

FACTCORRECTOR: A Graph-Inspired Approach to Long-Form Factuality Correction of Large Language Models

Javier Carnerero-Cano¹, Massimiliano Pronesti¹, Radu Marinescu¹, Tigran Tchrakian¹, James Barry¹, Jasmina Gajcin¹, Yufang Hou^{1,2}, Alessandra Pascale¹, Elizabeth Daly¹

¹IBM Research Europe - Ireland

²IT:U - Interdisciplinary Transformation University Austria

{javier.cano, massimiliano.pronesti, james.barry, jasmina.gajcin2, yufang.hou1}@ibm.com

{radu.marinescu, tigran, apascale, elizabeth.daly}@ie.ibm.com

Abstract

Large language models (LLMs) are widely used in knowledge-intensive applications but often generate factually incorrect responses. A promising approach to rectify these flaws is correcting LLMs using feedback. Therefore, in this paper, we introduce FACTCORRECTOR, a new post-hoc correction method that adapts across domains without retraining and leverages structured feedback about the factuality of the original response to generate a correction. To support rigorous evaluations of factuality correction methods, we also develop the VELI5 benchmark, a novel dataset containing systematically injected factual errors and ground-truth corrections. Experiments on VELI5 and several popular long-form factuality datasets show that the FACTCORRECTOR approach significantly improves factual precision while preserving relevance, outperforming strong baselines. We release our code at <https://ibm.biz/factcorrector>.

1 Introduction

Large Language Models (LLMs) have demonstrated remarkable capabilities in generating coherent and contextually relevant text across diverse domains (Brown et al., 2020; Chowdhery et al., 2023). However, their tendency to produce factually incorrect or hallucinated content remains a significant barrier to safe and reliable deployment in real-world applications (Zhang et al., 2023; Sahoo et al., 2024; Huang et al., 2025), where knowledge conflicts (Xu et al., 2024; Hou et al., 2024) are also common. Factuality is critical in high-stakes domains such as healthcare or finance, where even minor inaccuracies can lead to harmful decisions or financial loss (Tonmoy et al., 2024).

Prior research has explored several strategies to mitigate factual errors in LLM outputs. Specifi-

cally, training-time correction methods optimize the model behavior to avoid factual errors during learning by incorporating human or automated feedback. Approaches include direct optimization with human feedback (Glaese et al., 2022), reward modeling and RLHF (Ouyang et al., 2022; Zhong et al., 2025), and self-training with automated signals (Dubois et al., 2024). However, these strategies are often infeasible for closed-source or large scale models (Pan et al., 2024). Generation-time correction methods refine LLM outputs during decoding by leveraging critic models for guidance. Approaches such as generate-then-rank sample multiple candidate responses and select the best using a ranking model (Weng et al., 2023), while feedback-guided decoding integrates step-level feedback to steer generation in real time (Yao et al., 2023). The effectiveness of these methods hinges on the critic’s ability to provide high-quality intermediate feedback. In contrast, post-hoc correction methods evaluate the final LLM response to construct the feedback that is subsequently used to refine and correct the response. The feedback can be provided through self-critique (Shinn et al., 2023), external critics (e.g., RAC (Li and Flanagan, 2025), CRITIC (Gou et al., 2024)) or agentic debate (Cohen et al., 2023). However, current approaches are limited in that they accept feedback in a simple text format, which may miss the complex structure of relationships between statements to verify and their related evidence.

Contribution. In this paper, we present FACTCORRECTOR, a novel post-hoc method for long-form factuality correction that leverages structured feedback to produce factually accurate responses. The approach begins by applying a critic model to the generated output to assess its factuality.

This critic, built upon FACTREASONER—a recent long-form factuality assessor (Marinescu et al., 2025)—decomposes the response into atomic units, each representing a single fact or claim. For each atom, it retrieves supporting or contradicting evidence from external sources and estimates its truthfulness using a graphical model that captures entailment and contradiction relationships between atoms and contexts. The resulting feedback identifies the false atoms along with the evidence that refutes them. Subsequently, a refinement model integrates this feedback with the original response into a structured prompt to generate a corrected version of the initial response, ensuring that the model focuses on the parts of the response that were flagged as incorrect by the critic. The approach can be applied in a single pass or iteratively until the desired factuality accuracy is achieved.

We also introduce the VELI5 benchmark, a dataset specifically designed to enable rigorous, large-scale evaluation of long-form factuality correction methods and to support the training of specialized models for correcting LLM-generated responses. VELI5 builds upon the widely used ELI5 dataset (Fan et al., 2019), comprising curated pairs of human-authored questions and answers that have been explicitly verified and corrected for factual inaccuracies using FACTCORRECTOR.

Finally, we present an extensive empirical evaluation with our VELI5 and several popular long-form factuality datasets using a wide range of open-source LLMs. The results show clearly that FACTCORRECTOR consistently delivers the most reliable improvements in factuality across all datasets thus outperforming strong baselines and demonstrating strong generalization across models.

2 Related Work

Although LLMs have achieved remarkable performance across a wide range of tasks, they still struggle with hallucinations, unfaithful reasoning, or the propagation of bias and toxicity. In recent years, several techniques have been developed to address these issues by leveraging automated feedback – either generated by the LLM itself or provided by external systems – to refine and *correct* model outputs (Pan et al., 2024). These correction strategies can be broadly categorized into three types: training-time, generation-time, and post-hoc approaches.

Training-time correction strategies aim to improve

model behavior during the training phase by incorporating human feedback, reward models, or automated feedback. Common approaches include direct optimization using human feedback e.g., (Glaese et al., 2022; Scheurer et al., 2024; Chen et al., 2024; Liu et al., 2023; Gao et al., 2023), reward modeling and reinforcement learning from human feedback (RLHF) e.g., (Ouyang et al., 2022; Bai et al., 2022; Ganguli et al., 2023), and self-training with automated feedback e.g., (Zelikman et al., 2022; Bai et al., 2022; Dubois et al., 2024). However, these methods are often impractical for closed-source models or extremely large-scale models with billions of parameters.

Generation-time correction methods, such as generate-then-rank and feedback-guided decoding, aim to refine the output of large language models (LLMs) during the generation process itself. In generate-then-rank approaches, a large pool of candidate outputs is first sampled, and then ranked using a secondary critic model to select the most appropriate response (Weng et al., 2023; He et al., 2022; Ni et al., 2023; Chen et al., 2022). In contrast, feedback-guided decoding incorporates a step-level critic model that provides real-time feedback during generation, enabling more fine-grained control. This strategy underpins methods like Tree-of-Thought (Yao et al., 2023), GRACE (Khalifa et al., 2023), and RAP (Hao et al., 2023). These approaches primarily differ in the type of critic model employed—ranging from reward models trained with human feedback, to verifiers, external metrics or external knowledge sources.

The effectiveness of generation-time correction depends on the critic model’s ability to deliver high-quality feedback on intermediate outputs. In contrast, *post-hoc correction* methods apply both the critic and refinement models after the full output is generated, enabling richer and more varied feedback—from targeted diagnostics to general writing suggestions. These strategies are typically categorized into: self-correction, correction with external feedback and multi-agent debate. More specifically, in self-correction methods such as Self-Refine (Madaan et al., 2023), Self-Verification (Gero et al., 2023), or Reflexion (Shinn et al., 2023) a single LLM both generates and refines its output iteratively until an acceptable quality of the output is obtained. Additionally, multiple external tools can be leveraged to provide improved feedback, thus leading to methods such as RAC (Li and Flani-

gan, 2025), CRITIC (Gou et al., 2024), FACTOOL (Chern et al., 2023), REFINER (Paul et al., 2024) and others (Pan et al., 2024).

Our FACTCORRECTOR is a post-hoc factuality correction method, placing it in the same category as approaches like RAC and CRITIC. However, it distinguishes itself through its integration with FACTREASONER, a recently introduced pipeline designed for long-form factuality assessment (Marinescu et al., 2025) which serves as the critic model. Furthermore, the refinement process in FACTCORRECTOR is guided by the probabilistic graphical model derived from FACTREASONER.

3 Preliminaries

We begin by providing background on long-form factuality assessment for LLMs.

3.1 Long-Form Factuality Assessment

Let y be the long-form response generated by an LLM to a query x . We assume that y can be decomposed into a set of n atomic units (or atoms) that can be either true or false, denoted by $\mathcal{A}_y = \{a_1, a_2, \dots, a_n\}$ (Min et al., 2023; Song et al., 2024; Wei et al., 2024). An atomic unit $a_i \in \mathcal{A}_y$ is defined as a short sentence conveying one piece of information (e.g., a claim or a fact). Furthermore, given an external knowledge source \mathcal{K} ¹, we say that an atomic unit $a_i \in \mathcal{A}_y$ is *supported* by \mathcal{K} if there exists at least one piece of information in \mathcal{K} (e.g., a passage) called a *context* that undebatably supports a_i . Otherwise, the atomic unit is *not supported*.

The *factual precision* denoted by $Pr(y)$ of the response y with respect to a knowledge source \mathcal{K} is defined as: $Pr(y) = \frac{S(y)}{|\mathcal{A}_y|}$, where $S(y) = \sum_{i=1}^n \mathbb{I}[a_i \text{ is supported by } \mathcal{K}]$ is the number of supported atomic units. The *factual recall* $R_K(y)$ up to the K -th supported atomic unit is given by: $R_K(y) = \min(\frac{S(y)}{K}, 1)$. Finally, an F_1 measure for long-form factuality denoted by $F_1@K$ can be defined as: $F_1@K(y) = \frac{2 \cdot Pr(y) \cdot R_K(y)}{Pr(y) + R_K(y)}$ if $S(y) > 0$, and 0 otherwise (Wei et al., 2024).

3.2 The FACTREASONER Factuality Assessor

FACTREASONER (Marinescu et al., 2025) is a recent factuality assessor that, unlike previous prompt-based approaches (Min et al., 2023; Song

et al., 2024; Wei et al., 2024), leverages probabilistic reasoning to assess the factuality of the LLM generated response with respect to an external knowledge source \mathcal{K} .

Specifically, the FACTREASONER pipeline shown in Figure 5 (top) consists of four stages called *Atomizer*, *Reviser*, *Retriever* and *Evaluator*, respectively. The *Atomizer* decomposes the input response y into a set of n atomic units \mathcal{A}_y by applying any of the decomposition strategies proposed recently (Min et al., 2023; Bayat et al., 2025). The *Reviser* post-processes the atoms such that the pronouns, unknown entities, or incomplete names are replaced with their corresponding named entities in the response (Wei et al., 2024). Next, the *Retriever* is responsible for querying the external knowledge source \mathcal{K} to retrieve the contexts relevant to the response’s atoms (Song et al., 2024). Finally, the *Evaluator* constructs a graphical model representing a joint probability distribution over the atomic units in the response and their corresponding contexts in \mathcal{K} . For each atom a_i , it then computes the posterior marginal probability $P(a_i)$ which quantifies the likelihood that a_i is true (or supported) given the information available in \mathcal{K} . Finally, the factuality metrics defined previously, such as $Pr(y)$, can be readily calculated using the supported atoms. Note that both the *Atomizer* and *Reviser* stages use LLMs for their specific tasks, as previously shown in (Marinescu et al., 2025).

Graphical Model. FACTREASONER’s graphical model \mathcal{G} is defined by a tuple $\langle \mathbf{X}, \mathbf{D}, \mathbf{F} \rangle$ where \mathbf{X} is a set of Boolean variables A_i or C_j associated with each atom $a_i \in \mathcal{A}_y$ or context $c_j \in \mathcal{C}_y$, respectively, where \mathcal{C}_y is the set of retrieved contexts. The domains \mathbf{D} are $\{\text{true}, \text{false}\}$ for all variables (e.g., we denote a_i and $\neg a_i$ for $A_i = \text{true}$ and $A_i = \text{false}$, respectively) and \mathbf{F} is a set of factors. Each variable has a unary factor $f(\cdot)$ in \mathbf{F} encoding prior belief. Atoms have uniform priors: $f(a_i) = f(\neg a_i) = 0.5$. Contexts from reliable sources have high prior (e.g., $f(c_j) = 0.99$), while less reliable sources use smaller values. Binary factors $f(C_j, A_i)$ and $f(C_j, C_k)$ capture probabilistic logical relations between utterances. A relation model $p_\theta(\cdot|t, t')$ (either a specialized BERT model or an LLM) predicts the most likely relation from $\{\text{neutral}, \text{entail}, \text{contradict}\}$ together with its probability. Therefore, each factor $f(C_j, A_i)$ or $f(C_j, C_k)$ represents a probabilistic encoding of an entailment relationship (i.e., context C_j entails

¹For example, \mathcal{K} could be Wikipedia, Google Search, or a collection of documents embedded into a vector database.

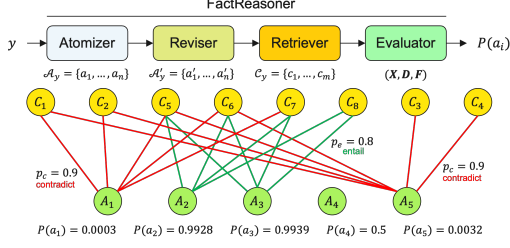


Figure 1: FACTREASONER pipeline and an example graphical model with 5 atoms and 8 context variables.

Algorithm 1: FACTCORRECTOR

Input: response y , refinement model \mathcal{R} , threshold θ
Output: corrected response y
 Evaluate factual precision $Pr(y)$ of y
while $Pr(y) < \theta$ **do**
 Decompose y into atoms $\mathcal{A}_y = \{a_1, \dots, a_n\}$
 Revise the atoms $\mathcal{A}_y = \{a_1, \dots, a_n\}$
 Retrieve contexts $\mathcal{C}_y = \{c_1, \dots, c_m\}$ from \mathcal{K}
 Use \mathcal{A}_y and \mathcal{C}_y to build $\mathcal{G} = \langle \mathbf{X}, \mathbf{D}, \mathbf{F} \rangle$
 for atom $a_i \in \mathcal{A}_y$ **do**
 Evaluate posterior probability $P(a_i)$ in \mathcal{G}
 if $P(a_i) > 0.5$ **then** Label a_i as True
 else if $P(a_i) < 0.5$ **then** Label a_i as False
 else Label a_i as Unverified
 $feedback \leftarrow \{a_i | a_i \text{ is False or Unverified} \} \cup \{c_j | \exists a_i \in \mathcal{G} \text{ s.t. } c_j \text{ connects to } a_i \text{ in } \mathcal{G} \text{ and } a_i \text{ is False} \}$
 $y' \leftarrow \mathcal{R}(y, feedback)$
 Evaluate factual precision $Pr(y')$ of y'
 if $Pr(y') > Pr(y)$ **then**
 Let $y \leftarrow y'$ and $Pr(y) \leftarrow Pr(y')$
 else break
return y

or supports atom A_i) or a contradiction relationship (i.e., context C_j contradicts atom A_i), respectively. Note that the neutral relationships are ignored.

