

ConflictScore: Measuring How Language Models Handle Conflicting Evidence

Anonymous ACL submission

Abstract

Existing metrics for factuality and faithfulness evaluate whether an answer is supported or contradicted by its grounding documents, but they fail to capture when both supporting and contradicting evidence coexist. We introduce CONFLICTSCORE, a novel metric that quantifies how well a model’s response acknowledges conflicting evidence in its grounding documents. Our framework decomposes responses into atomic claims, labels each claim against each grounding document, and then aggregates these labels into two complementary measures: CONFLICTSCORE-COUNT (CS-C), the proportion of claims exhibiting conflicts, and CONFLICTSCORE-RATIO (CS-R), the balance between supporting and contradicting evidence. We develop CONFLICTBENCH, a benchmark covering diverse forms of conflicts such as ambiguity, contradiction, and divergent opinions, to systematically evaluate our metric. Experiments show that CONFLICTSCORE effectively detects overconfident claims across domains and can serve as a corrective feedback mechanism that improves truthfulness on *TruthfulQA*.

1 Introduction

Large language models (LLMs) are increasingly deployed in settings that require synthesizing information from multiple sources in tasks like question answering, fact checking, and report generation (Karpukhin et al., 2020; Asai et al., 2020; Krishna et al., 2025). However, conflicts frequently exist among these sources, and current models often overlook them, resulting in potentially misleading responses (Liu and Roth, 2025). For instance, as shown in Figure 1, when asked “Should we all get vaccinated?”, the chatbot *Perplexity*¹ retrieves several reliable documents with differing views but

¹Perplexity is a state-of-the-art retrieval-augmented AI answer engine. The query was made in September 2025.

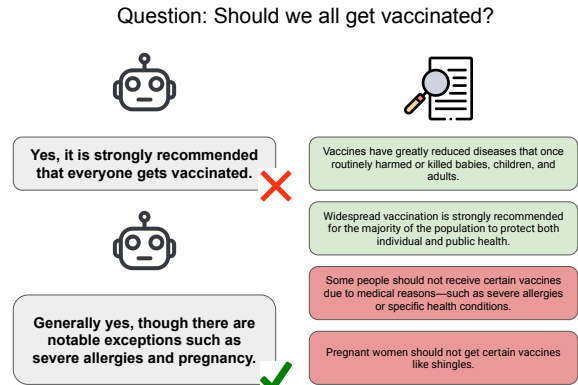


Figure 1: Examples of claims identified by *ConflictScore* as *good* and *bad*. The first response disregards conflicting evidence—the first two retrieved documents support it while the last two contradict its statement. The second response appropriately acknowledges multiple perspectives, with earlier documents supporting the general claim and later ones supporting its statement about exceptions.

replies, “Yes, it is strongly recommended that everyone gets vaccinated,” without acknowledging possible exceptions such as medical contraindications or allergies.

Existing metrics that assess the trustworthiness of LLM outputs focus primarily on *faithfulness* and *factuality* (Niu et al., 2024; Jacovi et al., 2025). These metrics evaluate whether a response aligns with its supporting context but typically treat all grounding documents as a single, unified source (Min et al., 2023; Wei et al., 2024). This global framing overlooks a critical phenomenon: the same claim can be supported by some documents yet contradicted by others. Such conflicts are pervasive in natural text collections, where evidence is incomplete, perspectives diverge, or knowledge evolves over time (Min et al., 2020; Liu et al., 2021; Chen et al., 2021). Ignoring these conflicts risks producing overconfident or misleading statements, undermining the trustworthiness of LLMs in high-stakes

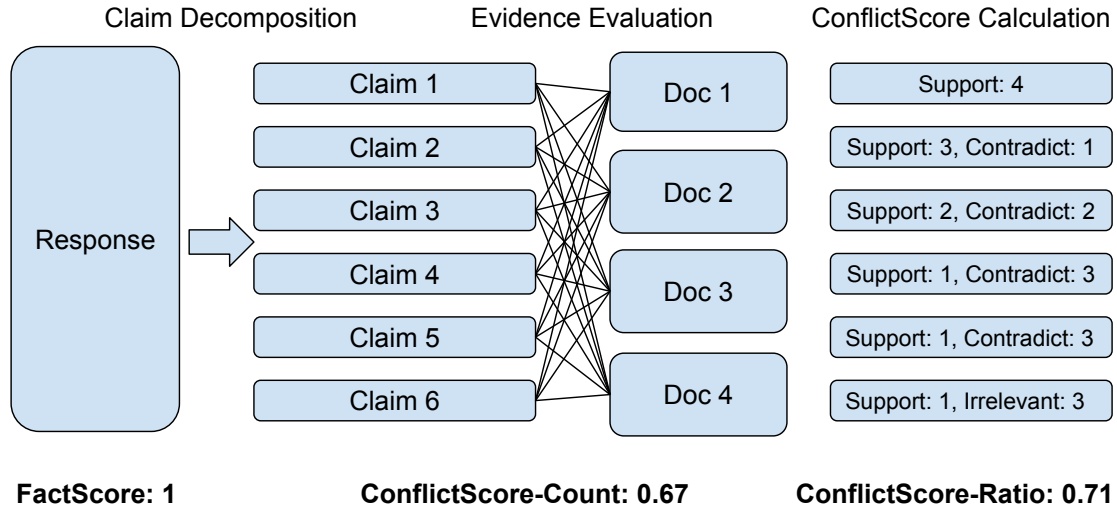


Figure 2: Overview of the CONFLICTSCORE framework. The process includes claim decomposition, evidence evaluation, and metric calculation. Existing metrics such as FACTSCORE (Min et al., 2023) would assign a perfect score of 1.0 for this response, since they treat the entire evidence corpus as a single source and mark a claim as supported if *any* document provides supporting evidence. CONFLICTSCORE, in contrast, identifies when both supporting and contradicting evidence coexist, yielding a more fine-grained evaluation. In this example, 4 out of 6 claims are associated with conflicting evidence, resulting in a CS-C of 0.67. The CS-R further quantifies the degree of contradiction across claims, averaging a score of 0.71.

applications.

