.

Context: who is the owner of reading football club

Passage A: Reading’s owner is Yongge Dai. Yongge Dai is also the president of Chinese company Dai Yongge Real Estate. Yongge’s son, Dai Xiu Li, is Reading’s vice-president.

Passage B: Reading’s owner is Dai Yongge. Yongge’s brother and sister pair behind the Reading FC takeover—Dai Yongge and Dai Xiu Li—has made their fortune through a massive property empire. Mr Dai, has been the chairman of Renhe Commercial since 1999, which is an organisation owned by his sister behind a vast network of underground shopping centres in China.

Rating: Somewhat similar. Passage A is about the same topic as Passage B, but differs substantially in level of detail or style of presentation. They may differ in factual details.

The above example should be rated “somewhat similar” because both passages are still about the same topic—Reading’s owner— but differ substantially in the information they discuss: Passage A includes information about (1a) who Reading’s owner is, (2a) which company they are the pres-

ident of and (3a) who their vice-president is. In contrast, while Passage B shares information about (1a), it also includes information about (2b) how the Reading owner made their fortune, (3b) their company position and how long they held it for and (4b) what the company also owns.

Context: what is the numbers of total elected member of indian parliment in present time

Passage A: The total number of elected members of the Lok Sabha is 543.

Passage B: The total number of elected members of the Rajya Sabha is 238.

Rating: Not at all similar. Passage A is about a significantly different topic than Passage B.

Even though the passages look very similar, the above example should be rated “not at all similar” because the two passages are about significantly different topics: “the Lok Sabha” vs “the Rajya Sabha”.

D Details on the model

Few-shot prompting with LLMs We implement many sub-tasks within RARR using *few-shot prompting* of LLMs (also known as *in-context learning* (Brown et al., 2020)) as follows:

1. For each sub-task, we manually author a small number of training examples:

$(\text{input}_j, \text{output}_j)$ for $j = 1, \dots, J$, where J ranges between 5 and 10 and where both the input and output are strings.

2. We form the following prompt: $\text{input}_1 \diamond \text{output}_1 \oplus \text{input}_2 \diamond \text{output}_2 \oplus \dots \oplus \text{input}_J \diamond \text{output}_J \oplus \text{new_input}$, where \diamond denotes a newline character and \oplus denotes a double newline character.
3. To perform inference on a new input, we condition the LLM on the prompt and sample continuations of the prompt up until the next double newline character.

All of our prompts are included in Figures 13, 14, and 15. The contextual version used for QReCC are in Figures 16, 17, and 18.

Model statistics We implemented most parts of RARR with the PALM model which has 540B parameters. We prompted PALM without any training or finetuning. We used a TPU v2-128 to run inference with PALM.

We manually wrote our prompts by eye-balling quality on a dozen of examples from a separate validation set. We tune our hyperparameters on the validation set as well. We used sampling temperature 0.7 for all generation tasks. For each input text, we sample 3 question generations, and for each question we retrieve 5 results. For agreement gate and editing, we only sample 1 generation. We reject an editing if the edit distance is more than 50 characters or more than half of the original text length.

E Details on the dataset

As explained in Section 5.1, we generated 150 development and 150 test passages for each of the 6 combinations of dataset and model: (NQ, PaLM), (SQA, PaLM), (QReCC, LaMDA), (NQ, GPT-3), (SQA, GPT-3), (QReCC, GPT-3). Figures 19, 20, 21, and 22 are the few-shot prompts used to generate the passages.

Following the corresponding datasets, all generated passages are in English. The authors have manually looked through most of the data and found no personal identifiers.

```

1 [web] I will check things you said and ask questions.
2
3 (1) You said: Your nose switches back and forth between nostrils. When you sleep, you switch about every 45 minutes. This
is to prevent a buildup of mucus. It's called the nasal cycle.
4 To verify it,
5 a) I googled: Does your nose switch between nostrils?
6 b) I googled: How often does your nostrils switch?
7 c) I googled: Why does your nostril switch?
8 d) I googled: What is nasal cycle?
9
10 (2) You said: The Stanford Prison Experiment was conducted in the basement of Encina Hall, Stanford's psychology building.
To verify it,
11 a) I googled: Where was Stanford Prison Experiment was conducted?
12
13 (3) You said: The Havel-Hakimi algorithm is an algorithm for converting the adjacency matrix of a graph into its adjacency
list. It is named after Vaclav Havel and Samih Hakimi.
14 To verify it,
15 a) I googled: What does Havel-Hakimi algorithm do?
16 b) I googled: Who are Havel-Hakimi algorithm named after?
17
18 (4) You said: "Time of My Life" is a song by American singer-songwriter Bill Medley from the soundtrack of the 1987 film
Dirty Dancing. The song was produced by Michael Lloyd.
19 To verify it,
20 a) I googled: Who sings "Time of My Life"?
21 b) I googled: Which film is "Time of My Life" from?
22 c) I googled: Who produced the song "Time of My Life"?
23
24 (5) You said: Kelvin Hopins was suspended from the Labor Party due to his membership in the Conservative Party.
To verify it,
25 a) I googled: Why was Kelvin Hopins suspended from Labor Party?
26
27 (6) You said: Social work is a profession that is based in the philosophical tradition of humanism. It is an intellectual
discipline that has its roots in the 1800s.
28 To verify it,
29 a) I googled: What philosophical tradition is social work based on?
30 b) I googled: What year does social work has its root in?
31
32 (7) You said: {text}
33 To verify it,
34 -----
35
36