Example 1. Figure 5 (bottom) shows an example graphical model with 5 atoms (A_1, \dots, A_5), 8 contexts (C_1, \dots, C_8) and 19 binary relations between atoms and contexts. Each edge in the graph is labeled by the probability of the corresponding entailment (green) or contradiction (red) relation.

4 The FACTCORRECTOR Pipeline

In this section, we present FACTCORRECTOR, a novel post-hoc method for correcting factual inaccuracies in long-form responses generated by LLMs. As a post-hoc correction strategy, FACTCORRECTOR employs a *refinement model* \mathcal{R} to revise an initial response y to a user query x , guided by structured feedback from a *critic model* \mathcal{C} that evaluates the factual quality of the response. Specifically, the critic model in FACTCORRECTOR builds on the FACTREASONER factuality assessor, which systematically identifies the incorrect parts within

the response, while the refinement model can be any instruction-following LLM.

Algorithm 1 describes the main steps of the proposed method which operate iteratively, starting from the initial response y . At each iteration, the critic decomposes the response into atomic units and retrieves relevant evidence from an external knowledge source \mathcal{K} . These atoms and contexts are subsequently used to build a graphical model \mathcal{G} that captures entailment and contradiction relationships between atoms and retrieved contexts. Then, each atom a_i is labeled as True, False, or Unverified if the posterior probability $P(a_i)$ is greater than, less than or equal to 0.5, respectively. The feedback includes all atoms labeled as False or Unverified, along with any contexts connected to them in \mathcal{G} and the corresponding relationship types, ensuring that corrections are grounded in explicit factual signals.

The refinement model then integrates this feedback with the original response in a structured prompt to produce a revised output. The process repeats until the corrected response meets a predefined factual precision threshold θ or there is no improvement.

Example 2. Figure 2 illustrates the FACTCORRECTOR workflow on a simple question–response pair. In this example, the initial response has a factual precision of 0.42 and is decomposed into five atoms such that a_2 and a_3 are labeled as True, a_1 and a_5 as False, and a_4 as Unverified, respectively. The corresponding graphical model is shown in Figure 5 (bottom). The feedback includes the False atoms $\{a_1, a_5\}$ along with their contradicting and entailing contexts $\{c_1, c_2, \dots, c_7\}$. Atom a_4 is also flagged for removal. The refinement model (Corrector) uses this feedback to generate a corrected response that omits unverified claims and aims to correct the False atoms. The resulting output contains only one incorrect claim (highlighted in red) and achieves a factual precision of 0.93.

FACTCORRECTOR is closely related to the recent RAC method (Li and Flanagan, 2025), but with a key difference: while RAC corrects each atom using only its retrieved documents, our approach uses the graphical model representing the atom–context relationships and aggregates all contexts connected to an incorrect atom into its feedback. This holistic view enables the refinement model to make more accurate corrections, reflected in the superior performance observed in our experiments.

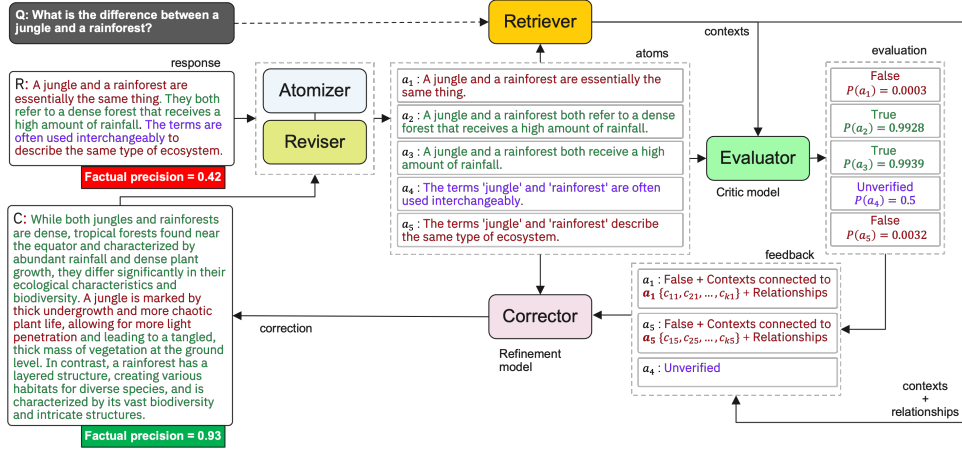


Figure 2: Illustration of the FACTCORRECTOR pipeline for long-form factuality correction of LLMs.

5 The VELI5 Benchmark Dataset

In this section, we introduce the VELI5 benchmark, a dataset specifically constructed to enable rigorous evaluation of long-form factuality correction methods and to facilitate the training of specialized models for correcting LLM generated responses.

The starting point for the VELI5 benchmark dataset is the ELI5-Category dataset², a curated subset of ELI5 (Fan et al., 2019) that augments long-form explanatory question-answer threads scrapped from the *r/explainlikeimfive* reddit forum with explicit topical annotations. This resource contains questions in which users request intuitive explanations of complex topics, each assigned by community moderators to one of 12 high-level categories (11 topical domains plus a *Repost* category) and paired with multiple candidate answers and their corresponding upvote scores.

For each question, we deterministically select the answer with the highest number of upvotes as the *canonical response*. While these answers are generally high quality and well articulated, they are not guaranteed to be factually correct and may include inaccuracies, outdated claims, or speculative statements. Addressing this gap between explanatory quality and factual reliability is the primary motivation behind VELI5. We therefore apply FACTCORRECTOR to each canonical response.

To further diversify the dataset, we intentionally introduce factually incorrect responses to user questions. These synthetic responses are generated by prompting a reasonably strong LLM, such as the mixtral-8x22b-instruct model. The ratio of

synthetic to human-authored responses is maintained at 50%, ensuring a balanced mix of realistic and adversarial content. The final VELI5 dataset comprises 17,522 instances, uniformly distributed across 12 categories. We further partition the dataset into training (14,017 samples), validation (1,752 samples), and test (1,753 samples) splits. Illustrative examples and additional details are provided in Appendix E.

6 Experiments

We empirically evaluate the FACTCORRECTOR (FC) pipeline for long-form factuality correction and compare it against state-of-the-art approaches on several popular long-form factuality datasets.

6.1 Baseline Correctors

We consider two recent state-of-the-art post-hoc correction methods with feedback: CRITIC (Gou et al., 2024) and RAC (Li and Flanagan, 2025). CRITIC is a prompt-based iterative approach that leverages LLMs and retrieved knowledge. At each step, the model generates a query to retrieve relevant information, revises its output based on revision history, and selects the most plausible answer. RAC, in contrast, decomposes the generated response into atomic facts, retrieves evidence from trusted sources (e.g., Google), verifies each fact, and corrects inaccuracies using reliable content. In addition, we also evaluate two prompting strategies: LLM1, which performs corrections using only the model’s internal knowledge, and LLM2, which uses contexts retrieved for the question only and ignores its internal knowledge (see Appendix D).

We instantiated the competing correctors with

²https://huggingface.co/datasets/rexarski/eli5_category

open-source LLMs from the IBM Granite (Research, 2024), Meta LLaMA (Touvron et al., 2023), MistralAI Mixtral (Jiang et al., 2024) and OpenAI (OpenAI et al., 2025) families, spanning a broad range of architectures and scales. All models run remotely on A100 80GB GPUs and are accessed via litellm APIs (1,500 prompts per sec).

6.2 Datasets

We use three datasets in our experiments: a reduced version of VELI5 (200 instances uniformly sampled from the test split), Biographies (BIO) (Min et al., 2023), and AskHistorians (ASKHIST) (Xu et al., 2023). The BIO dataset includes 183 biographical passages generated by ChatGPT for entities with Wikipedia pages. ASKHIST contains 200 questions from the r/AskHistorians reddit forum paired with long-form answers produced by the llama-3.3-70b-instruct model.

We also use the CONFLICTS dataset from (Marinescu et al., 2025), which comprises 100 atomic claims (or responses) sampled from Conflict-Bank (Su et al., 2024). Each claim, originally from Wikidata, is assumed true and is associated with both supporting and conflicting contexts. Since all claims are correct, they require no corrections, offering a controlled setting for our evaluation.

6.3 Measures of Performance

For each dataset \mathcal{D} and each competing corrector, we computed three factuality metrics: precision (Pr), recall at K ($R@K$), and $F_1@K$, averaged over all prompts in \mathcal{D} , where K is set to the median number of atoms. These metrics are evaluated using the FACTREASONER (FR) assessor with Google search results as external knowledge source (Marinescu et al., 2025). In addition, we computed two complementary metrics: verifiability (V) and comprehensiveness (C). Verifiability V is defined as the number of atoms in the response (or correction) that can be verified—i.e., atoms connected to either supporting or contradicting evidence in FR’s graph. Comprehensiveness C quantifies coverage and is given by $C = \frac{|\mathcal{A}_{in}|}{|\mathcal{A}_{in}|+|\mathcal{A}_{out}|}$, where \mathcal{A}_{in} denotes the atoms covered by the response (or correction), and \mathcal{A}_{out} denotes atoms that are *uncovered* or missing (Dejl et al., 2025).

For each factuality metric S (e.g., precision), we report its *relative gain*, denoted as $G(S)$ and defined by: $G(S) = \frac{2 \cdot (S_c - S_r)}{S_c + S_r}$, where S_r and S_c are the metrics corresponding to the original response and

| corrector | ROUGE \uparrow | BLEU \uparrow | BLEURT \uparrow | JUDGE \uparrow |
|------------------------|---------------------------------|---------------------------------|---------------------------------|------------------|
| mixtral-8x22b-instruct | | | | |
| CRITIC | 0.14 \pm 0.07 | 0.02 \pm 0.03 | -0.69 \pm 0.25 | 0.41 |
| RAC | 0.54 \pm 0.23 | 0.23 \pm 0.27 | 0.29 \pm 0.43 | 0.87 |
| LLM1 | 0.15 \pm 0.05 | 0.02 \pm 0.02 | -0.63 \pm 0.21 | 0.30 |
| LLM2 | 0.20 \pm 0.10 | 0.05 \pm 0.05 | -0.48 \pm 0.22 | 0.10 |
| FC (ours) | 0.89\pm0.26 | 0.87\pm0.32 | 0.73\pm0.47 | 0.87 |
| llama-3.3-70b-instruct | | | | |
| CRITIC | 0.32 \pm 0.27 | 0.15 \pm 0.28 | -0.36 \pm 0.56 | 0.31 |
| RAC | 0.87\pm0.20 | 0.77\pm0.32 | 0.65\pm0.43 | 0.77 |
| LLM1 | 0.19 \pm 0.11 | 0.03 \pm 0.05 | -0.62 \pm 0.33 | 0.17 |
| LLM2 | 0.21 \pm 0.10 | 0.04 \pm 0.04 | -0.46 \pm 0.22 | 0.06 |
| FC (ours) | <u>0.73\pm0.33</u> | <u>0.62\pm0.45</u> | <u>0.46\pm0.56</u> | <u>0.60</u> |
| granite-4.0-h-small | | | | |
| CRITIC | 0.15 \pm 0.11 | 0.03 \pm 0.10 | -0.69 \pm 0.26 | 0.20 |
| RAC | 0.74 \pm 0.23 | 0.48 \pm 0.37 | 0.49 \pm 0.46 | <u>0.73</u> |
| LLM1 | 0.15 \pm 0.07 | 0.02 \pm 0.02 | -0.68 \pm 0.23 | 0.17 |
| LLM2 | 0.42 \pm 0.24 | 0.15 \pm 0.17 | -0.10 \pm 0.44 | 0.13 |
| FC (ours) | 0.83\pm0.30 | 0.76\pm0.39 | 0.58\pm0.57 | 0.78 |

Table 1: Results obtained on the CONFLICTS dataset.

the correction, respectively. A positive $G(S)$ indicates that the correction outperforms the response, while a negative value means that the correction performs worse. By construction, $G(S)$ ranges from -2 to 2 and remains well defined even when either S_r or S_c equals zero.

6.4 Results for Post-hoc Correction

We next present the results obtained on our datasets using a range of LLMs. For consistency, the same LLM employed by FACTREASONER to evaluate the metrics reported in the tables was also used to instantiate the components of FACTCORRECTOR.