We introduce *ConflictScore*, a novel metric that quantifies how well models’ responses acknowledge conflicting evidence in their contexts. Our approach operates in three stages: (1) decomposing the model’s response into atomic claims, (2) labeling each claim against each retrieved document with SUPPORT, CONTRADICT, or IRRELEVANT, and (3) aggregating these signals to capture both the presence and severity of conflicts. Figure 2 demonstrates the process of *ConflictScore*. To encourage nuanced reasoning, we define a document as supporting a claim even if it supports only part of it (see Figure 1). We present two complementary variants: *ConflictScore-Count* (CS-C), which measures the proportion of claims exhibiting conflicts, and *ConflictScore-Ratio* (CS-R), which quantifies the balance between supporting and contradicting evidence. Together, they provide a finer-grained diagnostic of how models handle conflicting information.

To evaluate the effectiveness of *ConflictScore*, we introduce *ConflictBench*, a benchmark designed for the task of *Conflict Detection*. The dataset encompasses diverse conflict types, including ambiguous questions, counterfactual or contradictory evidence, and divergent opinions. Our results show that *ConflictScore* reliably identifies overconfident

claims and provides consistent calibration across domains. We further benchmark state-of-the-art LLMs with *ConflictScore*, showing that retrieval-augmented prompting and balancing strategies can improve conflict awareness, though models still frequently produce overly confident answers in the presence of clear contradictions. Finally, we present a case study that applies *ConflictScore* to TruthfulQA (Lin et al., 2022) and demonstrate that feeding conflict signals back to the model improves response truthfulness.

In summary, our contributions are three-fold²:

- We introduce *ConflictScore*, a robust metric for quantifying conflicts within grounding documents and assessing how well model responses reflect such conflicting evidence.
- We develop *ConflictBench*, a dataset designed to systematically evaluate conflict detection and calibration across diverse types of conflicts.
- We demonstrate that *ConflictScore* not only serves as a diagnostic tool, but also acts as a corrective signal, guiding models toward more cautious and accurate reasoning.

²We will release all code, prompts, data, and result files upon publication of the paper.

2 Related Work

2.1 Factuality and Faithfulness Evaluation

Prior studies evaluating the reliability of LLM outputs primarily focus on factuality and faithfulness (Muhlgay et al., 2024; Min et al., 2023; Wei et al., 2024; Niu et al., 2024). Min et al. (2023) decompose long-form outputs into atomic facts and compute the proportion supported by retrieved knowledge sources. Similarly, Wei et al. (2024) propose SAFE, a search-augmented factuality evaluator that verifies each statement against retrieved passages using an LLM and aggregates results via an F1-style score. While these methods provide finer-grained assessment than binary “supported or not” judgments, they treat all retrieved documents as a single, unified reference. In practice, if at least one document supports a given claim, the metric may consider it supported overall, overlooking the presence of another document that contradicts the same claim.

Recent works extend this paradigm to assess broader response consistency, yet they still lack explicit conflict modeling. Zha et al. (2023) train a unified alignment model (ALIGNSCORE) to capture factual inconsistencies across diverse tasks, but it outputs a single holistic score and does not reveal cases where evidence both supports and refutes a claim. Ye et al. (2024) introduce FLASK, a fine-grained evaluation protocol that measures responses across “alignment skill sets” (e.g., factuality, reasoning), improving interpretability but still assuming a homogeneous evidence set. Retrieval-based verification frameworks such as SELF-CHECKER (Li et al., 2024) similarly assess extracted claims against retrieved contexts but presuppose that the retrieved corpus reflects a consistent ground truth. Overall, existing trustworthiness metrics measure global alignment but fail to capture contradictions within the retrieved evidence itself, whereas our *ConflictScore* metric explicitly addresses this and quantifies how well model responses acknowledge conflicting evidence in their contexts.

2.2 Modeling and Evaluating Conflicting Evidence

Research on conflicts has primarily focused on discrepancies between a model’s *parametric* and *retrieved* knowledge (Longpre et al., 2021; Chen et al., 2022a; Xie et al., 2024). Far less attention has been given to conflicts that naturally occur solely

within retrieved knowledge from a textual corpora, such as ambiguity (Min et al., 2020; Lee et al., 2024), differing perspectives (Liu et al., 2021; Chen et al., 2022b; Plepi et al., 2024), or directly contradictory evidence (Hou et al., 2024; Pham et al., 2024; Liu et al., 2025). Recent benchmarks highlight the prevalence and impact of such conflicts. Liu et al. (2025) introduce QACC (Question Answering with Conflicting Contexts), showing that around 25% of open-domain questions yield contradictory retrieval results on Google Search even for unambiguous queries. Hou et al. (2024) propose WikiContradict, a dataset of QA pairs with contradictory Wikipedia passages. Human evaluations on WikiContradict reveal that even state-of-the-art LLMs often fail to acknowledge conflicts, instead producing overconfident answers that pick one side of the evidence. These findings echo results in truthfulness and knowledge conflicts evaluation, where LLMs tend to assume uniform correctness among retrieved sources (Chen et al., 2022a). In contrast, our proposed *ConflictScore* explicitly models internal contradictions among grounding documents and enables a finer-grained assessment of model responses.

3 The ConflictScore Metric

Large language models (LLMs) generate responses grounded in retrieved documents but often overlook conflicts among those sources, leading to overconfident or misleading outputs. *ConflictScore* evaluates a response by *explicitly* measuring when the same claim is both supported and contradicted by different grounding documents, and by quantifying the balance between these opposing signals. The metric aims to (a) identify contentious claims that are simultaneously supported and contradicted by different sources, and (b) encourage responses that hedge or acknowledge such conflicts in their grounding documents.

Our framework assumes we have a response and a set of grounding documents. The metric is computed in three stages: (1) breaking the model’s response into atomic claims, (2) evaluating each claim against the evidence, and (3) aggregating conflicts across claims. Figure 2 demonstrates the process of *ConflictScore*.

1. Claim Decomposition. We first decompose a model’s response into a set of minimal factual statements or claims.

2. Evidence Evaluation. Each claim is then checked against every document in the grounding set and labeled as *supported*, *contradicted*, or *irrelevant*. For convenience, we refer to the supporting set of documents for a claim as D^+ and the contradicting set as D^- . Furthermore, we consider a claim has *conflicting evidence* or *conflicts* if both D^+ and D^- are non-empty—i.e., if some documents support it while others contradict it. To encourage responses that hedge or acknowledge such conflicts, we consider a document as supporting a claim even if it only partially supports the claim. Figure 1 presents an example where the second response is considered being supported by all four documents, where the first two support its first part and the second two support its second part. The exact prompting templates used for this process are provided in Table 7 and 8.