```

Figure 13: Few-shot prompt for query generation. To increase diversity and coverage, we sample the model three times and combine the resulting lists of queries.

1 [web] I will check some things you said.
2
3 (1) You said: Your nose switches back and forth between nostrils. When you sleep, you switch about every 45 minutes. This is to prevent a buildup of mucus. It's called the nasal cycle.
4 I checked: How often do your nostrils switch?
5 I found this article: Although we don't usually notice it, during the nasal cycle one nostril becomes congested and thus contributes less to airflow, while the other becomes decongested. On average, the congestion pattern switches about every 2 hours, according to a small 2016 study published in the journal PLOS One.
6 Your nose's switching time is about every 2 hours, not 45 minutes.
7 This disagrees with what you said.
8
9 (2) You said: The Little House books were written by Laura Ingalls Wilder. The books were published by HarperCollins.
10 I checked: Who published the Little House books?
11 I found this article: These are the books that started it all - the stories that captured the hearts and imaginations of children and young adults worldwide. Written by Laura Ingalls Wilder and published by HarperCollins, these beloved books remain a favorite to this day.
12 The Little House books were published by HarperCollins.
13 This agrees with what you said.
14
15 (3) You said: The Stanford Prison Experiment was conducted in the basement of Jordan Hall, Stanford's psychology building.
16 I checked: Where was Stanford Prison Experiment conducted?
17 I found this article: Carried out August 15-21, 1971 in the basement of Jordan Hall, the Stanford Prison Experiment set out to examine the psychological effects of authority and powerlessness in a prison environment.
18 The Stanford Prison Experiment was conducted in Jordan Hall.
19 This agrees with what you said.
20
21 (4) You said: Social work is a profession that is based in the philosophical tradition of humanism. It is an intellectual discipline that has its roots in the 1800s.
22 I checked: When did social work have its roots?
23 I found this article: The Emergence and Growth of the Social work Profession
24 Social work's roots were planted in the 1880s, when charity organization societies (COS) were created to organize municipal voluntary relief associations and settlement houses were established.
25 Social work has its roots in the 1880s, not 1800s.
26 This disagrees with what you said.
27
28 (5) You said: The Havel-Hakimi algorithm is an algorithm for converting the adjacency matrix of a graph into its adjacency list. It is named after Vaclav Havel and Samih Hakimi.
29 I checked: What is the Havel-Hakimi algorithm?
30 I found this article: The Havel-Hakimi algorithm constructs a special solution if a simple graph for the given degree sequence exists, or proves that one cannot find a positive answer. This construction is based on a recursive algorithm. The algorithm was published by Havel (1955), and later by Hakimi (1962).
31 Havel-Hakimi algorithm is for constructing a special solution if a simple graph for the given degree sequence exists, or proving that one cannot find a positive answer, not converting the adjacency matrix of a graph into its adjacency list.
32 This disagrees with what you said.
33
34 (6) You said: "Time of My Life" is a song by American singer-songwriter Bill Medley from the soundtrack of the 1987 film Dirty Dancing. The song was produced by Michael Lloyd.
35 I checked: Who was the producer of "(I've Had) The Time of My Life"?
36 I found this article: On September 8, 2010, the original demo of this song, along with a remix by producer Michael Lloyd, was released as digital files in an effort to raise money for the Patrick Swayze Pancreas Cancer Research Foundation at Stanford University.
37 "Time of My Life" was produced by Michael Lloyd.
38 This agrees with what you said.
39
40 (7) You said: Kelvin Hopkins was suspended from the Labor Party because he had allegedly sexually harassed and behaved inappropriately towards a Labour Party activist, Ava Etemadzadeh.
41 I checked: Why was Kelvin Hopkins suspended from the Labor Party?
42 I found this article: A former Labour MP has left the party before an inquiry into sexual harassment allegations against him was able to be concluded, the party has confirmed. Kelvin Hopkins was accused in 2017 of inappropriate physical contact and was suspended by the Labour party pending an investigation. This agrees with what you said.
43 Kelvin Hopkins was suspended because he had allegedly sexually harassed and behaved inappropriately towards a Labour Party activist, Ava Etemadzadeh.
44 This agrees with what you said.
45
46 (8) You said: In the battles of Lexington and Concord, the British side was led by General Thomas Smith.
47 I checked: Who led the British side in the battle of Lexington and Concord?
48 I found this article: Interesting Facts about the Battles of Lexington and Concord. The British were led by Lieutenant Colonel Francis Smith. There were 700 British regulars.
49 The British side was led by Lieutenant Colonel Francis Smith, not General Thomas Hall.
50 This disagrees with what you said.
51
52 (9) You said: {text}
53 I checked: {query}
54 I found this article: {evidence}
55 -----

Figure 14: Few-shot prompt for the agreement model, which uses chain-of-thought prompting.