CONFLICTS. Table 1 summarizes the results obtained on our controlled experiment with the CONFLICTS dataset. In this case, we assess the similarity between correction and response using the ROUGE, BLEU, and BLEURT metrics (Lin, 2004; Papineni et al., 2002; Sellam et al., 2020) because the response is known to be true and doesn’t need correction. We also report JUDGE, which counts the number of instances where an LLM-as-a-Judge model (in our case DeepSeek-3.2), prompted appropriately, infers that the correction is equivalent to the original response. The table shows the mean and standard deviation for all metrics. The results indicate that FACTCORRECTOR achieves superior performance with the Mixtral and Granite models, while ranking second when using LLaMA, where RAC attains the best scores. In contrast, the LLM1 and LLM2 baselines perform rather poorly almost always producing corrections significantly different

| corrector | Pr \uparrow | R@K \uparrow | F1@K \uparrow | V \uparrow | C \uparrow |
|------------------------|---------------|----------------|-----------------|--------------|--------------|
| mixtral-8x22b-instruct | | | | | |
| CRITIC | 0.26 | 0.08 | 0.19 | 0.09 | 0.06 |
| RAC | 0.24 | 0.08 | 0.18 | 0.08 | 0.06 |
| LLM1 | 0.32 | 0.10 | <u>0.23</u> | 0.12 | 0.07 |
| LLM2 | <u>0.32</u> | 0.10 | <u>0.23</u> | <u>0.13</u> | 0.07 |
| FC (ours) | 0.34 | <u>0.09</u> | 0.24 | 0.14 | 0.07 |
| llama-3.3-70b-instruct | | | | | |
| CRITIC | 0.19 | 0.06 | 0.13 | 0.03 | 0.05 |
| RAC | 0.13 | 0.02 | 0.09 | 0.02 | 0.05 |
| LLM1 | 0.25 | <u>0.06</u> | 0.18 | 0.05 | 0.08 |
| LLM2 | <u>0.26</u> | 0.07 | 0.18 | 0.05 | 0.08 |
| FC (ours) | 0.27 | 0.07 | 0.18 | 0.05 | 0.08 |
| granite-4.0-h-small | | | | | |
| CRITIC | 0.17 | 0.04 | 0.12 | 0.05 | 0.06 |
| RAC | 0.20 | 0.09 | 0.16 | 0.06 | <u>0.07</u> |
| LLM1 | 0.31 | 0.11 | 0.23 | 0.11 | 0.08 |
| LLM2 | 0.20 | <u>0.10</u> | 0.16 | <u>0.09</u> | 0.05 |
| FC (ours) | <u>0.29</u> | <u>0.10</u> | <u>0.21</u> | 0.11 | <u>0.07</u> |
| gpt-oss-120b | | | | | |
| CRITIC | -0.04 | 0.01 | -0.02 | -0.01 | -0.03 |
| RAC | 0.34 | 0.15 | 0.27 | <u>0.06</u> | <u>0.13</u> |
| LLM1 | 0.21 | 0.16 | 0.19 | -0.01 | 0.13 |
| LLM2 | 0.33 | 0.20 | <u>0.28</u> | 0.05 | <u>0.13</u> |
| FC (ours) | 0.36 | <u>0.19</u> | 0.30 | 0.07 | 0.14 |

Table 2: Mean relative gains for the factuality metrics obtained on the VELI5 dataset.

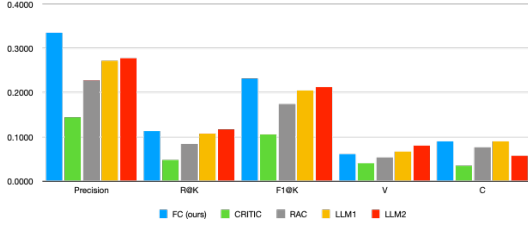


Figure 3: Mean relative gains for factuality metrics across models on the VELI5 dataset.

from the response, while CRITIC often struggles to follow its instructions also leading to poor results.

VELI5. In Table 2 we report the mean relative gains obtained for the factuality metrics on the VELI5 dataset. The results demonstrate that FACT-CORRECTOR consistently delivers strong and reliable improvements in factuality across all evaluated metrics and models. For example, when considering relative gains on precision, FC outperforms all baselines, achieving macro-averaged gains of 0.315 — substantially higher than the next-best corrector (see also Figure 3). Importantly, FC is often best or tied-best in many metric–model combinations, showing strong generalization across diverse LLM architectures. LLM1 and LLM2 show strong performance in this case, likely due to the provenance of the VELI5 questions which were originally derived from a public Reddit forum. Therefore, it is

| corrector | llama-3.3-70b | | mixtral-8x22b | | granite-4.0-small | | gpt-oss-120b | |
|-----------|---------------|-------------|---------------|-------------|-------------------|-------------|--------------|-------------|
| | before | after | before | after | before | after | before | after |
| CRITIC | 0.78 | 0.90 | 0.70 | 0.88 | 0.74 | 0.86 | 0.62 | 0.58 |
| RAC | 0.78 | 0.91 | 0.70 | 0.86 | 0.74 | 0.88 | 0.62 | 0.85 |
| LLM1 | 0.78 | 0.95 | 0.70 | <u>0.91</u> | 0.74 | 0.94 | 0.62 | 0.72 |
| LLM2 | 0.78 | <u>0.96</u> | 0.70 | <u>0.91</u> | 0.74 | 0.86 | 0.62 | 0.79 |
| FC (ours) | 0.78 | 0.97 | 0.70 | 0.93 | 0.74 | <u>0.93</u> | 0.62 | <u>0.84</u> |

Table 3: Mean factual precision before and after correction on the VELI5 dataset.

plausible that similar content was included in these models’ pre-training corpora which may confer a significant advantage on this dataset.

Table 3 reinforces these findings by examining factual precision before and after correction. FC achieves the highest post-correction precision on two out of four models (Llama, Mixtral) and ranks second on Granite and GPT-OSS, where it still delivers large relative gains. On average, FC improves precision by +0.21, outperforming all baselines, including LLM1 (+0.17) and RAC (+0.16). While RAC slightly surpasses FC on GPT-OSS in final precision, FC dominates in overall gains, suggesting a more robust correction strategy. In contrast, CRITIC exhibits inconsistent behavior, even reducing precision for GPT-OSS.

BIO. Table 4 presents the mean relative gains on the BIO dataset. FC exhibits robust performance across models, achieving its highest gains with Mixtral and Granite, and ranking second for Llama. In contrast, LLM1 performs poorly, particularly on Granite and GPT-OSS. Both RAC and LLM2 remain highly competitive on this dataset, while CRITIC, as observed previously, struggles to follow instructions, leading to inconsistent results.

ASKHIST. Table 5 shows the mean relative gains on the ASKHIST dataset. Unlike previous cases, the gains here are smaller. This is primarily because the factual precision of the original responses is already high (typically exceeding 89% on average) leaving limited room for improvement through correction. Interestingly, all other baselines exhibit negative gains, indicating that their corrections are less factual than the original responses. In contrast, FC maintains robust performance across all models, almost always producing corrections that are at least as factual as the original responses and, in many cases, slightly better.

6.5 Results for SFT Correction

We used the larger VELI5 dataset to train a LoRA adapter on the Granite-Guardian-5B model for gen-

| corrector | Pr \uparrow | R@K \uparrow | F1@K \uparrow | V \uparrow | C \uparrow |
|------------------------|---------------|----------------|-----------------|--------------|--------------|
| mixtral-8x22b-instruct | | | | | |
| CRITIC | 0.27 | 0.22 | 0.25 | 0.05 | 0.05 |
| RAC | 0.39 | 0.29 | 0.34 | 0.05 | 0.13 |
| LLM1 | 0.27 | 0.25 | 0.27 | 0.02 | 0.07 |
| LLM2 | 0.42 | 0.30 | 0.37 | 0.11 | 0.07 |
| FC (ours) | 0.46 | 0.33 | 0.41 | 0.12 | 0.11 |
| llama-3.3-70b-instruct | | | | | |
| CRITIC | 0.37 | 0.18 | 0.29 | 0.08 | 0.13 |
| RAC | 0.42 | 0.21 | 0.33 | 0.12 | 0.16 |
| LLM1 | 0.23 | 0.14 | 0.19 | 0.07 | 0.11 |
| LLM2 | 0.45 | <u>0.21</u> | 0.36 | 0.12 | 0.18 |
| FC (ours) | 0.44 | 0.22 | <u>0.35</u> | 0.12 | <u>0.17</u> |
| granite-4.0-h-small | | | | | |
| CRITIC | 0.12 | 0.02 | 0.08 | 0.04 | 0.02 |
| RAC | <u>0.22</u> | 0.08 | 0.16 | 0.10 | 0.04 |
| LLM1 | -0.03 | -0.01 | -0.02 | 0.02 | 0.00 |
| LLM2 | 0.19 | 0.13 | <u>0.16</u> | 0.13 | 0.03 |
| FC (ours) | 0.26 | <u>0.10</u> | 0.20 | <u>0.14</u> | 0.03 |
| gpt-oss-120b | | | | | |
| CRITIC | 0.00 | 0.00 | 0.00 | 0.00 | 0.01 |
| RAC | 0.51 | 0.28 | 0.42 | 0.11 | 0.28 |
| LLM1 | -0.28 | -0.10 | -0.22 | -0.24 | 0.06 |
| LLM2 | <u>0.48</u> | 0.31 | <u>0.41</u> | <u>0.10</u> | <u>0.27</u> |
| FC (ours) | 0.40 | <u>0.29</u> | 0.36 | <u>0.10</u> | 0.24 |

Table 4: Mean relative gains for the factuality metrics obtained on the BIO dataset.

| corrector | Pr \uparrow | R@K \uparrow | F1@K \uparrow | V \uparrow | C \uparrow |
|------------------------|---------------|----------------|-----------------|--------------|--------------|
| mixtral-8x22b-instruct | | | | | |
| CRITIC | -0.12 | -0.04 | -0.09 | -0.06 | -0.01 |
| RAC | <u>-0.01</u> | -0.01 | -0.01 | -0.01 | 0.00 |
| LLM1 | -0.02 | 0.00 | -0.01 | 0.00 | 0.00 |
| LLM2 | -0.03 | 0.00 | -0.02 | 0.00 | 0.00 |
| FC (ours) | 0.02 | 0.00 | 0.00 | 0.02 | 0.00 |
| llama-3.3-70b-instruct | | | | | |
| CRITIC | -0.03 | 0.00 | -0.01 | -0.01 | -0.01 |
| RAC | <u>0.03</u> | 0.01 | 0.02 | 0.02 | 0.00 |
| LLM1 | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| LLM2 | 0.01 | 0.00 | 0.01 | 0.00 | 0.00 |
| FC (ours) | 0.04 | 0.01 | 0.02 | 0.01 | 0.01 |
| granite-4.0-h-small | | | | | |
| CRITIC | -0.13 | -0.09 | -0.11 | -0.05 | -0.04 |
| RAC | <u>-0.06</u> | -0.01 | -0.04 | -0.01 | -0.02 |
| LLM1 | -0.08 | -0.03 | -0.05 | -0.03 | -0.02 |
| LLM2 | -0.07 | -0.01 | -0.04 | -0.03 | -0.01 |
| FC (ours) | 0.00 | 0.00 | 0.00 | 0.01 | 0.00 |
| gpt-oss-120b | | | | | |
| CRITIC | -0.03 | -0.01 | -0.02 | 0.01 | -0.01 |
| RAC | 0.12 | 0.03 | 0.08 | 0.06 | 0.02 |
| LLM1 | -0.49 | -0.26 | -0.41 | -0.19 | -0.08 |
| LLM2 | -0.03 | -0.01 | -0.03 | 0.03 | -0.04 |
| FC (ours) | <u>0.00</u> | 0.00 | 0.00 | 0.01 | 0.00 |

Table 5: Mean relative gains for the factuality metrics obtained on the ASKHIST dataset.

erating corrections³. Table 6 reports the mean relative gains on the VELI5 dataset when using

³LoRA adapter available at <https://ibm.biz/granite-guardian-3-2-5b-lora-factuality-correction>.

| corrector | Pr \uparrow | R@K \uparrow | F1@K \uparrow | V \uparrow | C \uparrow |
|------------------------|---------------|----------------|-----------------|--------------|--------------|
| mixtral-8x22b-instruct | | | | | |
| FC (ours) | 0.34 | 0.09 | 0.24 | 0.14 | 0.07 |
| SFT (ours) | 0.19 | 0.10 | 0.16 | 0.05 | 0.08 |
| llama-3.3-70b-instruct | | | | | |
| FC (ours) | 0.27 | 0.07 | 0.18 | 0.05 | 0.08 |
| SFT (ours) | 0.25 | 0.11 | 0.19 | 0.10 | 0.08 |
| granite-4.0-h-small | | | | | |
| FC (ours) | 0.29 | 0.10 | 0.21 | 0.11 | 0.07 |
| SFT (ours) | 0.22 | 0.11 | 0.18 | 0.09 | 0.07 |
| gpt-oss-120b | | | | | |
| FC (ours) | 0.36 | 0.19 | 0.30 | 0.07 | 0.14 |
| SFT (ours) | 0.08 | 0.08 | 0.08 | 0.07 | 0.02 |

Table 6: Mean relative gains for the factuality metrics on the VELI5 dataset using the SFT corrector.

the FACTCORRECTOR pipeline versus the LoRA-based correction approach (denoted as SFT). Overall, the SFT method produces substantially better corrections than the original responses, which is expected given that the evaluation instances originate from the same distribution as the training data. In contrast, on the out-of-distribution BIO and ASKHIST datasets the gains achieved by SFT are smaller than before; however, most gains remain non-negative, indicating that the SFT approach generalizes reasonably well (see Appendix C).

6.6 Human Evaluation of VELI5

We conducted a user study to evaluate the quality of corrections in the VELI5 dataset (see Appendix C.6). Table 7 summarizes results for 30 filtered tasks, each averaging 2.30 incorrect and 4.40 correct atoms. At the task level, original responses were highly relevant to the user question (90% fully, 10% partially), while corrections achieved even greater relevance (96.7% fully). For incorrect atoms, 67.4% were successfully corrected by FACTCORRECTOR, with 26.4% undercorrected and only minor proportions overcorrected (3.3%) or invalid (2.9%). Correct atoms were preserved in 60.9% of cases, with 27.7% not preserved and small fractions partially preserved or invalid. Relevance analysis showed both incorrect and correct atoms were predominantly relevant (75.0% and 72.4%), underscoring the effectiveness of the correction process.

7 Conclusion

The contributions of this paper are threefold. (1) We introduce FACTCORRECTOR, a novel post-hoc method for improving long-form factuality by leveraging feedback on incorrect response segments and retrieved contradicting evidence to generate refined

answers. (2) We develop VELI5, a benchmark dataset of curated question–answer pairs explicitly verified and corrected using FACTCORRECTOR, enabling rigorous evaluation. (3) We conduct extensive empirical studies across a wide range of open-source LLMs and datasets, demonstrating that FACTCORRECTOR consistently achieves the most reliable factuality improvements compared with strong state-of-the-art baselines.

Limitations

We acknowledge further limitations of the proposed FACTCORRECTOR framework.