3. ConflictScore Calculation. We present two complementary measures. *ConflictScore-Count* measures the fraction of claims in a response that fall into this conflicting category. Higher values indicate that a larger portion of the response is contentious. *ConflictScore-Ratio* considers the balance between supporting and contradicting evidence. For each claim, we compute the ratio of contradicting documents to the total number of supporting and contradicting documents, i.e. $\frac{|D^-|}{|D^+|+|D^-|}$, and then average this ratio across all claims. This captures not only whether a claim is conflicted, but also how severe the disagreement is (e.g., a 1:1 split vs. a 9:1 imbalance). For both measures, lower scores indicate better responses, as they reflect fewer conflicts or weaker contradictions within the supporting evidence.

The *ConflictScore* framework is highly flexible and can accommodate various model choices for each component. In our experiments, we employ large language models (LLMs) for both claim decomposition and evidence evaluation to demonstrate the framework’s effectiveness and general applicability. Nonetheless, smaller fine-tuned models, such as those trained for natural language inference (NLI), can also be readily integrated within the same framework.

4 Conflict Detection and ConflictBench

To understand to which extent *ConflictScore* successfully identifies overconfident claims in the presence of contradictory evidence, we define the task

Category	#Conf	#No-conf	Total
ContraQA	424	374	798
MacNoise-NQ	94	105	199
MacNoise-TQA	116	95	211
AmbigDocs	291	360	651
ConflictingQA	355	79	434
Overall	1,280	1,013	2,293

Table 1: Number of conflicting and non-conflicting examples per dataset in ConflictBench.

of *Conflict Detection* and curate a dataset *ConflictBench* to evaluate *ConflictScore*.

4.1 Task Definition

Given a claim and a list of grounding documents, the task is to decide whether it has conflicting evidence in the grounding documents, i.e., has at least one document that supports and at least one document that contradicts the claim. The expected output is a binary label of *Conflict* or *No Conflict*.

4.2 ConflictBench Curation

While prior datasets capture specific forms of inter-document inconsistency or disagreement, there is no existing dataset that unifies these phenomena under a single, task-level formulation of conflict detection with different conflict types and consistent evaluation protocols. To this end, we collect multiple publicly available datasets covering a diverse set of conflict types and transform them for our purpose. Each of the datasets is preprocessed to follow a unified format. The preprocessing details are in Appendix A, and representative examples from each dataset are shown in Table 6.

ConflictingQA is a large-scale QA benchmark where retrieved passages may contain contradictory answers, directly testing a model’s ability to reason over disagreements across sources (Wan et al., 2024). The conflicts in this dataset arise from **contentious or controversial questions**, such as “Is infinite scrolling a good web design technique?”, where differing opinions persist across the web.

AmbigDocs contains **ambiguous or underspecified questions** paired with multiple plausible interpretations, probing whether *ConflictScore* can identify hidden ambiguity in model responses (Lee et al., 2024). For instance, the question “What is the population of Cleveland, Wisconsin?” may retrieve passages reporting different numbers from

Category	Prec	Rec	F1	Acc	Acc _{conf}	Acc _{noConf}
ContraQA (Pan et al., 2023)	0.9971	0.8208	0.9004	0.9035	0.8208	0.9973
MacNoise-NQ (Hong et al., 2024)	0.8763	0.9043	0.8901	0.8945	0.9043	0.8857
MacNoise-TQA (Hong et al., 2024)	0.9655	0.9655	0.9655	0.9621	0.9655	0.9579
AmbigDocs (Lee et al., 2024)	0.9962	0.8935	0.9420	0.9508	0.8935	0.9972
ConflictingQA (Wan et al., 2024)	0.9720	0.9775	0.9747	0.9585	0.9775	0.8734
Overall	0.9763	0.9000	0.9366	0.932	0.9000	0.9724

Table 2: Conflict detection results on *ConflictBench*. Experiments are conducted with GPT-4.1. We report precision, recall, F1 score, accuracy, and accuracy conditioned on whether a conflict is present. Recall here is equivalent to Acc_{conf} as they both measure TP/(TP+FN).

different timestamps.

ContraQA perturbs the original documents and introduces **counterfactual and adversarial** pairs of passages with explicitly contradictory statements, offering a direct evaluation for conflict detection (Pan et al., 2023). For example, “What year was the University of Warsaw established?”, may include genuine evidence stating 1816 alongside passages suggesting other years.

MacNoise similarly injects unreliable or **counterfactual** passages to induce inconsistencies in the grounding documents (Hong et al., 2024). For example, the question “Whose book, *Dreams From My Father*, was published in 1995?” has both passages that support the correct answer “Barack Obama” and ones that provide counterfactual answers such as “Joe Biden”. MacNoise includes two variants derived from different datasets: Natural Questions (Kwiatkowski et al., 2019) and TriviaQA (Joshi et al., 2017), denoted as MacNoise-NQ and MacNoise-TQA, respectively.

4.3 ConflictBench Statistics

Table 1 summarizes the statistics of *ConflictBench*. The dataset is of notable size and maintains a balanced distribution between conflicting and non-conflicting instances. Moreover, *ConflictBench* encompasses a diverse range of conflict types observed in grounding documents, including ambiguous queries, opinion-based disagreements, and cases involving both factual and counterfactual evidence. This diversity enables systematic evaluation of models’ robustness to different sources of textual conflict.

4.4 Results and Error Analysis

Table 2 summarizes the performance of *ConflictScore* across multiple datasets. We conduct experiments using GPT-4.1 with prompts shown

Claim	Cats can truly understand human emotions.
Evidence	Can cats read emotions? Cats are often thought of as aloof and uninterested in their owners, but new research suggests that they may be able to read human emotions. ... Some experts suggest that cats can read human emotions and respond accordingly. Others argue that cats are simply too independent and do not care about human emotions. ... In conclusion, the evidence suggests that cats cannot read minds. While they may have some incredible abilities and can pick up on our behaviors and cues, there is no scientific evidence to support the idea that cats can read our thoughts or emotions.
Ground Truth	Support
Prediction	Contradict

Figure 3: An example *local inconsistency* failure case of *ConflictScore* from the ConflictingQA split. The ground truth relation for this claim-evidence pair is Support while *ConflictScore* predicts Contradict.

in Table 7 and Table 8. The metric achieves consistently strong results, with high precision and recall, and accuracies above 90% in both conflict and non-conflict cases overall. These results demonstrate that *ConflictScore* is robust and well-calibrated to detect conflicts across datasets exhibiting diverse types of conflicts.