1 [web] I will fix some things you said.
2
3 (1) You said: Your nose switches back and forth between nostrils. When you sleep, you switch about every 45 minutes. This is to prevent a buildup of mucus. It's called the nasal cycle.
4 I checked: How often do your nostrils switch?
5 I found this article: Although we don't usually notice it, during the nasal cycle one nostril becomes congested and thus contributes less to airflow, while the other becomes decongested. On average, the congestion pattern switches about every 2 hours, according to a small 2016 study published in the journal PLOS One.
6 This suggests 45 minutes switch time in your statement is wrong.
7 My fix: Your nose switches back and forth between nostrils. When you sleep, you switch about every 2 hours. This is to prevent a buildup of mucus. It's called the nasal cycle.
8
9 (2) You said: In the battles of Lexington and Concord, the British side was led by General Thomas Hall.
10 I checked: who led the British side in the battle of Lexington and Concord?
11 I found this article: Interesting Facts about the Battles of Lexington and Concord. The British were led by Lieutenant Colonel Francis Smith. There were 700 British regulars.
12 This suggests General Thomas Hall in your statement is wrong.
13 My fix: In the battles of Lexington and Concord, the British side was led by Lieutenant Colonel Francis Smith.
14
15 (3) You said: The Stanford Prison Experiment was conducted in the basement of Encina Hall, Stanford's psychology building.
16 I checked: where was Stanford Prison Experiment conducted.
17 I found this article: Carried out August 15-21, 1971 in the basement of Jordan Hall, the Stanford Prison Experiment set out to examine the psychological effects of authority and powerlessness in a prison environment.
18 This suggests Encina Hall in your statement is wrong.
19 My fix: The Stanford Prison Experiment was conducted in the basement of Jordan Hall, Stanford's psychology building.
20
21 (4) You said: Phoenix Mills Ltd., a diversified business conglomerate, was established in 1854. It has a history of over 160 years.
22 I checked: When was Phoenix Mills Ltd. founded?
23 I found this article: Phoenix Mills Ltd was incorporated in the year 1905. The company began their operations as a textile manufacturing company on 17.3 acres of land at Lower Parel in Mumbai. In the year 1959 the company was listed in the Bombay Stock Exchange.
24 This suggests the year of establishment 1854 in your statement is wrong.
25 My fix: Phoenix Mills Ltd., a diversified business conglomerate, was established in 1905. It has a history of over 160 years.
26
27 (5) You said: The Havel-Hakimi algorithm is an algorithm for converting the adjacency matrix of a graph into its adjacency list. It is named after Vaclav Havel and Samih Hakimi.
28 I checked: What is the Havel-Hakimi algorithm?
29 I found this article: The Havel-Hakimi algorithm constructs a special solution if a simple graph for the given degree sequence exists, or proves that one cannot find a positive answer. This construction is based on a recursive algorithm. The algorithm was published by Havel (1955), and later by Hakimi (1962).
30 This suggests the Havel-Hakimi algorithm's functionality in your statement is wrong.
31 My fix: The Havel-Hakimi algorithm constructs a special solution if a simple graph for the given degree sequence exists, or proves that one cannot find a positive answer. It is named after Vaclav Havel and Samih Hakimi
32
33 (6) You said: "Time of My Life" is a song by American singer-songwriter Bill Medley from the soundtrack of the 1987 film Dirty Dancing. The song was produced by Phil Ramone.
34 I checked: Who was the producer of "(I've Had) The Time of My Life"?
35 I found this article: On September 8, 2010, the original demo of this song, along with a remix by producer Michael Lloyd, was released as digital files in an effort to raise money for the Patrick Swayze Pancreas Cancer Research Foundation at Stanford University.
36 This suggests "Time of My Life" producer name in your statement is wrong.
37 My fix: "Time of My Life" is a song by American singer-songwriter Bill Medley from the soundtrack of the 1987 film Dirty Dancing. The song was produced by Michael Lloyd.
38
39 (7) You said: Phoenix Market City Pune is located on 21 acres of prime property in Pune. It is spread across four levels with approximately 1.4 million square feet of built-up space. The mall is owned and operated by Phoenix Mills Limited.
40 I checked: What is the area of Phoenix Market City in Pune?
41 I found this article: Phoenix Market City was opened in January 2013 and has the distinction of being the largest mall in the city of Pune, with the area of 3.4 million square feet. It is located in the Viman Nagar area of Pune.
42 This suggests the 1.4 million square feet of built-up space in your statement is wrong.
43 My fix: Phoenix Market City Pune is located on 21 acres of prime property in Pune. It is spread across four levels with approximately 3.4 million square feet of built-up space. The mall is owned and operated by Phoenix Mills Limited.
44
45 (8) You said: {text}
46 I checked: {query}
47 I found this article: {evidence}
48 This suggests _____