The Atomizer and Reviser components are sensitive to the quality of the prompt and few shot examples used as well as the LLM employed to perform the atomic unit decomposition and decontextualization of the response (and correction). In our work we only considered open-source models from the Mixtral, Llama, Granite and GPT-OSS families. Furthermore, decomposing a given response can be done at different granularities such as sentence level, paragraph level or the entire response level. Our implementation is limited to decomposing the response in one shot.

The quality of the contexts retrieved by the Retriever component for each atomic unit depends on both the implementation details of the component and the formulation of the query string it receives. In our approach, we leverage an LLM to generate a well-structured Google search query for each atomic unit, which is then used to retrieve the top k search results via the Serper API (serper.dev). Each search result typically includes a short snippet and a hyperlink; to enrich the retrieved context, we additionally extract up to 4,000 characters of text from the linked page.

The Evaluator component relies on the logical relationships between atoms and contexts. It also depends on the quality of the prompt and the underlying LLM. As before, we only used open-source models such as mixtral-8x22b-instruct, llama-3.3-70b-instruct, granite-4.0-h-small, and gpt-oss-120b with a fairly straightforward prompt. It is possible to craft better prompts that could lead to a better extraction of the relationships. Fine-tuning is another option to obtain a stronger relation model.

The Corrector component (i.e., refinement model) is also sensitive to the quality of the prompt used.

The version presented in Table 14 in Appendix D is the current version that works fairly well across different models. Clearly, improved prompts could lead to much better corrections and this is currently one direction of future research.

Finally, since the FACTCORRECTOR pipeline is tightly integrated with the FACTREASONER factuality assessor (Marinescu et al., 2025), it inherits the computational overhead of the latter. Specifically, the Evaluator component requires $O(n \cdot m)$ LLM calls to extract the logical relationships between atoms and contexts, where n is the number of atomic units in the response (or correction), m is the total number of non-duplicated contexts retrieved for the atoms. In contrast, the SFT approach for correction requires only a single LLM call per correction.

Ethical Statement

We recognize the positive and negative societal impacts of LLMs in general, including potential misuse of our work around factuality correction of LLM generated output. We note that the datasets considered are public and peer reviewed, there are no human subjects involved, and as far as we know, there are no obvious harmful consequences from our work. All creators and original owners of assets have been properly credited and licenses and terms of use have been respected. We have not conducted crowd-sourcing experiments or research with human subjects.

References

- Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamera Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. 2022. *Constitutional ai: Harmlessness from ai feedback*. Preprint, arXiv:2212.08073.
- Farima Fatahi Bayat, Lechen Zhang, Sheza Munir, and Lu Wang. 2025. *Factbench: A dynamic benchmark*

for in-the-wild language model factuality evaluation. *Preprint*, arXiv:2410.22257.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris 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 *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Angelica Chen, J  r  my Scheurer, Tomasz Korbak, Jon Ander Campos, Jun Shern Chan, Samuel R. Bowman, Kyunghyun Cho, and Ethan Perez. 2024. [Improving code generation by training with natural language feedback](#). *Preprint*, arXiv:2303.16749.

Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. 2022. [Codet: Code generation with generated tests](#). *Preprint*, arXiv:2207.10397.

I-Chun Chern, Steffi Chern, Shiqi Chen, Weizhe Yuan, Kehua Feng, Chunting Zhou, Junxian He, Graham Neubig, and Pengfei Liu. 2023. [Factool: Factuality detection in generative ai – a tool augmented framework for multi-task and multi-domain scenarios](#). *Preprint*, arXiv:2307.13528.

Aakanksha Chowdhery, Sharan Narang, and Jacob Devlin. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 4(1):1–113.

Roi Cohen, May Hamri, Mor Geva, and Amir Globerson. 2023. Lm vs lm: Detecting factual errors via cross examination. *arXiv preprint arXiv:2305.13281*.

Rina Dechter. 2003. *Constraint Processing*. Morgan Kaufmann Publishers.

Adam Dejl, James Barry, Alessandra Pascale, and Javier Carnerero Cano. 2025. [Comprehensiveness metrics for automatic evaluation of factual recall in text generation](#). *Preprint*, arXiv:2510.07926.

Yann Dubois, Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. 2024. AlpacaFarm: A simulation framework for methods that learn from human feedback. In *NeurIPS*.

Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Auli. 2019. [ELI5: Long form question answering](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3558–3567, Florence, Italy. Association for Computational Linguistics.

Deep Ganguli, Amanda Askell, Nicholas Schiefer, Thomas I. Liao, Kamil   Luko  i  t  , Anna Chen, Anna Goldie, Azalia Mirhoseini, Catherine Olsson, Danny

Hernandez, Dawn Drain, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jackson Kernion, Jamie Kerr, Jared Mueller, Joshua Landau, Kamal Ndousse, Karina Nguyen, Liane Lovitt, Michael Sellitto, Nelson Elhage, Noemi Mercado, Nova DasSarma, Oliver Rausch, Robert Lasenby, Robin Larson, Sam Ringer, Sandipan Kundu, Saurav Kadavath, Scott Johnston, Shauna Kravec, Sheer El Showk, Tamera Lanham, Timothy Telleen-Lawton, Tom Henighan, Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, Christopher Olah, Jack Clark, Samuel R. Bowman, and Jared Kaplan. 2023. [The capacity for moral self-correction in large language models](#). *Preprint*, arXiv:2302.07459.

Ge Gao, Hung-Ting Chen, Yoav Artzi, and Eunsol Choi. 2023. Continually improving extractive qa via human feedback. In *EMNLP*, pages 406–417.

Zelalem Gero, Chandan Singh, Hao Cheng, Tristan Naumann, Michel Galley, Jianfeng Gao, and Hoifung Poon. 2023. [Self-verification improves few-shot clinical information extraction](#). *Preprint*, arXiv:2306.00024.

Amelia Glaese, Nat McAleese, Maja Tr  bacz, John Aslanides, Vlad Firoiu, Timo Ewalds, Maribeth Rauh, Laura Weidinger, Martin Chadwick, Phoebe Thacker, Lucy Campbell-Gillingham, Jonathan Uesato, Po-Sen Huang, Ramona Comanescu, Fan Yang, Abigail See, Sumanth Dathathri, Rory Greig, Charlie Chen, Doug Fritz, Jaume Sanchez Elias, Richard Green, So  a Mok  r  , Nicholas Fernando, Boxi Wu, Rachel Foley, Susannah Young, Iason Gabriel, William Isaac, John Mellor, Demis Hassabis, Koray Kavukcuoglu, Lisa Anne Hendricks, and Geoffrey Irving. 2022. [Improving alignment of dialogue agents via targeted human judgements](#). *Preprint*, arXiv:2209.14375.

Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong Shen, Yujia Yang, Nan Duan, and Weizhu Chen. 2024. Critic: Large language models can self-correct with tool-interactive critiquing. In *ICLR*.

Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. 2023. Reasoning with language model is planning with world model. In *EMNLP*.

Hangfeng He, Hongming Zhang, and Dan Roth. 2022. [Rethinking with retrieval: Faithful large language model inference](#). *Preprint*, arXiv:2301.00303.

Yufang Hou, Alessandra Pascale, Javier Carnerero Cano, Tigran Tchrakian, Radu Marinescu, Elizabeth Daly, Inkit Padhi, and Prasanna Sattigeri. 2024. [Wiki-contradict: A benchmark for evaluating llms on real-world knowledge conflicts from wikipedia](#). In *Advances in Neural Information Processing Systems*, volume 37, pages 109701–109747. Curran Associates, Inc.

Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. 2025. [A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions](#). *ACM Trans. Inf. Syst.*, 43(2).

- Albert Q. Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, Gianna Lengyel, Guillaume Bour, Guillaume Lample, L  lio Renard Lavaud, Lucile Saulnier, Marie-Anne Lachaux, Pierre Stock, Sandeep Subramanian, Sophia Yang, Szymon Antoniak, Teven Le Scao, Th  ophile Gervet, Thibaut Lavril, Thomas Wang, Timoth  e Lacroix, and William El Sayed. 2024. [Mixtral of experts](#). *Preprint*, arXiv:2401.04088.
- Muhammad Khalifa, Lajanugen Logeswaran, Moon-tae Lee, Honglak Lee, and Lu Wang. 2023. Grace: Discriminator-guided chain-of-thought reasoning. In *EMNLP*.
- Daphne Koller and Nir Friedman. 2009. *Probabilistic Graphical Models: Principles and Techniques*. MIT Press.
- Changmao Li and Jeffrey Flanigan. 2025. [Rac: Efficient llm factuality correction with retrieval augmentation](#). *Preprint*, arXiv:2410.15667.
- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Hao Liu, Carmelo Sferrazza, and Pieter Abbeel. 2023. [Chain of hindsight aligns language models with feedback](#). *Preprint*, arXiv:2302.02676.
- Qiang Liu and Alexander Ihler. 2011. Bounding the partition function using Holder’s inequality. In *International Conference on Machine Learning (ICML)*, pages 849–856.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. 2023. [Self-refine: Iterative refinement with self-feedback](#). *Preprint*, arXiv:2303.17651.
- Radu Marinescu, Debarun Bhattacharjya, Junkyu Lee, Tigran Tchrakian, Javier Carnerero Cano, Yufang Hou, Elizabeth Daly, and Alessandra Pascale. 2025. Factreasoner: A probabilistic approach to long-form factuality assessment for large language models. In *EMNLP*.
- 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 *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12076–12100.
- Ansong Ni, Srini Iyer, Dragomir Radev, Ves Stoyanov, Wen tau Yih, Sida I. Wang, and Xi Victoria Lin. 2023. Lever: Learning to verify language-to-code generation with execution. In *ICML*.
- OpenAI, :, Sandhini Agarwal, Lama Ahmad, Jason Ai, Sam Altman, Andy Applebaum, Edwin Arbus, Rahul K. Arora, Yu Bai, Bowen Baker, Haiming Bao, Boaz Barak, Ally Bennett, Tyler Bertao, Nivedita Brett, Eugene Brevdo, Greg Brockman, Sebastien Bubeck, Che Chang, Kai Chen, Mark Chen, Enoch Cheung, Aidan Clark, Dan Cook, Marat Dukhan, Casey Dvorak, Kevin Fives, Vlad Fomenko, Timur Garipov, Kristian Georgiev, Mia Glaese, Tarun Gogineni, Adam Goucher, Lukas Gross, Katia Gil Guzman, John Hallman, Jackie Hehir, Johannes Heidecke, Alec Helyar, Haitang Hu, Romain Huet, Jacob Huh, Saachi Jain, Zach Johnson, Chris Koch, Irina Kofman, Dominik Kundel, Jason Kwon, Volodymyr Kyrlyov, Elaine Ya Le, Guillaume Leclerc, James Park Lennon, Scott Lessans, Mario Lezcano-Casado, Yuanzhi Li, Zhuohan Li, Ji Lin, Jordan Liss, Lily, Liu, Jiancheng Liu, Kevin Lu, Chris Lu, Zoran Martinovic, Lindsay McCallum, Josh McGrath, Scott McKinney, Aidan McLaughlin, Song Mei, Steve Mostovoy, Tong Mu, Gideon Myles, Alexander Neitz, Alex Nichol, Jakub Pachocki, Alex Paino, Dana Palmie, Ashley Pantuliano, Giambattista Parascandolo, Jongsoo Park, Leher Pathak, Carolina Paz, Ludovic Peran, Dmitry Pimenov, Michelle Pokrass, Elizabeth Proehl, Huida Qiu, Gaby Raila, Filippo Raso, Hongyu Ren, Kimmy Richardson, David Robinson, Bob Rotsted, Hadi Salman, Suvansh Sanjeev, Max Schwarzer, D. Sculley, Harshit Sikchi, Kendal Simon, Karan Singhal, Yang Song, Dane Stuckey, Zhiqing Sun, Philippe Tillet, Sam Toizer, Foivos Tsimpouras, Nikhil Vyas, Eric Wallace, Xin Wang, Miles Wang, Olivia Watkins, Kevin Weil, Amy Wendling, Kevin Whinnery, Cedric Whitney, Hannah Wong, Lin Yang, Yu Yang, Michihiro Yasunaga, Kristen Ying, Wojciech Zaremba, Wenting Zhan, Cyril Zhang, Brian Zhang, Eddie Zhang, and Shengjia Zhao. 2025. [gpt-oss-120b & gpt-oss-20b model card](#). *Preprint*, arXiv:2508.10925.
- Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback. In *NeurIPS*.
- Liangming Pan, Michael Saxon, Wenda Xu, Deepak Nathani, Xinyi Wang, and William Yang Wang. 2024. [Automatically correcting large language models: Surveying the landscape of diverse automated correction strategies](#). *Transactions of the Association for Computational Linguistics*, 12:484–506.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: a method for automatic evaluation of machine translation](#). In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Debjit Paul, Mete Ismayilzada, Maxime Peyrard, Beatriz Borges, Antoine Bosselut, Robert West, and Boi Faltings. 2024. Refiner: Reasoning feedback on intermediate representations. In *EACL*.
- Judea Pearl. 1988. *Probabilistic Reasoning in Intelli-*

gent Systems. Morgan Kaufmann.

IBM Research. 2024. [Granite foundation models](#).

Pranab Sahoo, Prabhash Meharia, Akash Ghosh, Sriparna Saha, Vinija Jain, and Aman Chadha. 2024. [A comprehensive survey of hallucination in large language, image, video and audio foundation models](#). *Preprint*, arXiv:2405.09589.

Jérémy Scheurer, Jon Ander Campos, Tomasz Korbak, Jun Shern Chan, Angelica Chen, Kyunghyun Cho, and Ethan Perez. 2024. [Training language models with language feedback at scale](#). *Preprint*, arXiv:2303.16755.

Thibault Sellam, Dipanjan Das, and Ankur P. Parikh. 2020. [BLEURT: learning robust metrics for text generation](#). *CoRR*, abs/2004.04696.

Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. 2023. [Reflexion: Language agents with verbal reinforcement learning](#). *Preprint*, arXiv:2303.11366.

Yixiao Song, Yekyung Kim, and Mohit Iyyer. 2024. [VeriScore: Evaluating the factuality of verifiable claims in long-form text generation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 9447–9474, Miami, Florida, USA. Association for Computational Linguistics.