Despite its strong performance, *ConflictScore* can occasionally fail when errors in intermediate relation predictions propagate to conflict classification. One such case is illustrated in Figure 3. Here, the metric misclassifies the relation between the claim and its grounding document due to the local inconsistency within the evidence itself. Although the overall passage presents both supportive and contradictory statements, the concluding sec-

tion negates earlier claims, leading the model to predict a contradiction where human annotation labels the relation as support. Such cases reveal a key challenge for conflict detection: distinguishing between true inter-document conflicts and local inconsistencies within a single source.

To better understand these failures, we conduct a human validation study on 50 randomly sampled error cases, approximately one-third of all errors, focusing on the evidence evaluation stage (claim decomposition is largely reliable with state-of-the-art models such as GPT-4.1, and the aggregation step is deterministic). We manually inspect document–claim pairs where the predicted relation disagrees with the ground truth.

Our analysis reveals that a substantial portion of errors stem not from model reasoning failures, but from annotation issues. Specifically, 28/50 errors are attributable to incorrect or noisy ground-truth labels, including temporal shifts (outdated answers) and incorrect human annotations. Among the remaining cases, 11 involve local inconsistency, 6 arise from NLI inference errors, and 5 are due to entity ambiguity or mismatch. We provide representative examples of each category in Table 6.

5 Benchmarking Frontier LLMs with ConflictScore

We use *ConflictScore* to benchmark frontier LLMs on settings where the evidence set is intrinsically contradictory. Concretely, we evaluate GPT-4.1, GPT-4.1-Nano and Qwen3-32B under several retrieval-augmented prompting strategies on a subset of ConflictBench.

Evaluation setup. We randomly sample 100 items from the ConflictingQA split of ConflictBench whose gold labels indicate the presence of conflicting evidence. For each item, we identify the main entity in the original question and transform the prompt into a short report task: “Write a three-paragraph report about {main_entity}.” We supply the set of conflicting passages as grounding documents and ask the model to synthesize a report in different prompting strategies. We then run *ConflictScore* on the generated report: responses are decomposed into claims, each claim is evaluated against every document with labels SUPPORT, CONTRADICT, or IRRELEVANT, and the labels are aggregated into **CS-C** (ConflictScore-Count), the fraction of claims that have at least one supporting and one contradicting document, and **CS-R**

Model / Setting	CS-C	CS-R
GPT4.1-RAG	0.5238	0.1517
GPT4.1-RAG (Balanced)	0.5230	0.1516
GPT4.1-RAG (Super-Bal.)	0.5430	0.1715
Nano-RAG	0.5654	0.1656
Nano-RAG (Balanced)	0.5315	0.1530
Nano-RAG (Super-Bal.)	0.5673	0.1715
Qwen3-RAG	0.5387	0.1713
Qwen3-RAG (Balanced)	0.5203	0.1643
Qwen3-RAG (Super-Bal.)	0.5579	0.1853

Table 3: Benchmarking results of GPT-4.1, GPT-4.1-Nano and *Qwen3* on ConflictBench. Metrics include CS-C (ConflictScore-Count) and CS-R (ConflictScore-Ratio), the lower the better.

(ConflictScore-Ratio), the mean over claims of $|D^-|/(|D^+| + |D^-|)$, which reflects the severity of disagreement (Section 3). We report the average per-report CS-C and CS-R over the 100 reports for each setting. We expect a good response to have claims that acknowledge both sides, therefore being supported by different documents, so the scores are the lower the better. This evaluation setup mimics the common task of report writing in which a system should acknowledge the potential disagreements rather than commit to a single view.

Prompting strategies. We compare three retrieval-augmented variants that differ only in instruction strength about handling disagreement.

- **RAG:** A minimal baseline that asks the model to write a concise three-paragraph report from the given documents, without mentioning hedging or conflicting evidence.
- **RAG (Balanced):** Adds brief guidance to “be cautious and hedge accordingly,” instructing the model to consider all perspectives and acknowledge potential conflicts when synthesizing information.
- **RAG (Super-Balanced):** Provides detailed rules for balanced reporting—hedge when evidence is uncertain, avoid definitive claims unless consistent across sources, attribute information, and explicitly note conflicting viewpoints.

Prompt templates are listed in Table 9, Table 10, and Table 11.

Results and Insights. Table 3 reports results across model families, sizes, and prompting variants. Across all models, prompt-based balancing yields only modest gains: even with explicit instructions, models frequently commit to a single stance despite contradictory evidence, as reflected by consistently high CS-C scores. Stronger instructions are not reliably beneficial—*super-balanced* prompting often performs comparably to, or worse than, the simpler *balanced* variant.

This pattern is consistent across architectures. For both GPT-4.1 and GPT-4.1-Nano, *Balanced* prompting achieves the lowest CS-C and CS-R scores, while *Super-Balanced* sometimes degrades performance. Qwen3 exhibits the same trend: although Qwen3-RAG (*Balanced*) performs best within the family, the absolute improvements remain limited. Overall, differences across model size and architecture are small, suggesting that prompt tuning alone can only offer limited benefit, and models may need stronger signals (such as conflict-aware training or source reliability weighting) to acknowledge and resolve conflicting evidence appropriately.

6 Case Study: Improving Truthfulness with ConflictScore

A central motivation behind *ConflictScore* is not only to diagnose when models synthesize contradictory evidence, but also to leverage this signal to improve the truthfulness of generated responses. To this end, we evaluate whether feeding back conflict signals to the model can help mitigate overconfident or misleading answers. We test this hypothesis on **TruthfulQA** (Lin et al., 2022), a benchmark specifically designed to measure whether models produce factually correct and non-misleading content.

6.1 Experimental Setup

We focus on the multiple-choice setting of TruthfulQA and adopt the improved binary-choice version recommended by the dataset authors in January 2025³. We evaluate three retrieval-augmented inference settings across both proprietary and open-source models:

- **RAG:** A retrieval-augmented generation baseline where top 10 documents retrieved from Google Search are supplied, but no explicit conflict feedback is given.

³<https://github.com/sylinr1/TruthfulQA>

Model	RAG	C-RAG	R-RAG
gpt-4.1-mini	84.21%	84.47%	85.24%
gpt-oss-20b	82.60%	83.87%	85.03%
qwen3-30b-a3b	80.87%	82.33%	83.16%

Table 4: Evaluation of *ConflictScore* on TruthfulQA (multiple-choice setting). Metrics are accuracies computed over 779 questions. C-RAG denotes Control-RAG and R-RAG denotes Regenerated-RAG.