Figure 15: Few-shot prompt for the revision model, which uses chain-of-thought prompting.


```

1 [web] I will read the context and check only the last thing you said by asking questions.
2
3 (1) Context: Your nose switches back and forth between nostrils. When you sleep, you switch about every 45 minutes.
4 You said: This is to prevent a buildup of mucus. It's called the nasal cycle.
5 To verify what you just said,
6 a) I googled: Why does your nostril switch during sleep?
7 b) I googled: What is nasal cycle?
8 c) I googled: What is the nostril switching during sleep called?
9
10 (2) Context: The Stanford Prison Experiment was conducted in the basement of Encina Hall, Stanford's psychology building.
11 You said: It is a psychological study to observe the behaviors of conflict and violence that happen between inmates and
12 prisoners in real prisons.
13 To verify what you just said,
14 a) I googled: What type of experiment was the Stanford Prison Experiment?
15 b) I googled: What was the objective of the Stanford Prison Experiment?
16
17 (3) Context: The Havel-Hakimi algorithm is an algorithm for converting the adjacency matrix of a graph into its adjacency
18 list.
19 You said: It is named after Vaclav Havel and Samih Hakimi.
20 To verify what you just said,
21 a) I googled: Who are Havel-Hakimi algorithm named after?
22
23 (4) Context: "Time of My Life" is a song by American singer-songwriter Bill Medley from the soundtrack of the 1987 film
24 Dirty Dancing.
25 You said: The song was produced by Michael Lloyd in the same year.
26 To verify what you just said,
27 a) I googled: Who produced the song "Time of My Life"?
28 b) I googled: When was the song "Time of My Life" by Bill Medley produced?
29
30 (5) Context: The Late Show with Stephen Colbert is an American late-night talk show hosted by Stephen Colbert, which
31 premiered on September 8, 2015.
32 You said: Produced by Spartina Productions and CBS Television Studios, it is the second iteration of CBS' Late Show
33 franchise.
34 To verify what you just said,
35 a) I googled: Who produces "The Late Show with Stephen Colbert"?
36 b) I googled: What are the iterations of CBS' Late Show franchise?
37
38 (6) Context: Super Mario Sunshine was released on GameCube in 2002. In the game, Mario uses a tool strapped to his back
39 called FLUDD, which stands for The Flash Liquidizer Ultra Dousing Device.
40 You said: It can be used to spray water at objects or enemies. This allows Mario to change his movements, kill enemies,
41 or clean up hazards on the floor.
42 To verify what you just said,
43 a) I googled: What is the main function of FLUDD in Super Mario Sunshine?
44 b) I googled: What can FLUDD in Super Mario Sunshine be used on?
45 c) I googled: In Super Mario Sunshine, can Mario change movement with FLUDD?
46 d) I googled: In Super Mario Sunshine, can Mario kill enemies with FLUDD?
47 e) I googled: In Super Mario Sunshine, can Mario clean up hazards on the floor with FLUDD?
48
49 (7) Context: {context}
50 You said: {text}
51 To verify what you just said,
52 -----