Zhaochen Su, Jun Zhang, Xiaoye Qu, Tong Zhu, Yanshu Li, Jiashuo Sun, Juntao Li, Min Zhang, and Yu Cheng. 2024. [Conflictbank: A benchmark for evaluating the influence of knowledge conflicts in llm](#). *arXiv preprint arXiv:2408.12076*.

S. M Towhidul Islam Tonmoy, S M Mehedi Zaman, Vinija Jain, Anku Rani, Vipula Rawte, Aman Chadha, and Amitava Das. 2024. [A comprehensive survey of hallucination mitigation techniques in large language models](#). *Preprint*, arXiv:2401.01313.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. [Llama: Open and efficient foundation language models](#). *Preprint*, arXiv:2302.13971.

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

Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Shengping Liu, Bin Sun, Kang Liu, and Jun Zhao. 2023. Large language models are better reasoners with self-verification. In *EMNLP*, pages 2550–2575.

Fangyuan Xu, Yixiao Song, Mohit Iyyer, , and Eunsol Choi. 2023. A critical evaluation of evaluations for long-form question answering. In *In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3225–3245.

Rongwu Xu, Zehan Qi, Zhijiang Guo, Cunxiang Wang, Hongru Wang, Yue Zhang, and Wei Xu. 2024. Knowledge conflicts for llms: A survey. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 8541–8565.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. In *NeurIPS*.

Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D. Goodman. 2022. Star: Bootstrapping reasoning with reasoning. In *NeurIPS*.

Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. 2023. [Siren’s song in the ai ocean: A survey on hallucination in large language models](#). *Preprint*, arXiv:2309.01219.

Jialun Zhong, Wei Shen, Yanzeng Li, Songyang Gao, Hua Lu, Yicheng Chen, Yang Zhang, Wei Zhou, Jinjie Gu, and Lei Zou. 2025. [A comprehensive survey of reward models: Taxonomy, applications, challenges, and future](#). *Preprint*, arXiv:2504.12328.

A Background on Graphical Models

Graphical models such as Bayesian or Markov networks provide a powerful framework for reasoning about conditional dependency structures over many variables (Pearl, 1988; Koller and Friedman, 2009).

A *graphical model* is a tuple $\mathcal{M} = \langle \mathbf{X}, \mathbf{D}, \mathbf{F} \rangle$, where $\mathbf{X} = \{X_1, \dots, X_n\}$ is a set of variables, $\mathbf{D} = \{D_1, \dots, D_n\}$ is the set of their finite domains of values and $\mathbf{F} = \{f_1, \dots, f_m\}$ is a set of discrete positive real-valued functions. Each function f_i (also called *factor*) is defined on a subset of variables $\mathbf{S}_i \subseteq \mathbf{X}$ called its *scope* and denoted by $\text{vars}(f_i)$. The model \mathcal{M} defines a factorized probability distribution on \mathbf{X} :

$$P(\mathbf{x}) = \frac{1}{Z} \prod_{j=1}^m f_j(\mathbf{x}) \text{ s.t. } Z = \sum_{\mathbf{x} \in \Omega(\mathbf{X})} \prod_{j=1}^m f_j(\mathbf{x}) \quad (1)$$

where the normalization constant Z is known as the *partition function* and $\Omega(\mathbf{X})$ denotes the Cartesian product of the variables domains.

The function scopes of a model \mathcal{M} define a *primal graph* whose vertices are the variables and its edges connect any two variables that appear in the scope of the same function.

A common inference task over graphical models is to compute the posterior marginal distributions over all variables. Namely, for each variable $X_i \in \mathbf{X}$ and domain value $x_i \in D_i$, compute:

$$P(x_i) = \sum_{\mathbf{x} \in \Omega(\mathbf{X})} \delta_{x_i}(\mathbf{x}) \cdot P(\mathbf{x}) \quad (2)$$

where $\delta_{x_i}(\mathbf{x})$ is 1 if X_i is assigned x_i in \mathbf{x} and 0 otherwise (Koller and Friedman, 2009).

Example 3. Figure 4 shows a graphical model with 3 bi-valued variables X_1 , X_2 and X_3 and 3 binary functions $f_1(X_1, X_2)$, $f_2(X_1, X_3)$ and $f_3(X_2, X_3)$. The joint probability distribution is given by $P(X_1, X_2, X_3) = \frac{1}{Z} \cdot f_1(X_1, X_2) \cdot f_2(X_1, X_3) \cdot f_3(X_2, X_3)$. In this case, the posterior marginal distribution of X_1 is: $P(X_1 = 0) = 0.46$ and $P(X_1 = 1) = 0.54$, respectively.

Equation 2 can be solved using any probabilistic inference algorithm for graphical models, such as

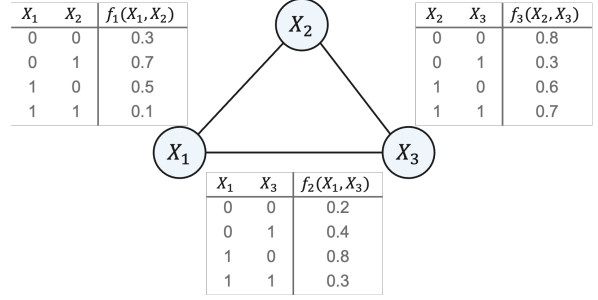


Figure 4: A graphical model with three bi-valued variables X_1 , X_2 and X_3 , and three binary functions.

variable elimination (Dechter, 2003), belief propagation (Pearl, 1988), or variational inference (Liu and Ihler, 2011). In our implementation, we employed the Weighted Mini-Buckets (WMB) algorithm (Liu and Ihler, 2011). WMB is parameterized by an i-bound, which controls the trade-off between computational complexity and inference accuracy. For our experiments, we selected an i-bound of 6, which enabled us to solve all inference problems efficiently. Notably, WMB proved highly effective in practice, solving each inference instance in under 0.05 seconds across our benchmark datasets.

B Background on Factuality Assessment and Correction with Feedback

B.1 Long-Form Factuality Assessment

Let y be the long-form response generated by an LLM to a query x . We assume that y can be decomposed into a set of n *atomic units* (or *atoms*) that can be either true or false, denoted by $\mathcal{A}_y = \{a_1, a_2, \dots, a_n\}$ (Min et al., 2023; Song et al., 2024; Wei et al., 2024). An atomic unit $a_i \in \mathcal{A}_y$ is defined as a short sentence conveying one piece of information (e.g., a claim or a fact). Furthermore, given an external knowledge source \mathcal{K} ⁴, we say that an atomic unit $a_i \in \mathcal{A}_y$ is *supported* by \mathcal{K} if there exists at least one piece of information in \mathcal{K} (e.g., a passage) called a *context* that undebatably supports a_i . Otherwise, the atomic unit is *not supported*.

The *factual precision* denoted by $Pr(y)$ of the response y with respect to a knowledge source \mathcal{C} is defined as: $Pr(y) = \frac{S(y)}{|\mathcal{A}_y|}$, where $S(y) = \sum_{i=1}^n \mathbb{I}[a_i \text{ is supported by } \mathcal{K}]$ is the number of supported atomic units. The *factual recall* $R_K(y)$ up to the K -th supported atomic unit is given by: $R_K(y) = \min(\frac{S(y)}{K}, 1)$. Finally, an F_1 measure for long-form factuality denoted by $F1@K$ can be de-

⁴For example, \mathcal{K} could be Wikipedia, Google Search, or a collection of documents embedded into a vector database.

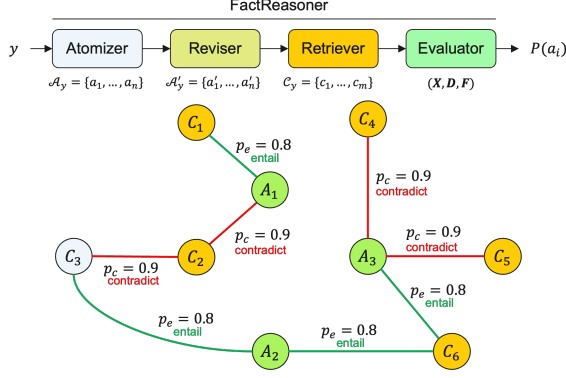


Figure 5: FACTREASONER pipeline and an example graphical model with 3 atom and 6 context variables.

defined as: $F_1@K(y) = \frac{2 \cdot Pr(y) \cdot R_K(y)}{Pr(y) + R_K(y)}$ if $S(y) > 0$, and 0 otherwise (Wei et al., 2024).

B.2 The FACTREASONER Factuality Assessor

Existing long-form factuality assessors such as FactScore (Min et al., 2023), VeriScore (Song et al., 2024) and others (Wei et al., 2024; Bayat et al., 2025) are prompt-based approaches that essentially prompt an LLM to determine if each atomic unit of the response is supported (factually correct) or not by the retrieved evidence. FACTREASONER (Marinescu et al., 2025) is a recent factuality assessor that, unlike the prompt-based approaches, leverages probabilistic reasoning to assess the factuality of the generated response with respect to an external knowledge source \mathcal{K} .

Specifically, the FACTREASONER pipeline for factuality assessment shown in Figure 5 (top) consists of four main stages called *Atomizer*, *Reviser*, *Retriever* and *Evaluator*, respectively. The *Atomizer* prompts an LLM to decompose the response y into a set of n atomic units \mathcal{A}_y by applying any of the decomposition strategies proposed recently (Min et al., 2023; Bayat et al., 2025). Subsequently, the *Reviser* also uses an LLM to revise the atoms such that the pronouns, unknown entities, or incomplete names are replaced with their corresponding named entities in the response (Wei et al., 2024). Next, the *Retriever* is responsible for querying an external knowledge source \mathcal{K} to retrieve the contexts relevant to the response’s atoms (Song et al., 2024). Finally, the *Evaluator* constructs a graphical model representing a joint probability distribution over the atomic units in the response and their corresponding contexts in \mathcal{C} . For each atom a_i , it then computes the posterior marginal probability $P(a_i)$ which quantifies the likelihood the a_i is true (or supported) given the information available in \mathcal{K} .

Finally, the factuality measures, such as $Pr(y)$ or $F_1@K(y)$, can be readily calculated using the supported atoms (Marinescu et al., 2025).

Graphical Model. FACTREASONER’s graphical model is defined by a tuple $\langle \mathbf{X}, \mathbf{D}, \mathbf{F} \rangle$ where \mathbf{X} is a set of Boolean variables A_i or C_j associated with each atom $a_i \in \mathcal{A}_y$ or context $c_j \in \mathcal{C}_y$, respectively. The domains \mathbf{D} are $\{\text{true}, \text{false}\}$ for all variables (e.g., we denote a_i and $\neg a_i$ for $A_i = \text{true}$ and $A_i = \text{false}$, respectively). Each variable has a unary factor $f(\cdot)$ in \mathbf{F} encoding prior belief. Atoms have uniform priors: $f(a_i) = f(\neg a_i) = 0.5$. Contexts from reliable sources have high prior (e.g., $f(c_j) = 0.99$), while less reliable sources use smaller values. Binary factors $f(A_i, C_j)$ and $f(C_j, C_k)$ capture probabilistic logical relations between utterances. A relation model $p_\theta(\cdot|t, t')$ (either a specialized BERT model or an LLM) predicts the most likely relation from $\{\text{neutral}, \text{entail}, \text{contradict}\}$ together with its probability. Therefore, each factor $f(A_i, C_j)$ or $f(C_j, C_k)$ represents a probabilistic encoding of an entailment relationship (i.e., context C_j entails or supports atom A_i) or a contradiction relationship (i.e., context C_j contradicts atom A_i), respectively. Note that the neutral relationships are ignored.

Example 4. Figure 5 (bottom) shows an example graphical model with 3 atoms (A_1, A_2, A_3), 6 contexts (C_1, \dots, C_6) and 8 binary relations between atoms and contexts. Note that each edge in the graph is labeled by the probability of the corresponding entailment or contradiction relationship.

B.3 Post-hoc Correction with Feedback

In general, a *post-hoc correction strategy* leverages a Refine Model \mathcal{R} to correct the response \hat{y} generated by a Language Model \mathcal{M} to a user query x using the feedback provided by a Critic Model \mathcal{C} regarding the quality of the response.

Formally, let $\mathcal{M} : \mathcal{X} \rightarrow \mathcal{Y}$ be a language model that performs a specific task by mapping the user input $x \in \mathcal{X}$ to an output text $\hat{y} \in \mathcal{Y}$. A wide range of NLP tasks adhere to this formulation, such as QA, summarization, translation, and more. The initial generation \hat{y} may have problems such as hallucinations, factual errors, or incorrect reasoning. The critic model $\mathcal{C} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{F}$ learns to generate feedback $x, \hat{y} \rightarrow c$ where $\hat{y} \sim \mathcal{M}(x)$ is the output of the language model, and c is the feedback of some format such as a scalar value or natural language. A simple example is binary feed-

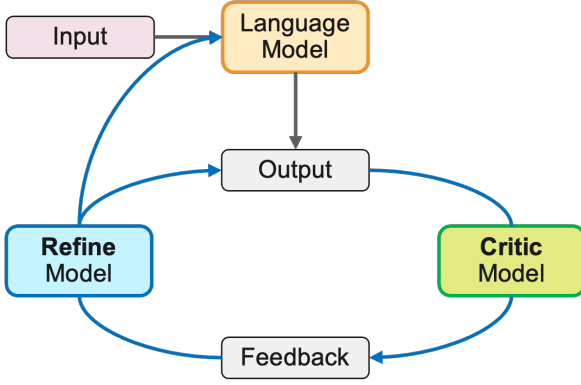


Figure 6: Post-hoc correction with feedback.

back of whether the output is good or bad given the input $\mathcal{C} : \mathcal{X} \times \mathcal{Y} \rightarrow \{0, 1\}$. The refine model $\mathcal{R} : \mathcal{X} \times \mathcal{Y} \times \mathcal{F} \rightarrow \mathcal{Y}$ learns to repair/correct an output $x, \hat{y}, c \rightarrow y_{new}$ based on the feedback c , where y_{new} is the revised output (Pan et al., 2024).