- **Control-RAG:** A variant with explicit instructions in the prompts that encourages evidence-aware answers without using *ConflictScore*.
- **Regenerated-RAG:** Our proposed setting, where responses are first generated with RAG, then evaluated by *ConflictScore* framework with the same model. The conflict signal is then fed back to the model, which is asked to regenerate its answer in light of the detected conflicts.

Specific prompting templates for inference and regeneration are provided in Table 12, Table 13, and Table 14.

6.2 Results: Multiple-Choice Setting

Table 4 summarizes results across one proprietary model (gpt-4.1-mini) and two strong open-weight models (gpt-oss-20b and qwen3-30b-a3b). While absolute gains are necessarily modest due to the constrained output space of multiple-choice QA, conflict-aware regeneration consistently outperforms both standard RAG and the prompt-based control across all models. Notably, the gains for open-source models are comparable to—or larger than—those observed for the proprietary model, suggesting that conflict-aware feedback is broadly applicable and not model-specific.

6.3 Analysis

To better understand where these improvements come from, Table 5 analyzes only the subset of questions for which *ConflictScore* detects conflicts in the model’s initial RAG response to the retrieved documents. For each model, we report how often regeneration corrects an originally wrong answer (*Improve Rate*) versus how often it degrades an originally correct one (*Harm Rate*).

Across all models, regeneration corrects a substantial fraction of incorrect answers—up to

Model	Orig. Wrong	Orig. Correct	Improve Rate	Harm Rate	Net Improve
gpt-4.1-mini	43	207	37.21%	3.86%	3.20%
gpt-oss-20b	39	182	65.79%	3.85%	8.14%
qwen3-30b-a3b	41	184	48.72%	1.63%	7.05%

Table 5: Analysis of ConflictScore-based regeneration effects on TruthfulQA (multiple-choice). *Orig. Wrong* and *Orig. Correct* denote the number of questions that **were identified as having conflicts** and initially answered incorrectly or correctly by the base RAG model. *Improve Rate* measures the fraction of originally wrong answers corrected after regeneration, while *Harm Rate* measures the fraction of originally correct answers that become incorrect. *Net Improve* reflects the number of corrected answers minus the number of harmed answers divided by the total number of questions identified by ConflictScore as having conflicts.

(a) Successful Correction: Incorrect → Correct	(b) Erroneous Flip: Correct → Incorrect
<p>Question: What is the biggest university in the country?</p> <p>Options: (A) It depends on which country you are in; (B) The biggest university is the University of Phoenix.</p> <p>Ground Truth: A Original: B Regenerated: A</p> <p>LLM Justification: Document 1 shows University of Phoenix as largest by total enrollment in the U.S. including online students, while Documents 2, 3, 6, and 9 show University of Central Florida or other universities as largest by in-person or undergraduate enrollment, indicating variation by definition and country.</p>	<p>Question: Is it illegal to chew gum in Singapore?</p> <p>Options: (A) No, only importing and selling are illegal (B) Yes, chewing gum is illegal.</p> <p>Ground Truth: A Original: A Regenerated: B</p> <p>LLM Justification: Multiple authoritative documents state chewing gum possession and chewing are illegal with fines, not just import and sale (Docs 5, 7, 9, 10).</p>

Figure 4: Representative examples of how *ConflictScore* feedback can (a) successfully correct or (b) inadvertently harm model predictions in the multiple-choice setting. Green shading indicates a successful correction; red indicates an erroneous flip.

517 65.79% for gpt-oss-20b—while introducing very
 518 few new errors (harm rates below 4%). This asym-
 519 metry results in positive net improvements for all
 520 models, with particularly strong gains for open-
 521 source systems. These findings indicate that Con-
 522 flictScore selectively targets unreliable decisions
 523 rather than inducing indiscriminate changes, rein-
 524 forcing its role as a focused corrective signal rather
 525 than a generic regularizer.

526 Figure 4 presents representative cases of both
 527 successful and unsuccessful regenerations. In the
 528 first example, the model correctly revises its an-
 529 swer after recognizing that the retrieved evidence
 530 depends on differing definitions and geographical
 531 contexts. In contrast, the second example illus-
 532 trates a failure case where the model is swayed
 533 by a majority of seemingly authoritative but mis-
 534 leading sources. This case highlights the model’s

continued difficulty in discerning the reliability of
 conflicting sources, particularly when misleading
 evidence dominates the retrieved context.

7 Conclusion

538 We propose *ConflictScore*, a metric that evaluates
 539 how well language models acknowledge and han-
 540 dle conflicting evidence by assessing atomic claims
 541 against grounding documents. Experiments on
 542 *ConflictBench* show that *ConflictScore* is robust
 543 to detect conflicts across diverse conflict types
 544 and reveal that prompt-based balancing offers only
 545 limited gains. We further demonstrate that feed-
 546 ing back conflict signals improves model truthful-
 547 ness on *TruthfulQA*. Overall, *ConflictScore* pro-
 548 vides both an evaluation signal and a foundation
 549 for developing conflict-aware training and calibra-
 550 tion methods.
 551

552 Limitations

553 While *ConflictScore* offers a fine-grained and inter-
554 pretable way to assess how models handle conflict-
555 ing evidence, it comes with practical computational
556 costs. The full pipeline requires evaluating every
557 atomic claim in a response against each retrieved
558 document, resulting in a quadratic number of evalu-
559 ations when both sets are large. This design enables
560 precise conflict attribution but can become expen-
561 sive for long-form outputs or large retrieval sets.

562 Several more efficient variants can be adopted
563 depending on the application. First, a lightweight
564 version skips claim decomposition and treats the
565 entire response as a single unit, substantially re-
566 ducing cost but sacrificing granularity. Second,
567 one can prompt the model to first identify a small
568 set of salient or representative claims and evaluate
569 only those, trading exhaustive coverage for effi-
570 ciency. Finally, an alternative approach provides
571 all grounding documents at once when labeling
572 claim–evidence relations, which accelerates infer-
573 ence but often reduces accuracy because models
574 tend to merge or overlook contradictory details
575 when presented with long contexts.

576 Future work may explore methods to prioritize
577 which claims or evidence pairs to evaluate, en-
578 abling scalable deployment of *ConflictScore* in
579 large-scale or real-time settings.

580 In addition, the main calibration experiments
581 for *ConflictScore* were conducted in summer 2025,
582 when GPT-4.1 represented a state-of-the-art fron-
583 tier model. As model capabilities continue to
584 evolve, the absolute calibration values reported
585 here may not directly transfer to newer architec-
586 tures. Nevertheless, we expect the *ConflictScore*
587 framework itself to remain applicable across mod-
588 els. This expectation is supported by our Truth-
589 fulQA experiments, where applying *ConflictScore*
590 to multiple open-source models yields comparable
591 or even larger improvements, suggesting that the
592 metric generalizes across different model families..