```

Figure 16: Contextual version of the query generation prompt. The prompt works well for dialog contexts from QReCC even though the few-shot examples are not formatted as such.

```

1 [web] I will check some things you said.
2
3 (1) Context: Your nose switches back and forth between nostrils. It's called the nasal cycle. This is to prevent a buildup
4 of mucus.
5 You said: When you sleep, you switch about every 45 minutes.
6 I checked: How often do your nostrils switch?
7 I found this article: Although we don't usually notice it, during the nasal cycle one nostril becomes congested and thus
8 contributes less to airflow, while the other becomes decongested. On average, the congestion pattern switches about every
9 2 hours, according to a small 2016 study published in the journal PLOS One.
10 Your nose's switching time is about every 2 hours, not 45 minutes.
11 This disagrees with what you said.
12
13 (2) Context: The Little House books is a series of American children's novels.
14 You said: The books were published by HarperCollins.
15 I checked: Who published the Little House books?
16 I found this article: These are the books that started it all - the stories that captured the hearts and imaginations of
17 children and young adults worldwide. Written by Laura Ingalls Wilder and published by HarperCollins, these beloved books
18 remain a favorite to this day.
19 The Little House books were published by HarperCollins.
20 This agrees with what you said.
21
22 (3) Context: The Stanford Prison Experiment is a psychological study to observe the behaviors of conflict and violence
23 that happen between inmates and prisoners in real prisons.
24 You said: It was conducted in the basement of Jordan Hall, Stanford's psychology building.
25 I checked: Where was Stanford Prison Experiment conducted?
26 I found this article: Carried out August 15-21, 1971 in the basement of Jordan Hall, the Stanford Prison Experiment set
27 out to examine the psychological effects of authority and powerlessness in a prison environment.
28 The Stanford Prison Experiment was conducted in Jordan Hall.
29 This agrees with what you said.
30
31 (4) Context: Social work is a profession that is based in the philosophical tradition of humanism.
32 You said: It is an intellectual discipline that has its roots in the 1800s.
33 I checked: When did social work have its roots?
34 I found this article: The Emergence and Growth of the Social work Profession<br><br> Social work's roots were planted in
35 the 1880s, when charity organization societies (COS) were created to organize municipal voluntary relief associations and
36 settlement houses were established.
37 Social work has its roots in the 1880s, not 1800s.
38 This disagrees with what you said.
39
40 (5) Context: The Havel-Hakimi algorithm is named after Vaclav Havel and Samih Hakimi.
41 You said: It is an algorithm for converting the adjacency matrix of a graph into its adjacency list.
42 I checked: What is the Havel-Hakimi algorithm?
43 I found this article: The Havel-Hakimi algorithm constructs a special solution if a simple graph for the given degree
44 sequence exists, or proves that one cannot find a positive answer. This construction is based on a recursive algorithm.
45 The algorithm was published by Havel (1955), and later by Hakimi (1962).
46 Havel-Hakimi algorithm is for constructing a special solution if a simple graph for the given degree sequence exists, or
47 proving that one cannot find a positive answer, not converting the adjacency matrix of a graph into its adjacency list.
48 This disagrees with what you said.
49
50 (6) Context: "Time of My Life" is a song by American singer-songwriter Bill Medley from the soundtrack of the 1987 film
51 Dirty Dancing.
52 You said: The song was produced by Michael Lloyd in the same year.
53 I checked: Who was the producer of "(I've Had) The Time of My Life"?
54 I found this article: On September 8, 2010, the original demo of this song, along with a remix by producer Michael Lloyd,
55 was released as digital files in an effort to raise money for the Patrick Swayze Pancreas Cancer Resarch Foundation at
56 Stanford University.
57 The song "Time of My Life" was produced by Michael Lloyd.
58 This agrees with what you said.
59
60 (7) Context: Super Mario Sunshine was released on GameCube in 2002. In the game, Mario uses a tool strapped to his back
61 called FLUDD.
62 You said: FLUDD stands for Functional Language in a Unified Design Discipline. It can be used to spray water at objects
63 or enemies. This allows Mario to change his movements, kill enemies, or clean up hazards on the floor.
64 I checked: What does FLUDD stands for in Super Mario Sunshine?
65 I found this article: The Flash Liquidizer Ultra Dousing Device, abbreviated and better known as FLUDD or F.L.U.D.D.,
66 is a multipurpose water pack from Super Mario Sunshine invented by Professor Elvin Gadd, indicated by the Gadd Science,
67 Incorporated logo at the base of its nozzle exclusively during the cutscene at Pinna Park.
68 In Super Mario Sunshine, FLUDD stands for the Flash Liquidizer Ultra Dousing Device, not Functional Language in a Unified
69 Design Discipline.
70 This disagrees with what you said.
71
72 (8) Context: {context}
73 You said: {text}
74 I checked: {query}
75 I found this article: {evidence}
76 -----