Figure 6 shows how the components of a post-hoc correction pipeline interact to repair the output of a language model.

C Additional Details on Experiments

In this section, we report additional details on our experiments.

C.1 Baseline Correctors

For our purpose, we consider two recent state-of-the-art post-hoc correction with feedback approaches, namely CRITIC (Gou et al., 2024) and RAC (Li and Flanigan, 2025). Specifically, CRITIC is a prompt-based iterative correction method that uses LLMs with retrieved knowledge. At each iteration, CRITIC prompts an LLM to generate a query to retrieve relevant information from a knowledge base and subsequently revise the output based on the history of revising and decide the most possible answer. RAC is a post-hoc factuality correction method designed to improve the accuracy of LLMs responses without additional fine-tuning. Its key steps are: the model’s generated text is broken down into individual factual statements, relevant information is fetched from trusted sources (e.g., Wikipedia or document collections), and each atomic fact is checked against retrieved evidence using a verification model. Incorrect statements are then corrected based on reliable content.

In addition, we consider two prompting strategies, denoted LLM1 and LLM2. The first, LLM1, performs response correction using only the model’s internal knowledge, without incorporating any re-

trieved context. In contrast, LLM2 uses the contexts retrieved for the question only and instructs the model to ignore its internal knowledge to refine the response (see Appendix D for the actual prompts).

In our experiments, we instantiated the competing correctors with open-source LLMs belonging to the IBM Granite⁵, Meta LLaMA⁶, MistralAI Mixtral⁷ and OpenAI⁸ families, namely: granite-4.0-h-small, llama-3.3-70b-instruct, mixtral-8x22b-instruct, and gpt-oss-120b, respectively. All our LLMs are hosted remotely on compute nodes with A100 80GB GPUs and accessed via `litellm` APIs capable of serving 1500 prompts per second.

External Knowledge Source. We consider Google Search results as our external knowledge source \mathcal{K} . For a given atom, the top k results are retrieved as contexts from google.com using the Serper API⁹. In this case, a context is a tuple (t, l, s, d) , where t is the title of the web-page, l is the link, s is a short text snippet or summary and d is the content retrieved from l (but capped at max 4000 characters). We used $k = 3$ for the Google Search results (Min et al., 2023; Wei et al., 2024).

C.2 Datasets

We experimented with the following datasets: a reduced version of the VELI5 dataset presented in the previous section and consisting of 200 instances sampled uniformly at random across all categories, Biographies (BIO) (Min et al., 2023), and AskHistorians (ASKHIST) (Xu et al., 2023). These datasets have been widely adopted in prior work and are considered representative benchmarks for long-form factuality assessment, as they encompass a diverse range of topics and tasks, including creative writing, history, astronomy, chemistry, and more.

The BIO dataset contains 183 biographic passages spanning up to two paragraphs that were generated by ChatGPT for various person entities that have a Wikipedia page (Min et al., 2023). The ASKHIST datasets contains 200 questions sourced from the r/AskHistorians Reddit forum together with a long-form response of up to two paragraphs in length

⁵<https://huggingface.co/ibm-granite>

⁶<https://huggingface.co/meta-llama>

⁷<https://huggingface.co/mistralai>

⁸<https://openai.com>

⁹<https://serper.dev>

generated using the llama-3.3-70b-instruct model (Marinescu et al., 2025).

Additionally, we include results on the CONFLICTS dataset introduced in (Marinescu et al., 2025). This dataset comprises 100 claims (atomic units) randomly drawn from the ConflictBank benchmark (Su et al., 2024). Each claim, originally sourced from Wikidata, is assumed to be true (i.e., supported). For our purpose, we treat each claim as a *response* and provide both a supporting context (as in ConflictBank) and a conflicting context (representing misinformation). Importantly, since all responses are true, they require no corrections, offering a controlled setting for our evaluation.

C.3 Measures of Performance

For each dataset \mathcal{D} and each competing corrector, we compute three factuality metrics: precision (Pr), recall at K ($R@K$), and $F_1@K$, averaged over all prompts in \mathcal{D} , where K is set to the median number of atoms. These metrics are assessed using the FACTREASONER framework (Marinescu et al., 2025). In addition, we report two complementary measures: verifiability (V) and comprehensiveness (C). The verifiability score V is defined as the number of atoms in the response (or correction) that can be verified—i.e., atoms connected to either supporting or contradicting evidence in FR’s graph. Comprehensiveness C quantifies coverage and is given by $C = \frac{|\mathcal{A}_{in}|}{|\mathcal{A}_{in}| + |\mathcal{A}_{out}|}$, where \mathcal{A}_{in} denotes the atoms covered by the response (or correction), and \mathcal{A}_{out} denotes atoms that are *uncovered* or missing (Dejl et al., 2025).

For each score S (e.g., precision), we also report its *relative gain*, denoted as $G(S)$ and defined by:

$$G(S) = \frac{2 \cdot (S_c - S_r)}{S_c + S_r}$$

where S_r and S_c are the scores of the original response and the correction, respectively. A positive $G(S)$ indicates that the correction outperforms the response, while a negative value means that the correction performs worse. By construction, $G(S)$ ranges from -2 to 2 and remains well defined even when either S_r or S_c equals zero.

C.4 Detailed Results

Number of Atoms. Table 15 presents the mean number of atoms and their standard deviations for both initial responses and subsequent corrections generated by the different correctors across datasets

and base models. Notably, in some cases, the corrections contain substantially more atoms than the original responses. This increase may contribute to the observed decline in relative gains for factuality metrics, including the negative gains reported in Tables 9, 11, and 13.

CONFLICTS. Table 8 summarizes the results obtained on our controlled experiment with the CONFLICTS dataset. In this setting, since the original is known to be true and doesn’t require correction, we assess the similarity between correction and response using the ROUGE, BLEU, and BLEURT metrics (Lin, 2004; Papineni et al., 2002; Sellam et al., 2020). Additionally, we report JUDGE, which counts the number of instances where an LLM-as-a-Judge model (in our case DeepSeek-3.2), prompted appropriately, infers that the correction is equivalent to the original response. The table shows the mean and standard deviation for all metrics. The results indicate that FACTCORRECTOR achieves superior performance with the Mixtral and Granite models, while ranking second when using LLaMA and GPT-OSS models, where RAC attains the best scores. Therefore, both FACTCORRECTOR and RAC correctly identify that no correction is needed. In contrast, the LLM1 and LLM2 baselines perform rather poorly almost always producing corrections significantly different from the response, while CRITIC often struggles to follow its instructions also leading to poor results.

VELI5. In Table 9 we report the mean relative gains and standard deviations obtained for the factuality metrics on the VELI5 dataset using the mixtral-8x22b-instruct, llama-3.3-70b-instruct, granite-4.0-h-small and gpt-oss-120b models. The results demonstrate that FACTCORRECTOR consistently delivers strong and reliable improvements in factuality across all evaluated metrics and models. For example, when considering relative gains on precision and $F1@K$, FC outperforms all baselines, achieving macro-averaged gains of 0.315 for precision, and 0.2325 for $F1@K$ — substantially higher than the next-best corrector (see also Figure 3). Importantly, FC is often best or tied-best in many metric–model combinations, showing strong generalization across diverse LLM architectures. LLM1 and LLM2 show strong performance in this case, likely due to the provenance of the VELI5 questions which were originally derived from public Reddit forums. Therefore, it is plausible that similar content was included in these models’ pre-

training corpora which may confer a significant advantage on this dataset.

Table 10 reinforces these findings by examining factual precision before and after correction. FC achieves the highest post-correction precision on two out of four models (Llama, Mixtral) and ranks second on Granite and GPT-OSS, where it still delivers large relative gains. On average, FC improves precision by +0.21, outperforming all baselines, including LLM1 (+0.17) and RAC (+0.16). While RAC slightly surpasses FC on GPT-OSS in final precision, FC dominates in overall gains, suggesting a more robust correction strategy. In contrast, CRITIC exhibits inconsistent behavior, even reducing precision for GPT-OSS. Finally, Figure 7 plots the factual precision of both the original response and its correction generated by FACTCORRECTOR highlighting the substantial positive gains achieved by the corrector.

BIO. Table 11 presents the mean relative gains and standard deviations obtained on the BIO dataset. FC exhibits robust performance across models, achieving its highest gains with Mixtral and Granite, and ranking second for Llama. In contrast, LLM1 performs poorly, particularly on Granite and GPT-OSS. Both RAC and LLM2 remain highly competitive on this dataset, while CRITIC, as observed previously, struggles to follow instructions, leading to inconsistent results. Figure 8 further illustrates the macro-gains across models, showing that FC slightly outperforms RAC and LLM2. Notably, the BIO dataset is sourced from Wikipedia (Min et al., 2023) — a common component of the pre-training corpora for most modern LLMs, including those evaluated in our experiments. This explains the strong performance of RAC and LLM2 which rely heavily on their internal knowledge. Looking at Table 12 which reports the mean factual precision before and after correction we can see that FC is the best performing corrector on two models (Mixtral, Granite) and second best on Llama. In contrast, GPT-OSS seems to be better suited for RAC and LLM2.

ASKHIST. Table 13 summarizes the mean relative gains and standard deviations obtained on the ASKHIST dataset. Unlike previous experiments, the gains here are considerably smaller. This is primarily because the factual precision of the original responses is already high—typically exceeding 89% on average—leaving limited room for im-

provement through correction. Interestingly, all other baselines exhibit negative gains, indicating that their corrections are less factual than the original responses. In contrast, FC maintains robust performance across all models, almost always producing corrections that are at least as factual as the original responses and, in many cases, substantially better. The results in Table 14 corroborate this observation and further validate the effectiveness of our proposed FC corrector.

In summary, FACTCORRECTOR (FC) consistently delivers the most reliable improvements in factuality, across all datasets. Specifically, on VELI5, FC achieves substantial gains across precision, recall@K, and F1@K, outperforming all baselines and demonstrating strong generalization across models. For BIO, FC remains robust, slightly surpassing RAC and LLM2 despite their advantage from Wikipedia-based pre-training, and achieves top or near-top factual precision on most models. On ASKHIST, where original responses are already highly factual, FC still maintains positive gains, unlike other baselines that often degrade performance. These results highlight FC’s effectiveness and adaptability across diverse datasets and model architectures.

Statistical Significance Tests. We performed one-sided t-tests on the mean relative gain in precision achieved by FC compared to competing methods. The corresponding ppp-values are reported in Table 16. The near-zero ppp-values for the precision metric provide strong statistical evidence that FC consistently and significantly outperforms all baselines across datasets and underlying model configurations. These results underscore the robustness of FC and its ability to deliver substantially more accurate factual corrections than existing approaches.

C.5 Results for SFT Correction

We used the larger VELI5 dataset to train a LoRA adapter on the Granite-Guardian-5B model for generating corrections. The hyperparameters of the LoRA adapter were the following: $r = 64$, $\alpha = 128$, bias = none, lr = 1e-5, schedule = cosine, batch size = 16 and we used all linear layers.

Table 6 reports the mean relative gains on the VELI5 dataset when using the FACTCORRECTOR pipeline versus the LoRA-based correction approach (denoted as SFT). Overall, the SFT method

produces substantially better corrections than the original responses, which is expected given that the evaluation instances originate from the same distribution as the training data.

Tables 17 and 18 present results on the BIO and ASKHIST datasets, which are out-of-distribution relative to the SFT training data. As expected, the gains achieved by SFT are smaller than before; however, most gains remain non-negative, indicating that the SFT approach generalizes reasonably well. This is noteworthy because SFT is significantly more cost-efficient, requiring only a single LLM call per correction, compared to the more complex FACTCORRECTOR pipeline, which involves multiple LLM calls for each instance.

C.6 Human Evaluation of VELI5

We conducted a human evaluation to assess the quality and correctness of VELI5 corrections. Annotators were presented with the original question, the LLM response, a highlighted atom text extracted from the response, relevant contextual information, and the correction. For each instance, annotators made three judgments: (1) whether the incorrect atom had been appropriately corrected given the provided context (labeled as corrected, undercorrected (which also includes not corrected and wrongly corrected), rightly overcorrected, or invalid); (2) whether the factual content of a correct atom was preserved in the correction (labeled as preserved, partially preserved, not preserved, or invalid), and separately rate each atom’s relevance to the question as Relevant, Partially relevant, Irrelevant, Invalid (hallucinated atom), or Invalid (other reasons). Clear annotation guidelines and examples were provided to ensure consistency. This evaluation setup allowed us to measure both factual alignment with external context and faithfulness to the original atomic information, providing a fine-grained assessment of correction quality beyond overall response accuracy.

We design a Label Studio interface to evaluate revisions at multiple granularities. For incorrect or unverified response atoms, annotators judge whether each atom is Corrected, Undercorrected, Overcorrected, Invalid (hallucinated atom), or Invalid (other reasons) (with a free-text rationale when “Invalid (other reasons)” is selected), and separately rate each atom’s relevance to the question as Relevant, Partially relevant, Irrelevant, Invalid (hallucinated atom), or Invalid (other reasons) (again with

| Metric | % |
|-------------------------------------|-----------------|
| Dataset Summary | |
| Tasks (filtered) | 30 |
| Incorrect atoms / task | 2.30 ± 1.39 |
| Correct atoms / task | 4.40 ± 2.66 |
| Task-level relevance | |
| Response relevant | 90.0% |
| Response partially relevant | 10.0% |
| Response irrelevant | 0.0% |
| Correction relevant | 96.7% |
| Correction partially relevant | 3.3% |
| Correction irrelevant | 0.0% |
| Incorrect atoms — correction | |
| Corrected | 67.4 ± 38.5 |
| Undercorrected | 26.4 ± 36.8 |
| Overcorrected | 3.3 ± 18.3 |
| Invalid | 2.9 ± 12.6 |
| Incorrect atoms — relevance | |
| Relevant | 75.0 ± 43.1 |
| Partially relevant | 3.3 ± 18.3 |
| Irrelevant | 5.0 ± 20.1 |
| Correct atoms — preservation | |
| Preserved | 60.9 ± 34.9 |
| Partially preserved | 9.2 ± 14.8 |
| Not preserved | 27.7 ± 29.5 |
| Invalid | 2.3 ± 5.3 |
| Correct atoms — relevance | |
| Relevant | 72.4 ± 40.7 |
| Partially relevant | 7.9 ± 17.1 |
| Irrelevant | 3.0 ± 13.2 |

Table 7: Atom-level evaluation and task-level relevance (percent mean \pm std).

an optional rationale for the latter). For correct response atoms, annotators assess information preservation in the correction as Preserved, Partially preserved, Not preserved, Invalid (hallucinated atom), or Invalid (other reasons) (with rationale text for “Invalid (other reasons)”), and also label each correct atom’s relevance to the question using the same five-way relevance scale. Finally, at the response level, annotators label model-response relevance and correction relevance to the question as Relevant, Partially Relevant, or Irrelevant, providing an end-to-end check that the revision remains on-topic while atom-level labels capture whether factual errors were fixed and correct content was retained.