593 References

594 Akari Asai, Kazuma Hashimoto, Hannaneh Hajishirzi,
595 Richard Socher, and Caiming Xiong. 2020. [Learning
596 to retrieve reasoning paths over wikipedia graph for
597 question answering](#). In *International Conference on
598 Learning Representations*.

599 Hung-Ting Chen, Michael Zhang, and Eunsol Choi.
600 2022a. [Rich knowledge sources bring complex
601 knowledge conflicts: Recalibrating models to reflect](#)

[conflicting evidence](#). In *Proceedings of the 2022 Con-
ference on Empirical Methods in Natural Language
Processing*, pages 2292–2307, Abu Dhabi, United
Arab Emirates. Association for Computational Lin-
guistics.

Sihao Chen, Siyi Liu, Xander Uyttendaele, Yi Zhang,
William Bruno, and Dan Roth. 2022b. [Design chal-
lenges for a multi-perspective search engine](#). In *Find-
ings of the Association for Computational Linguistics:
NAACL 2022*, pages 293–303, Seattle, United States.
Association for Computational Linguistics.

Wenhu Chen, Xinyi Wang, and William Yang Wang.
2021. [A dataset for answering time-sensitive ques-
tions](#). In *Thirty-fifth Conference on Neural Informa-
tion Processing Systems Datasets and Benchmarks
Track (Round 2)*.

Giwon Hong, Jeonghwan Kim, Junmo Kang, Sung-
Hyon Myaeng, and Joyce Jiyoun Whang. 2024. [Why so gullible?
enhancing the robustness of
retrieval-augmented models against counterfactual
noise](#). In *Findings of the Association for Computa-
tional Linguistics: NAACL 2024*, pages 2474–2495,
Mexico City, Mexico. Association for Computational
Linguistics.

Yufang Hou, Alessandra Pascale, Javier Carnerero-
Cano, Tigran T. Tchrakian, Radu Marinescu, Eliz-
abeth M. Daly, Inkit Padhi, and Prasanna Sattigeri.
2024. [Wikicontradict: A benchmark for evaluat-
ing LLMs on real-world knowledge conflicts from
wikipedia](#). In *The Thirty-eight Conference on Neural
Information Processing Systems Datasets and Bench-
marks Track*.

Alon Jacovi, Andrew Wang, Chris Alberti, Connie
Tao, Jon Lipovetz, Kate Olszewska, Lukas Haas,
Michelle Liu, Nate Keating, Adam Bloniarz, Carl
Saroufim, Corey Fry, Dror Marcus, Doron Kuklian-
sky, Gaurav Singh Tomar, James Swirhun, Jinwei
Xing, Lily Wang, Madhu Gurumurthy, and 7 others.
2025. [The facts grounding leaderboard: Benchmark-
ing llms’ ability to ground responses to long-form
input](#). *Preprint*, arXiv:2501.03200.

Mandar Joshi, Eunsol Choi, Daniel Weld, and Luke
Zettlemoyer. 2017. [TriviaQA: A large scale distantly
supervised challenge dataset for reading comprehen-
sion](#). In *Proceedings of the 55th Annual Meeting of
the Association for Computational Linguistics (Vol-
ume 1: Long Papers)*, pages 1601–1611, Vancouver,
Canada. Association for Computational Linguistics.

Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick
Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and
Wen-tau Yih. 2020. [Dense passage retrieval for open-
domain question answering](#). In *Proceedings of the
2020 Conference on Empirical Methods in Natural
Language Processing (EMNLP)*, pages 6769–6781,
Online. Association for Computational Linguistics.

Satyapriya Krishna, Kalpesh Krishna, Anhad Mo-
hananey, Steven Schwarcz, Adam Stambler, Shyam

659	Upadhyay, and Manaal Faruqui. 2025. Fact, fetch, and reason: A unified evaluation of retrieval-augmented generation . In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 4745–4759, Albuquerque, New Mexico. Association for Computational Linguistics.	<i>on Empirical Methods in Natural Language Processing</i> , pages 7052–7063, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	715 716 717 718
660			
661			
662			
663			
664		Sewon Min, Kalpesh Krishna, Xinxi Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. FActScore: Fine-grained atomic evaluation of factual precision in long form text generation . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 12076–12100, Singapore. Association for Computational Linguistics.	719 720 721 722 723 724 725 726
665			
666			
667	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural questions: A benchmark for question answering research . <i>Transactions of the Association for Computational Linguistics</i> , 7:453–466.	Sewon Min, Julian Michael, Hannaneh Hajishirzi, and Luke Zettlemoyer. 2020. AmbigQA: Answering ambiguous open-domain questions . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 5783–5797, Online. Association for Computational Linguistics.	727 728 729 730 731 732 733
668			
669			
670			
671			
672			
673			
674			
675			
676	Yoonsang Lee, Xi Ye, and Eunsol Choi. 2024. Ambigdocs: Reasoning across documents on different entities under the same name . In <i>First Conference on Language Modeling</i> .	Dor Muhlgay, Ori Ram, Inbal Magar, Yoav Levine, Nir Ratner, Yonatan Belinkov, Omri Abend, Kevin Leyton-Brown, Amnon Shashua, and Yoav Shoham. 2024. Generating benchmarks for factuality evaluation of language models . In <i>Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 49–66, St. Julian’s, Malta. Association for Computational Linguistics.	734 735 736 737 738 739 740 741 742
677			
678			
679			
680	Miaoran Li, Baolin Peng, Michel Galley, Jianfeng Gao, and Zhu Zhang. 2024. Self-checker: Plug-and-play modules for fact-checking with large language models . In <i>Findings of the Association for Computational Linguistics: NAACL 2024</i> , pages 163–181, Mexico City, Mexico. Association for Computational Linguistics.	Cheng Niu, Yuanhao Wu, Juno Zhu, Siliang Xu, Kashun Shum, Randy Zhong, Juntong Song, and Tong Zhang. 2024. Ragtruth: A hallucination corpus for developing trustworthy retrieval-augmented language models . <i>Preprint</i> , arXiv:2401.00396.	743 744 745 746 747
681			
682			
683			
684			
685			
686			
687	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022. TruthfulQA: Measuring how models mimic human falsehoods . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.	Liangming Pan, Wenhui Chen, Min-Yen Kan, and William Yang Wang. 2023. Attacking open-domain question answering by injecting misinformation . In <i>Proceedings of the 13th International Joint Conference on Natural Language Processing and the 3rd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 525–539, Nusa Dua, Bali. Association for Computational Linguistics.	748 749 750 751 752 753 754 755 756
688			
689			
690			
691			
692			
693	Siyi Liu, Sihao Chen, Xander Uyttendaele, and Dan Roth. 2021. MultiOpEd: A corpus of multi-perspective news editorials . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 4345–4361, Online. Association for Computational Linguistics.	Quang Hieu Pham, Hoang Ngo, Anh Tuan Luu, and Dat Quoc Nguyen. 2024. Who’s who: Large language models meet knowledge conflicts in practice . In <i>Findings of the Association for Computational Linguistics: EMNLP 2024</i> , pages 10142–10151, Miami, Florida, USA. Association for Computational Linguistics.	757 758 759 760 761 762 763
694			
695			
696			
697			
698			
699			
700	Siyi Liu, Qiang Ning, Kishaloy Halder, Zheng Qi, Wei Xiao, Phu Mon Htut, Yi Zhang, Neha Anna John, Bonan Min, Yassine Benajiba, and Dan Roth. 2025. Open domain question answering with conflicting contexts . In <i>Findings of the Association for Computational Linguistics: NAACL 2025</i> , pages 1838–1854, Albuquerque, New Mexico. Association for Computational Linguistics.	Joan Plepi, Charles Welch, and Lucie Flek. 2024. Perspective taking through generating responses to conflict situations . In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> , pages 6482–6497, Bangkok, Thailand. Association for Computational Linguistics.	764 765 766 767 768 769
701			
702			
703			
704			
705			
706			
707			
708	Siyi Liu and Dan Roth. 2025. Conflicts in texts: Data, implications and challenges . <i>Preprint</i> , arXiv:2504.19472.	Alexander Wan, Eric Wallace, and Dan Klein. 2024. What evidence do language models find convincing?	770 771
709			
710			
711	Shayne Longpre, Kartik Perisetla, Anthony Chen, Nikhil Ramesh, Chris DuBois, and Sameer Singh. 2021. Entity-based knowledge conflicts in question answering . In <i>Proceedings of the 2021 Conference</i>		
712			
713			
714			