```

Figure 17: Contextual version of the agreement model prompt.

```

1 [web] I will fix some things you said.
2
3 (1) Context: Your nose switches back and forth between nostrils. It's called the nasal cycle. This is to prevent a buildup
4 of mucus.
5 You said: When you sleep, you switch about every 45 minutes.
6 I checked: How often do your nostrils switch?
7 I found this article: Although we don't usually notice it, during the nasal cycle one nostril becomes congested and thus
8 contributes less to airflow, while the other becomes decongested. On average, the congestion pattern switches about every
9 2 hours, according to a small 2016 study published in the journal PLOS One.
10 This suggests 45 minutes switch time in your statement is wrong.
11 My fix: When you sleep, you switch about every 2 hours.
12
13 (2) Context: The Little House books is a series of American children's novels.
14 You said: The books were published by Amberjack Publishing.
15 I checked: Who published the Little House books?
16 I found this article: These are the books that started it all - the stories that captured the hearts and imaginations of
17 children and young adults worldwide. Written by Laura Ingalls Wilder and published by HarperCollins, these beloved books
18 remain a favorite to this day.
19 This suggests Amberjack Publishing in your statement is wrong.
20 My fix: The books were published by HarperCollins.
21
22 (3) Context: The Stanford Prison Experiment is a psychological study to observe the behaviors of conflict and violence
23 that happen between inmates and prisoners in real prisons.
24 You said: It was conducted in the basement of Encina Hall, Stanford's psychology building.
25 I checked: where was Stanford Prison Experiment conducted.
26 I found this article: Carried out August 15-21, 1971 in the basement of Jordan Hall, the Stanford Prison Experiment set
27 out to examine the psychological effects of authority and powerlessness in a prison environment.
28 This suggests Encina Hall in your statement is wrong.
29 My fix: It was conducted in the basement of Jordan Hall, Stanford's psychology building.
30
31 (4) Context: The Havel-Hakimi algorithm is named after Vaclav Havel and Samih Hakimi.
32 You said: It is an algorithm for converting the adjacency matrix of a graph into its adjacency list.
33 I checked: What is the Havel-Hakimi algorithm?
34 I found this article: The Havel-Hakimi algorithm constructs a special solution if a simple graph for the given degree
35 sequence exists, or proves that one cannot find a positive answer. This construction is based on a recursive algorithm.
36 The algorithm was published by Havel (1955), and later by Hakimi (1962).
37 This suggests the Havel-Hakimi algorithm's functionality in your statement is wrong.
38 My fix: It constructs a special solution if a simple graph for the given degree sequence exists, or proves that one cannot
39 find a positive answer.
40
41 (5) Context: "Time of My Life" is a song by American singer-songwriter Bill Medley from the soundtrack of the 1987 film
42 Dirty Dancing.
43 You said: The song was produced by Phil Ramone in the same year.
44 I checked: Who was the producer of "(I've Had) The Time of My Life"?
45 I found this article: On September 8, 2010, the original demo of this song, along with a remix by producer Michael Lloyd,
46 was released as digital files in an effort to raise money for the Patrick Swayze Pancreas Cancer Resarch Foundation at
47 Stanford University.
48 This suggests "Time of My Life" producer name in your statement is wrong.
49 My fix: The song was produced by Michael Lloyd in the same year.
50
51 (6) Context: Phoenix Market City Pune is located on 21 acres of prime property in Pune. The mall is owned and operated by
52 Phoenix Mills Limited.
53 You said: It is spread across four levels with approximately 1.4 million square feet of built-up space.
54 I checked: What is the area of Phoenix Market City in Pune?
55 I found this article: Phoenix Market City was opened in January 2013 and has the distinction of being the largest mall
56 in the city of Pune, with the area of 3.4 million square feet. It is located in the Viman Nagar area of Pune.
57 This suggests the 1.4 million square feet of built-up space in your statment is wrong.
58 My fix: It is spread across four levels with approximately 3.4 million square feet of built-up space.
59
60 (7) Context: {context}
61 You said: {text}
62 I checked: {query}
63 I found this article: {evidence}
64 This suggests -----

```

Figure 18: Contextual version of the revision model prompt.

```

1 [web] I will think step by step and answer your question.
2
3 Question: is growing seedless cucumber good for a gardener with entomophobia
4 Explanation: Entomophobia is a fear of insects. Plants need insects to pollinate them. Seedless fruits such as seedless
5 cucumbers do not require pollination, so seedless fruits do not require insects. This makes good for people with
6 entomophobia.
7 Answer: Yes
8
9 Question: Who was british pm and viceroy during quit india movement?
10 Explanation: The Quit India Movement was launched in 8th August 1942. The british PM at that time was Winston Churchill.
11 The british viceroy during the movement is Victor Hope, usually referred to as Lord Linlithgow.
12 Answer: Winston Churchill and Lord Linlithgow
13
14 Question: Which year does game over man come out on netflix?
15 Explanation: Game Over, Man! is an action-comedy movie released in March 23, 2018, on Netflix. It's director is Kyle
16 Newacheck. The movie stars Anders Holm, Adam DeVine, and Blake Anderson.
17 Answer: March 23, 2018
18
19 Question: would it be very difficult for Nuno Gomes to dive to the Red Sea's deepest point?
20 Explanation: Nuno Gomes' deepest dive in the Red Sea to date is 317 metres. The Red Sea has a maximum depth of over 3,000
21 metres. So it would be difficult for Nuno Gomes to the deepest point of the Red sea.
22 Answer: Yes
23
24 Question: Are chinchillas cold-blooded?
25 Explanation: Chinchillas are rodents, which are mammals. All mammals are warm-blooded.
26 Answer: No
27
28 Question: what are the chemo drugs for gastric cancer?
29 Explanation: Chemotherapy can be given for gastric cancer before or after surgery. Usually you have a combination of 2 or
30 3 drugs. The most common types are: epirubicin, cisplatin and capecitabine (ECX) ; epirubicin, cisplatin and fluorouracil
31 (ECF); or fluorouracil, folinic acid, oxaliplatin and docetaxel (FLOT).
32 Answer: epirubicin, cisplatin and capecitabine (ECX) ; epirubicin, cisplatin and fluorouracil (ECF); or fluorouracil,
33 folinic acid, oxaliplatin and docetaxel (FLOT)
34
35 Question: How many sacks does Clay Matthews have in his career?
36 Explanation: Clay Matthews III had 383 solo tackles, 136 assists, and 91.5 sacks. His father, Clay Matthews Jr., also
37 palyed in NFL and had 69.5 sacks in his career.
38 Answer: 91.5 (Clay Matthews III) and 69.5 (Clay Matthews Jr.)
39
40 Question: In the U.S., the most deaths are caused by
41 Explanation: The leading cause of death in the U.S. is heart disease, causing nearly 700,000 deaths annually. The second
42 most common cause of deaths is cancer.
43 Answer: heart disease
44
45 Question: the rank of indian economy in terms of nominal gdp is
46 Explanation: India's GDP is economy is $2.94 trillion. It is the fifth-largest in the world. The top GDP contries are
47 United States, China, Japan, Germany and India.
48 Answer: 5
49
50 Question: {question}
51 Explanation: -----