Table 7 summarizes the results of our human evaluation conducted on 30 filtered tasks (or instances) from our VELI5 dataset. Each task comprises on average 2.30 incorrect atoms and 4.40 correct atoms. At the task level, the original responses demonstrated high relevance with respect to the user question, with 90% classified as fully relevant and 10% as partially relevant, while corrections achieved even greater relevance to the question,

| corrector | ROUGE \uparrow | BLEU \uparrow | BLEURT \uparrow | JUDGE \uparrow |
|------------------------|---------------------------------|---------------------------------|---------------------------------|------------------|
| mixtral-8x22b-instruct | | | | |
| CRITIC | 0.14 \pm 0.07 | 0.02 \pm 0.03 | -0.69 \pm 0.25 | 0.41 |
| RAC | 0.54 \pm 0.23 | 0.23 \pm 0.27 | 0.29 \pm 0.43 | 0.87 |
| LLM1 | 0.15 \pm 0.05 | 0.02 \pm 0.02 | -0.63 \pm 0.21 | 0.30 |
| LLM2 | 0.20 \pm 0.10 | 0.05 \pm 0.05 | -0.48 \pm 0.22 | 0.10 |
| FC (ours) | 0.89\pm0.26 | 0.87\pm0.32 | 0.73\pm0.47 | 0.87 |
| llama-3.3-70b-instruct | | | | |
| CRITIC | 0.32 \pm 0.27 | 0.15 \pm 0.28 | -0.36 \pm 0.56 | 0.31 |
| RAC | 0.87\pm0.20 | 0.77\pm0.32 | 0.65\pm0.43 | 0.77 |
| LLM1 | 0.19 \pm 0.11 | 0.03 \pm 0.05 | -0.62 \pm 0.33 | 0.17 |
| LLM2 | 0.21 \pm 0.10 | 0.04 \pm 0.04 | -0.46 \pm 0.22 | 0.06 |
| FC (ours) | <u>0.73\pm0.33</u> | <u>0.62\pm0.45</u> | <u>0.46\pm0.56</u> | <u>0.60</u> |
| granite-4.0-h-small | | | | |
| CRITIC | 0.15 \pm 0.11 | 0.03 \pm 0.10 | -0.69 \pm 0.26 | 0.20 |
| RAC | 0.74 \pm 0.23 | 0.48 \pm 0.37 | 0.49 \pm 0.46 | <u>0.73</u> |
| LLM1 | 0.15 \pm 0.07 | 0.02 \pm 0.02 | -0.68 \pm 0.23 | 0.17 |
| LLM2 | 0.42 \pm 0.24 | 0.15 \pm 0.17 | -0.10 \pm 0.44 | 0.13 |
| FC (ours) | 0.83\pm0.30 | 0.76\pm0.39 | 0.58\pm0.57 | 0.78 |
| gpt-oss-120b | | | | |
| CRITIC | 0.12 \pm 0.07 | 0.01 \pm 0.01 | -0.89 \pm 0.22 | 0.09 |
| RAC | 0.85\pm0.17 | 0.60\pm0.33 | 0.74\pm0.29 | 0.89 |
| LLM1 | 0.08 \pm 0.03 | 0.00 \pm 0.01 | -0.83 \pm 0.15 | 0.37 |
| LLM2 | 0.16 \pm 0.06 | 0.02 \pm 0.02 | -0.54 \pm 0.18 | 0.19 |
| FC (ours) | <u>0.72\pm0.37</u> | <u>0.66\pm0.46</u> | <u>0.44\pm0.67</u> | <u>0.66</u> |
| granite-4.0-h-tiny | | | | |
| CRITIC | 0.55 \pm 0.43 | 0.49 \pm 0.47 | 0.08 \pm 0.81 | 0.77 |
| RAC | 0.83 \pm 0.18 | 0.59 \pm 0.33 | 0.67 \pm 0.36 | <u>0.84</u> |
| LLM1 | 0.22 \pm 0.19 | 0.08 \pm 0.19 | -0.44 \pm 0.40 | 0.40 |
| LLM2 | 0.38 \pm 0.29 | 0.15 \pm 0.21 | -0.23 \pm 0.47 | 0.29 |
| FC (ours) | 0.97\pm0.15 | 0.95\pm0.19 | 0.86\pm0.28 | 0.97 |

Table 8: Results obtained on the CONFLICTS dataset. with 96.7% fully relevant. For incorrect atoms, 67.4% were successfully corrected by FACTCORRECTOR, though 26.4% remained undercorrected, and only minor proportions were overcorrected (3.3%) or invalid (2.9%). Correct atoms in the original response were preserved by the correction in 60.9% of cases, with 27.7% not preserved and small fractions partially preserved or invalid. Relevance analysis indicated that both incorrect and correct atoms were predominantly relevant to the user question (75.0% and 72.4%, respectively), underscoring the overall effectiveness of the correction process.

D Prompts

In this section, we present the prompts used in our experiments. Specifically, Figure 11 shows the prompt employed by the LLM-as-a-Judge model used in our experiments with the CONFLICTS

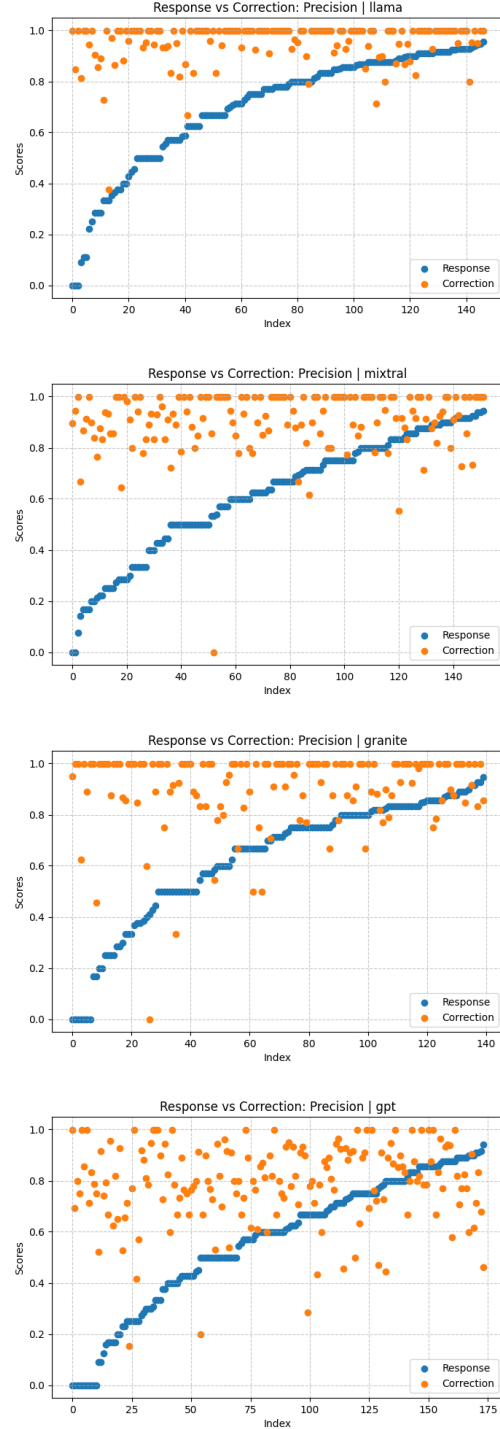


Figure 7: Factual precision: response vs correction by FACTCORRECTOR on the VELI5 dataset.

dataset. Figures 12 and 13 show the prompts employed by the LLM1 and LLM2 baselines. Figure 14 shows the prompt used by FACTCORRECTOR’s refinement model (i.e., Corrector) to generate the corrected response based on the critic’s feedback. The latter consists of the incorrect atoms, the contexts connected to them by an edge in the graphical model, together with the relationship types

| corrector | Pr \uparrow | R@K \uparrow | F1@K \uparrow | V \uparrow | C \uparrow |
|------------------------|---------------------------------|---------------------------------|---------------------------------|---------------------------------|---------------------------------|
| mixtral-8x22b-instruct | | | | | |
| CRITIC | 0.26 \pm 0.49 | 0.08 \pm 0.33 | 0.19 \pm 0.41 | 0.09 \pm 0.32 | 0.06 \pm 0.18 |
| RAC | 0.24 \pm 0.49 | 0.08 \pm 0.33 | 0.18 \pm 0.41 | 0.08 \pm 0.29 | 0.06 \pm 0.20 |
| LLM1 | 0.32 \pm 0.45 | 0.10 \pm 0.30 | 0.23 \pm 0.37 | 0.12 \pm 0.30 | 0.07 \pm 0.17 |
| LLM2 | 0.32 \pm 0.45 | 0.10 \pm 0.30 | 0.23 \pm 0.38 | 0.13 \pm 0.31 | 0.07 \pm 0.17 |
| FC (ours) | 0.34\pm0.46 | 0.09 \pm 0.34 | 0.24\pm0.40 | 0.14\pm0.29 | 0.07\pm0.16 |
| SFT (ours) | 0.19\pm0.51 | 0.10 \pm 0.41 | 0.16\pm0.45 | 0.05\pm0.33 | 0.08\pm0.23 |
| llama-3.3-70b-instruct | | | | | |
| CRITIC | 0.19 \pm 0.41 | 0.06 \pm 0.29 | 0.13 \pm 0.34 | 0.03 \pm 0.17 | 0.05 \pm 0.25 |
| RAC | 0.13 \pm 0.37 | 0.02 \pm 0.22 | 0.09 \pm 0.29 | 0.02 \pm 0.08 | 0.05 \pm 0.20 |
| LLM1 | 0.25 \pm 0.40 | 0.06 \pm 0.29 | 0.18 \pm 0.34 | 0.05 \pm 0.16 | 0.08 \pm 0.24 |
| LLM2 | 0.26 \pm 0.40 | 0.07 \pm 0.29 | 0.18 \pm 0.34 | 0.05 \pm 0.16 | 0.08 \pm 0.24 |
| FC (ours) | 0.27\pm0.39 | 0.07\pm0.29 | 0.18\pm0.34 | 0.05\pm0.16 | 0.08\pm0.23 |
| SFT (ours) | 0.25\pm0.49 | 0.11\pm0.41 | 0.19\pm0.44 | 0.10\pm0.32 | 0.08\pm0.23 |
| granite-4.0-h-small | | | | | |
| CRITIC | 0.17 \pm 0.55 | 0.04 \pm 0.43 | 0.12 \pm 0.49 | 0.05 \pm 0.38 | 0.06 \pm 0.23 |
| RAC | 0.20 \pm 0.49 | 0.09 \pm 0.35 | 0.16 \pm 0.41 | 0.06 \pm 0.26 | 0.07 \pm 0.23 |
| LLM1 | 0.31\pm0.48 | 0.11\pm0.40 | 0.23\pm0.44 | 0.11\pm0.31 | 0.08\pm0.22 |
| LLM2 | 0.20 \pm 0.53 | 0.10 \pm 0.42 | 0.16 \pm 0.47 | 0.09 \pm 0.32 | 0.05 \pm 0.25 |
| FC (ours) | 0.29 \pm 0.49 | 0.10 \pm 0.43 | 0.21 \pm 0.45 | 0.11 \pm 0.32 | 0.07 \pm 0.23 |
| SFT (ours) | 0.22 \pm 0.50 | 0.11 \pm 0.40 | 0.18 \pm 0.45 | 0.09 \pm 0.32 | 0.07 \pm 0.24 |
| gpt-oss-120b | | | | | |
| CRITIC | -0.04 \pm 0.62 | 0.01 \pm 0.54 | -0.02 \pm 0.57 | -0.01 \pm 0.30 | -0.03 \pm 0.33 |
| RAC | 0.34 \pm 0.60 | 0.15 \pm 0.48 | 0.27 \pm 0.53 | 0.06 \pm 0.31 | 0.13 \pm 0.34 |
| LLM1 | 0.21 \pm 0.64 | 0.16 \pm 0.51 | 0.19 \pm 0.57 | -0.01 \pm 0.30 | 0.13 \pm 0.34 |
| LLM2 | 0.33 \pm 0.64 | 0.20 \pm 0.52 | 0.28 \pm 0.57 | 0.05 \pm 0.29 | 0.13 \pm 0.34 |
| FC (ours) | 0.36\pm0.58 | 0.19\pm0.51 | 0.30\pm0.54 | 0.07\pm0.26 | 0.14 \pm 0.33 |
| SFT (ours) | 0.08\pm0.57 | 0.08\pm0.45 | 0.08\pm0.51 | 0.07\pm0.34 | 0.02 \pm 0.26 |

Table 9: Relative Gains for the factuality metrics obtained on the VELI5 dataset. We show the mean and standard deviation as well as highlight the best performance.

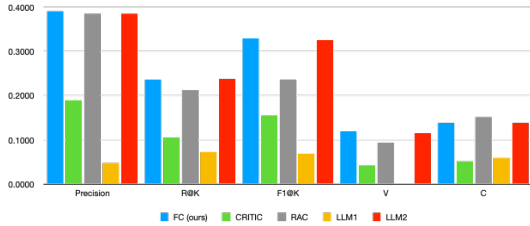


Figure 8: Mean relative gains (macro) for factuality metrics across models on the B10 dataset.

corresponding to those edges (i.e., entailment or contradiction).