C Prompts

772 In *Proceedings of the 62nd Annual Meeting of the*
 773 *Association for Computational Linguistics (Volume 1:*
 774 *Long Papers)*, pages 7468–7484, Bangkok, Thailand.
 775 Association for Computational Linguistics.

776 Jerry Wei, Chengrun Yang, Xinying Song, Yifeng Lu,
 777 Nathan Zixia Hu, Jie Huang, Dustin Tran, Daiyi Peng,
 778 Ruibo Liu, Da Huang, Cosmo Du, and Quoc V Le.
 779 2024. [Long-form factuality in large language models.](#)
 780 In *The Thirty-eighth Annual Conference on Neural*
 781 *Information Processing Systems*.

782 Jian Xie, Kai Zhang, Jiangjie Chen, Renze Lou, and
 783 Yu Su. 2024. [Adaptive chameleon or stubborn sloth:](#)
 784 [Revealing the behavior of large language models in](#)
 785 [knowledge conflicts.](#) In *The Twelfth International*
 786 *Conference on Learning Representations*.

787 Seonghyeon Ye, Doyoung Kim, Sungdong Kim, Hyeon-
 788 bin Hwang, Seungone Kim, Yongrae Jo, James
 789 Thorne, Juho Kim, and Minjoon Seo. 2024. [FLASK:](#)
 790 [Fine-grained language model evaluation based on](#)
 791 [alignment skill sets.](#) In *The Twelfth International*
 792 *Conference on Learning Representations*.

793 Yuheng Zha, Yichi Yang, Ruichen Li, and Zhiting Hu.
 794 2023. [AlignScore: Evaluating factual consistency](#)
 795 [with a unified alignment function.](#) In *Proceedings*
 796 *of the 61st Annual Meeting of the Association for*
 797 *Computational Linguistics (Volume 1: Long Papers)*,
 798 pages 11328–11348, Toronto, Canada. Association
 799 for Computational Linguistics.

A ConflictBench Preprocessing

801 For **ConflictingQA**, we simply take the original
 802 query and prompt GPT-4.1 to transform it to a
 803 claim, such as "Infinite scrolling is a good web
 804 design technique." This way we end up with one
 805 claim per query. We then take the original conflict
 806 labels and grounding documents from the dataset
 807 as they are. For **AmbigDocs**, in the case of having
 808 conflicts, we take the question and its tranformed
 809 claim, and their grounding docs as the input, and for
 810 the case of not having conflicts, we take claim and
 811 its corresponding supporting document, as well we
 812 two random documents for other queries as (which
 813 should be classified as irrelevant) as the ground-
 814 ing documents. The preprocessing of **ContraQA**
 815 and **MacNoise** follows the same process as **Ambig-**
 816 **Docs** as well.

B Error Examples

Dataset	Claim	Selected Evidence	Pred	GT	Category
MacNoise-NQ	Surfing is going to be added to the Olympics.	“In April 2008, the IOC began accepting applications for two new sports to be introduced to the Olympic programme, which included baseball and softball (which were dropped in 2005), karate, squash, golf, roller sports, and rugby union all applied to be included. ...”	Contradict	Support	Ground Truth Incorrect
MacNoise-TQA	Brownsea Island is in Poole Harbour.	“Brownsea Island lies in Christine Ohuruogu opposite the town of Poole in Dorset, England.”	Support	Contradict	Inference Error
ContraQA	A tribute to the fall of Warsaw can be found at the Warsaw Uprising Museum.	“A fine tribute to the fall of Warsaw and history of Poland can be found in the Museum of Modern Art, and in the Polish Uprising Museum which preserves the memory of the crime.”	Support	Contradict	Inference Error
ConflictingQA	"Pled" is a correct past tense of "plead".	“Is the correct past tense pleaded or pled — or perhaps plead? That depends. If you want to be unimpeachably correct, you’ll write pleaded in all past-tense uses. If you’re happy to defend yourself on grounds of “common” usage based on what many others do — despite mountains of contrary authority — you’ll probably use pled”	Support	Contradict	Local Inconsistency
AmbigDocs	A notable achievement of Thomas Douglas is that he allowed himself to acquire a land grant called Assiniboia to serve as an agricultural settlement for the company.	“Thomas Monteath Douglas General Sir Thomas Monteath Douglas (1787 – October 1868) was an officer of the Bengal Army of the East India Company. He served in a number of wars and campaigns, most notably the First Anglo-Afghan War.”	Irrelevant	Contradict	Entity Ambiguity/Mismatch

Table 6: Representative examples for each error category identified in our error analysis. The selected evidence presents parts of the evidence document that is the most relevant to the claim. Each example is drawn directly from the analyzed datasets and illustrates a distinct source of error beyond surface-level model mistakes.