```

Figure 19: The PaLM prompt for generating long-form answers to questions from NQ and SQA.

1 I will think step by step and answer your question.
2
3 1. Question: Is growing seedless cucumber good for a gardener with entomophobia?
4 2. Explanation: Entomophobia is a fear of insects. Plants need insects to pollinate them. Seedless fruits such as seedless cucumbers do not require pollination so seedless fruits do not require insects. This is good for people with entomophobia.
5 3. Answer: Yes.
6
7 1. Question: Who was British PM and Biceroy during Quit India Movement?
8 2. Explanation: The Quit India Movement was launched in 8th August 1942. The British PM at that time was Winston Churchill. The British Biceroy during the movement was Victor Hope, usually referred to as Lord Linlithgow.
9 3. Answer: Winston Churchil and Lord Linlithgow.
10
11 1. Question: Which year does Game Over Man come out on Netflix?
12 2. Explanation: Game Over, Man! is an action-comedy movie. Its director is Kyle Newacheck. The movie stars Anders Holm, Adam DeVine, and Blake Anderson. The movie was released March 23, 2018 on Netflix.
13 3. Answer: March 23, 2018.
14
15 1. Question: Would it be very difficult for Nuno Gomes to dive to the Red Sea's deepest point?
16 2. Explanation: Nuno Gomes' deepest dive in the Red Sea to date is 317 meters. The Red Sea has a maximum depth of over 3,000 meters. So it would be difficult for Nuno Gomes to dive to the deepest point of the Red Sea.
17 3. Answer: Yes.
18
19 1. Question: Are chinchillas cold-blooded?
20 2. Explanation: Chinchillas are rodents. Rodents are mammals. All mammals are warm-blooded.
21 3. Answer: No.
22
23 1. Question: What are the chemo drugs for gastric cancer?
24 2. Explanation: Chemotherapy can be given for gastric cancer before or after surgery. Usually you have a combination of 2 or 3 drugs. The most common types are: epirubicin, cisplatin and capecitabine (ECX) ; epirubicin, cisplatin and fluorouracil (ECF); or fluorouracil, folinic acid, oxaliplatin and docetaxel (FLOT).
25 3. Answer: Epirubicin, cisplatin and capecitabine (ECX) ; epirubicin, cisplatin and fluorouracil (ECF); or fluorouracil, folinic acid, oxaliplatin and docetaxel (FLOT).
26
27 1. Question: How many sacks does Clay Matthews have in his career?
28 2. Explanation: Clay Matthews has been pro linebacker for 11 seasons. Ten of the seasons he played for Green Bay Packers, while his last season was with Los Angeles Rams, dedicated to his hometown. Clay Matthews III had 383 solo tackles, 136 assists, and 91.5 sacks. His father, Clay Matthews Jr., also played in NFL and had 69.5 sacks in his career.
29 3. Answer: 91.5 (Clay Matthews III) and 69.5 (Clay Matthews Jr.).
30
31 1. Question: In the U.S., the most deaths are caused by?
32 2. Explanation: According to the Centers for Disease Control and Prevention, there were 3 million deaths in 2018. The leading cause of death in the U.S. is heart disease, causing nearly 700,000 deaths annually. The second most common cause of deaths is cancer.
33 3. Answer: Heart disease.
34
35 1. Question: The rank of Indian economy in terms of nominal GDP is?
36 2. Explanation: The top GDP contries are United States, China, Japan, Germany, and India. India's GDP is economy is \$2.94 trillion. It is the fifth-largest economy in the world.
37 3. Answer: 5.
38
39 1. Question: Who won season 16 on Dancing with the Stars?
40 2. Explanation: Dancing with the Stars is an American dance competition television series that premiered on June 1, 2005, on ABC. The 16th season premiered on March 18, 2013. It was won by Kellie Pickler.
41 3. Answer: Kellie Pickler.
42
43 1. Question: {question}
44 2. Explanation: _____