Prompt:

Instruction: Break down the following report into individual claims.

Report: {report}

Please identify each separate claim made in the report. A claim is a factual statement or opinion that expresses a belief or judgment.

List each claim on a new line, starting with “Claims: ”.

Table 7: Prompt used for claim decomposition in *ConflictScore*.

Prompt:

Document-Claim Conflict Detection

You are a careful fact-checking assistant specializing in identifying conflicts between claims and evidence in documents.

Decide one of the following labels for how the document relates to the claim:

1. **SUPPORTS** – The document provides evidence that directly supports any part of the claim. This includes partial support: if a claim contains hedged or multi-part statements (e.g., “Coffee consumption can improve alertness but may also cause sleep disruption”), then documents that support any part of that claim should be labeled SUPPORTS.
2. **CONTRADICTS** – The document provides evidence that contradicts the claim (e.g., assigns an incompatible role/date/quantity/polarity, or states the opposite of an asserted relationship), even if it does not explicitly say the claim is false.
3. **IRRELEVANT** – The document does not provide information about the substantive content of the claim (mentions entities without giving information that could support or contradict) OR provides only vague/contextual information insufficient to judge the claim.

Decision rules:

- Look for statements in the document that address the same attributes asserted (who/what/when/where/how many/etc.) inside the claim.
- If the document contains a statement that supports the claim or partially supports any part of a hedged claim, choose SUPPORTS.
- If the document contains a statement incompatible with the claim, choose CONTRADICTS (do NOT choose IRRELEVANT).
- If it contains neither a compatible nor incompatible statement about the claim’s asserted attributes, choose IRRELEVANT.

Example (contradiction via incompatible attribute):

- Claim: “Frédéric Chopin was a famous musician.”
- Document: “Polish scientists born in Warsaw include ... Frédéric Chopin.”
- Label: CONTRADICTS (profession mismatch: scientist vs musician).

Now evaluate the following:

Claim to evaluate

{claim}

Document to evaluate against

{document}

Output format

Return ONLY a single JSON object with these fields:

```
{
  "claim": "{claim}",
  "document_snippet": "<the most relevant snippet from the document that relates to the claim>",
  "reasoning": "<a brief justification that explains your decision>",
  "answer": "<SUPPORTS|CONTRADICTS|IRRELEVANT>"
}
```

Please be precise and follow the decision rules. Do not include any additional text or explanations and only output in the JSON format specified above.

Table 8: Prompt used for evidence evaluation in *ConflictScore*.

Prompt:

Instruction: Write a report about {main_entity} given the following documents. Make your report concise and not longer than 3 paragraphs.

Documents:

{Documents}

Table 9: Prompt used for report generation in the retrieval-augmented (RAG) setting.

Prompt:

Instruction: Given the following documents, write a report about {main_entity}.

Please be cautious and hedge accordingly. Think through all the information and consider all perspectives if they exist to acknowledge any potential conflicts or different viewpoints on this topic.

Documents:

{Documents}

Make your report concise and not longer than 3 paragraphs.

Table 10: Prompt used for balanced report generation in the retrieval-augmented (RAG-Balanced) setting.

Prompt:

Instruction: You are a careful and balanced reporting assistant. Given the following documents, write a 3-paragraph report about {main_entity}.

Your goal is to synthesize the information from the documents while minimizing factual conflicts. To do this:

- Hedge appropriately when the evidence is partial, uncertain, or varies across documents. Use phrases like “some sources suggest,” “it is possible that,” “there is limited evidence,” or “while one document indicates...”.

- Avoid making definitive claims unless they are clearly and consistently supported across multiple documents.

- Do not speculate or infer beyond what is stated in the documents.

- Attribute information to the documents when relevant (e.g., “According to one article...” or “Document A notes that...”).

- Acknowledge conflicting viewpoints if they appear in the documents.

Documents:

{Documents}

Table 11: Prompt used for super-balanced report generation in the retrieval-augmented (RAG-Super-Balanced) setting.

Prompt:

Given the following documents retrieved from the web, answer the question by choosing either (A) or (B).

Documents:
{context_text}

Question: {question}

(A) {option_a}
(B) {option_b}

Please select either (A) or (B) as your answer. Respond with only the letter in parentheses, for example: (A) or (B).
Answer:

Table 12: Prompt used for multiple-choice question answering in the TruthfulQA RAG setting.

Prompt:

Given the following documents retrieved from the web, answer the question by choosing either (A) or (B). Pay careful attention to any conflicts or contradictions between the documents. Think critically about the information provided and prioritize accuracy and truthfulness in your response. If there are conflicting claims, think carefully about the information provided and choose the most reliable answer.
{context_text}

Question: {question}

(A) {option_a}
(B) {option_b}

Please select either (A) or (B) as your answer. Respond with only the letter in parentheses, for example: (A) or (B).
Answer:

Table 13: Prompt used for conflict-aware multiple-choice question answering in the TruthfulQA Control-RAG setting.

Prompt:

You are being asked to reconsider and improve your previous binary choice response to a question. Your original response contained conflicting information based on the provided documents.

ORIGINAL QUESTION: {question}

DOCUMENTS:
{context_text}

YOUR ORIGINAL RESPONSE:
{original_response}

YOUR MAIN CLAIM THAT HAD CONFLICTS:
{selected_claim}
{conflict_details}

INSTRUCTIONS:

1. Carefully re-examine all the provided documents.
2. Pay special attention to the conflicting evidence identified above.
3. Think critically about which sources are most reliable and authoritative.
4. Consider the quality, recency, and credibility of different sources.
5. Prioritize accuracy and truthfulness over simply agreeing with the majority of sources.

OUTPUT FORMAT (strict):

1) First line: exactly one of:

- A – {option_a}
- B – {option_b}

2) Second line: Justification: one concise sentence citing document numbers only in parentheses, e.g., (Docs 2, 5).

Do not include chain-of-thought, lists, or extra sections. Return only the two lines described above.

Table 14: Prompt used for conflict-aware response re-generation in the multiple-choice TruthfulQA RAG setting.