Figure 20: The GPT-3 prompt for generating long-form answers to questions from NQ and SQA.

```

1 Hi, I will think step by step and answer your question.
2 Is growing seedless cucumber good for a gardener with entomophobia?
3 Yes. Entomophobia is a fear of insects. Plants need insects to pollinate them. Seedless cucumber fruit does not require
  pollination. So seedless cucumber does not require insects.
4 Here is another question. Which year does game over man come out on netflix?
5 Game Over, Man! is an American action-comedy movie. It was released on March 23, 2018, on Netflix. It is directed by Kyle
  Newacheck, written by Anders Holm, and starring himself, Adam DeVine, and Blake Anderson.
6 Here is another question. What are the chemo drugs for gastric cancer?
7 Chemotherapy can be given for gastric cancer before or after surgery. Usually you have a combination of 2 or 3 drugs. The
  most common types are: epirubicin, cisplatin and capecitabine (ECX) ; epirubicin, cisplatin and fluorouracil (ECF); or
  fluorouracil, folinic acid, oxaliplatin and docetaxel (FLOT).
8 Here is another question. Was N.W.A a album?
9 N.W.A was an American hip hop group from Compton, California.
10 Were they successful?
11 They has sold over 10 million units in the United States alone.
12 How many singles did it have?
13 The group NWA released 8 singles.
14 Here is another question. {Q1}
15 {A1}
16 ...
17 {Qk}
18 -----

```

Figure 21: The LaMDA prompt for generating answers to questions from QReCC. Each line is a conversation turn. The dialog context from QReCC contains rounds of questions and answers ($Q_1, A_1, Q_2, A_2, \dots, Q_k$).

```

1 I will think step by step and answer your question.
2
3 Is growing seedless cucumber good for a gardener with entomophobia?
4 Yes. Entomophobia is a fear of insects. Plants need insects to pollinate them. Seedless cucumber fruit does not require
  pollination. So seedless cucumber does not require insects.
5
6 Which year does game over man come out on netflix?
7 Game Over, Man! is an American action-comedy movie. It was released on March 23, 2018, on Netflix. It is directed by Kyle
  Newacheck, written by Anders Holm, and starring himself, Adam DeVine, and Blake Anderson.
8
9 What are the chemo drugs for gastric cancer?
10 Chemotherapy can be given for gastric cancer before or after surgery. Usually you have a combination of 2 or 3 drugs. The
  most common types are: epirubicin, cisplatin and capecitabine (ECX) ; epirubicin, cisplatin and fluorouracil (ECF); or
  fluorouracil, folinic acid, oxaliplatin and docetaxel (FLOT).
11
12 Was N.W.A an album?
13 N.W.A was an American hip hop group from Compton, California.
14 Were they successful?
15 They has sold over 10 million units in the United States alone.
16 How many singles did they have?
17 N.W.A had eight singles, including "Straight Outta Compton", "Express Yourself", "Gangsta Gangsta", "Dopeman" and "Always
  Into Somethin'".
18
19 {Q1}
20 {A1}
21 ...
22 {Qk}
23 -----

```

Figure 22: The GPT-3 prompt for generating answers to questions from QReCC.

ACL 2023 Responsible NLP Checklist

A For every submission:

- A1. Did you describe the limitations of your work?
8
- A2. Did you discuss any potential risks of your work?
8,9
- A3. Do the abstract and introduction summarize the paper’s main claims?
1
- A4. Have you used AI writing assistants when working on this paper?
Left blank.

B Did you use or create scientific artifacts?

2,3,5

- B1. Did you cite the creators of artifacts you used?
2,3,5
- B2. Did you discuss the license or terms for use and / or distribution of any artifacts?
E
- B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?
5
- B4. Did you discuss the steps taken to check whether the data that was collected / used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect / anonymize it?
E
- B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?
5,E
- B6. Did you report relevant statistics like the number of examples, details of train / test / dev splits, etc. for the data that you used / created? Even for commonly-used benchmark datasets, include the number of examples in train / validation / test splits, as these provide necessary context for a reader to understand experimental results. For example, small differences in accuracy on large test sets may be significant, while on small test sets they may not be.
5

C Did you run computational experiments?

5,6

- C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?
5,D

The Responsible NLP Checklist used at ACL 2023 is adopted from NAACL 2022, with the addition of a question on AI writing assistance.

- C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?
5,D
- C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run?
5
- C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation), did you report the implementation, model, and parameter settings used (e.g., NLTK, Spacy, ROUGE, etc.)?
5,B,D
- D** **Did you use human annotators (e.g., crowdworkers) or research with human participants?**
5,C
- D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.?
C
- D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)?
C
- D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating? For example, if you collected data via crowdsourcing, did your instructions to crowdworkers explain how the data would be used?
C
- D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?
Not applicable. Left blank.
- D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?
C