E VELI5 Benchmark Examples and Details

Figure 15 shows the prompt we used to generate factually incorrect responses for our VELI5 benchmark dataset.

Figures 16 illustrate the word count distributions for human responses versus corrections (top) and synthetic responses versus corrections (middle).

Figure 16 (bottom) shows the overall distribution for responses compared to their corrections. Across all cases, corrections tend to be longer than the original responses.

Tables 19 and 20 show two illustrative examples from the VELI5 dataset comprising a human-authored as well as a synthetic response together with their corresponding corrections.

| corrector | llama-3.3-70b-instruct | | mixtral-8x22b-instruct | | granite-4.0-h-small | | gpt-oss-120b | |
|-----------|------------------------|------------------|------------------------|------------------|---------------------|------------------|--------------|------------------|
| | before | after | before | after | before | after | before | after |
| CRITIC | 0.78±0.24 | 0.90±0.15 | 0.70±0.25 | 0.88±0.18 | 0.74±0.26 | 0.86±0.24 | 0.62±0.28 | 0.58±0.26 |
| RAC | 0.78±0.24 | 0.91±0.14 | 0.70±0.25 | 0.86±0.18 | 0.74±0.26 | 0.88±0.17 | 0.62±0.28 | 0.85±0.20 |
| LLM1 | 0.78±0.24 | 0.95±0.07 | 0.70±0.25 | 0.91±0.13 | 0.74±0.26 | 0.94±0.10 | 0.62±0.28 | 0.72±0.18 |
| LLM2 | 0.78±0.24 | 0.96±0.07 | 0.70±0.25 | 0.91±0.11 | 0.74±0.26 | 0.86±0.17 | 0.62±0.28 | 0.79±0.18 |
| FC (ours) | 0.78±0.24 | 0.97±0.07 | 0.70±0.25 | 0.93±0.11 | 0.74±0.26 | 0.93±0.13 | 0.62±0.28 | 0.82±0.16 |

Table 10: Mean factual precision and standard deviation before and after correction on the VELI5 dataset.

| corrector | Pr ↑ | R@K ↑ | F1@K ↑ | V ↑ | C ↑ |
|------------------------|------------------|------------------|------------------|------------------|------------------|
| mixtral-8x22b-instruct | | | | | |
| CRITIC | 0.27±0.83 | 0.22±0.80 | 0.25±0.82 | 0.05±0.32 | 0.05±0.41 |
| RAC | 0.39±0.68 | 0.29±0.68 | 0.34±0.68 | 0.05±0.38 | 0.13±0.35 |
| LLM1 | 0.27±0.66 | 0.25±0.64 | 0.27±0.65 | 0.02±0.33 | 0.07±0.33 |
| LLM2 | 0.42±0.72 | 0.30±0.72 | 0.37±0.72 | 0.11±0.29 | 0.07±0.41 |
| FC (ours) | 0.46±0.70 | 0.33±0.72 | 0.41±0.70 | 0.12±0.27 | 0.11±0.36 |
| SFT (ours) | 0.05±0.42 | 0.00±0.34 | 0.02±0.38 | 0.05±0.27 | 0.02±0.11 |
| llama-3.3-70b-instruct | | | | | |
| CRITIC | 0.37±0.61 | 0.18±0.59 | 0.29±0.59 | 0.08±0.31 | 0.13±0.50 |
| RAC | 0.42±0.58 | 0.21±0.58 | 0.33±0.57 | 0.12±0.29 | 0.17±0.44 |
| LLM1 | 0.23±0.64 | 0.14±0.62 | 0.19±0.62 | 0.07±0.30 | 0.11±0.44 |
| LLM2 | 0.45±0.57 | 0.21±0.57 | 0.36±0.56 | 0.12±0.28 | 0.18±0.44 |
| FC (ours) | 0.44±0.56 | 0.22±0.57 | 0.35±0.56 | 0.12±0.28 | 0.17±0.44 |
| SFT (ours) | 0.01±0.39 | -0.02±0.29 | -0.01±0.34 | 0.06±0.28 | -0.03±0.18 |
| granite-4.0-h-small | | | | | |
| CRITIC | 0.12±0.43 | 0.02±0.29 | 0.08±0.36 | 0.04±0.24 | 0.02±0.12 |
| RAC | 0.22±0.41 | 0.08±0.35 | 0.16±0.37 | 0.10±0.31 | 0.04±0.13 |
| LLM1 | -0.03±0.31 | -0.01±0.25 | -0.02±0.28 | 0.02±0.21 | 0.00±0.12 |
| LLM2 | 0.19±0.58 | 0.13±0.44 | 0.16±0.50 | 0.13±0.39 | 0.03±0.15 |
| FC (ours) | 0.26±0.44 | 0.10±0.38 | 0.20±0.40 | 0.14±0.33 | 0.03±0.13 |
| SFT (ours) | 0.05±0.41 | 0.00±0.36 | 0.02±0.38 | 0.06±0.27 | -0.02±0.26 |
| gpt-oss-120b | | | | | |
| CRITIC | 0.00±0.38 | 0.00±0.34 | 0.00±0.36 | 0.00±0.28 | 0.01±0.35 |
| RAC | 0.51±0.66 | 0.28±0.69 | 0.42±0.66 | 0.11±0.24 | 0.28±0.60 |
| LLM1 | -0.28±0.92 | -0.10±0.86 | -0.22±0.89 | -0.24±0.42 | 0.06±0.71 |
| LLM2 | 0.48±0.67 | 0.31±0.68 | 0.41±0.67 | 0.10±0.30 | 0.27±0.60 |
| FC (ours) | 0.40±0.71 | 0.29±0.70 | 0.36±0.70 | 0.10±0.30 | 0.24±0.61 |
| SFT (ours) | -0.04±0.42 | 0.00±0.34 | -0.02±0.37 | 0.0±0.31 | -0.04±0.21 |

Table 11: Relative Gains for the factuality metrics obtained on the BIO dataset. We show the mean and standard deviation as well as highlight the best performance.

| corrector | llama-3.3-70b-instruct | | mixtral-8x22b-instruct | | granite-4.0-h-small | | gpt-oss-120b | |
|-----------|------------------------|------------------|------------------------|------------------|---------------------|------------------|--------------|------------------|
| | before | after | before | after | before | after | before | after |
| CRITIC | 0.65±0.28 | 0.86±0.19 | 0.64±0.32 | 0.75±0.25 | 0.72±0.23 | 0.80±0.21 | 0.56±0.28 | 0.56±0.28 |
| RAC | 0.65±0.28 | 0.91±0.10 | 0.64±0.32 | 0.84±0.19 | 0.72±0.23 | 0.86±0.14 | 0.56±0.28 | 0.86±0.16 |
| LLM1 | 0.65±0.28 | 0.76±0.23 | 0.64±0.32 | 0.72±0.22 | 0.72±0.23 | 0.71±0.22 | 0.56±0.28 | 0.33±0.17 |
| LLM2 | 0.65±0.28 | 0.93±0.07 | 0.64±0.32 | 0.86±0.18 | 0.72±0.23 | 0.79±0.18 | 0.56±0.28 | 0.79±0.11 |
| FC (ours) | 0.65±0.28 | 0.91±0.11 | 0.64±0.32 | 0.87±0.15 | 0.72±0.23 | 0.90±0.12 | 0.56±0.28 | 0.75±0.18 |

Table 12: Mean factual precision and standard deviation before and after correction on the BIO dataset.

| corrector | Pr \uparrow | R@K \uparrow | F1@K \uparrow | V \uparrow | C \uparrow |
|---------------------------|----------------------------------|---------------------------------|----------------------------------|----------------------------------|---------------------------------|
| mixtral-8x22b-instruct | | | | | |
| CRITIC | -0.12 \pm 0.36 | -0.04 \pm 0.21 | -0.09 \pm 0.29 | -0.06 \pm 0.25 | -0.01 \pm 0.09 |
| RAC | -0.01 \pm 0.23 | -0.01 \pm 0.07 | -0.01 \pm 0.15 | -0.01 \pm 0.18 | 0.00 \pm 0.06 |
| LLM1 | -0.02 \pm 0.18 | 0.00 \pm 0.06 | -0.01 \pm 0.12 | 0.00 \pm 0.13 | 0.00 \pm 0.06 |
| LLM2 | -0.03 \pm 0.21 | 0.00 \pm 0.08 | -0.02 \pm 0.15 | 0.00 \pm 0.16 | 0.00 \pm 0.07 |
| FC (ours) | 0.02\pm0.19 | 0.00\pm0.10 | 0.00\pm0.14 | 0.02\pm0.14 | 0.00\pm0.06 |
| SFT (ours) | -0.09\pm0.21 | 0.00\pm0.07 | -0.05\pm0.14 | -0.08\pm0.16 | 0.01\pm0.08 |
| llama-3.3-70b-instruct | | | | | |
| CRITIC | -0.03 \pm 0.20 | 0.00 \pm 0.05 | -0.01 \pm 0.13 | -0.01 \pm 0.13 | -0.01 \pm 0.07 |
| RAC | 0.03 \pm 0.19 | 0.01 \pm 0.05 | 0.02 \pm 0.12 | 0.02 \pm 0.09 | 0.00 \pm 0.05 |
| LLM1 | 0.00 \pm 0.16 | 0.00 \pm 0.05 | 0.00 \pm 0.10 | 0.00 \pm 0.07 | 0.00 \pm 0.06 |
| LLM2 | 0.01 \pm 0.17 | 0.00 \pm 0.05 | 0.01 \pm 0.11 | 0.00 \pm 0.08 | 0.00 \pm 0.05 |
| FC (ours) | 0.04\pm0.15 | 0.01\pm0.04 | 0.02\pm0.10 | 0.01\pm0.07 | 0.01\pm0.05 |
| SFT (ours) | -0.01\pm0.20 | 0.00\pm0.05 | 0.00\pm0.13 | 0.02\pm0.11 | 0.00\pm0.08 |
| granite-4.0-h-small (32b) | | | | | |
| CRITIC | -0.13 \pm 0.45 | -0.09 \pm 0.42 | -0.11 \pm 0.42 | -0.05 \pm 0.29 | -0.04 \pm 0.19 |
| RAC | -0.06 \pm 0.27 | -0.01 \pm 0.08 | -0.04 \pm 0.18 | -0.01 \pm 0.12 | -0.02 \pm 0.10 |
| LLM1 | -0.08 \pm 0.30 | -0.03 \pm 0.24 | -0.05 \pm 0.27 | -0.03 \pm 0.16 | -0.02 \pm 0.15 |
| LLM2 | -0.07 \pm 0.20 | -0.01 \pm 0.05 | -0.04 \pm 0.13 | -0.03 \pm 0.12 | -0.01 \pm 0.07 |
| FC (ours) | 0.00\pm0.16 | 0.00\pm0.05 | 0.00\pm0.11 | 0.01\pm0.09 | 0.00\pm0.06 |
| SFT (ours) | 0.00\pm0.18 | 0.00\pm0.04 | 0.00\pm0.12 | 0.00\pm0.14 | 0.00\pm0.09 |
| gpt-oss-120b | | | | | |
| CRITIC | -0.03 \pm 0.33 | -0.01 \pm 0.24 | -0.02 \pm 0.28 | 0.01 \pm 0.24 | -0.01 \pm 0.18 |
| RAC | 0.12\pm0.37 | 0.03\pm0.24 | 0.08\pm0.30 | 0.06\pm0.26 | 0.02\pm0.13 |
| LLM1 | -0.49 \pm 0.54 | -0.26 \pm 0.48 | -0.41 \pm 0.52 | -0.19 \pm 0.31 | -0.08 \pm 0.20 |
| LLM2 | -0.03 \pm 0.48 | -0.01 \pm 0.34 | -0.03 \pm 0.42 | 0.03 \pm 0.25 | -0.04 \pm 0.22 |
| FC (ours) | 0.00 \pm 0.16 | 0.00 \pm 0.05 | 0.00 \pm 0.11 | 0.01 \pm 0.09 | 0.00 \pm 0.06 |
| SFT (ours) | -0.28 \pm 0.34 | -0.06 \pm 0.21 | -0.19 \pm 0.28 | -0.09 \pm 0.19 | -0.05 \pm 0.10 |

Table 13: Relative Gains for the factuality metrics obtained on the ASKHIST dataset. We show the mean and standard deviation as well as highlight the best performance.

| corrector | llama-3.3-70b-instruct | | mixtral-8x22b-instruct | | granite-4.0-h-small | | gpt-oss-120b | |
|-----------|------------------------|-----------------|------------------------|-----------------|---------------------|-----------------|-----------------|-----------------|
| | before | after | before | after | before | after | before | after |
| CRITIC | 0.89 \pm 0.12 | 0.87 \pm 0.14 | 0.84 \pm 0.13 | 0.77 \pm 0.23 | 0.89 \pm 0.13 | 0.82 \pm 0.22 | 0.69 \pm 0.20 | 0.67 \pm 0.21 |
| RAC | 0.89 \pm 0.12 | 0.92 \pm 0.11 | 0.84 \pm 0.13 | 0.84 \pm 0.15 | 0.89 \pm 0.13 | 0.85 \pm 0.16 | 0.69 \pm 0.20 | 0.77 \pm 0.20 |
| LLM1 | 0.89 \pm 0.12 | 0.89 \pm 0.12 | 0.84 \pm 0.13 | 0.82 \pm 0.14 | 0.89 \pm 0.13 | 0.84 \pm 0.16 | 0.69 \pm 0.20 | 0.42 \pm 0.18 |
| LLM2 | 0.89 \pm 0.12 | 0.89 \pm 0.09 | 0.84 \pm 0.13 | 0.82 \pm 0.15 | 0.89 \pm 0.13 | 0.83 \pm 0.13 | 0.69 \pm 0.20 | 0.67 \pm 0.20 |
| FC (ours) | 0.89 \pm 0.12 | 0.92 \pm 0.10 | 0.84 \pm 0.13 | 0.87 \pm 0.15 | 0.89 \pm 0.13 | 0.90 \pm 0.14 | 0.69 \pm 0.20 | 0.67 \pm 0.20 |

Table 14: Mean factual precision and standard deviation before and after correction on the ASKHIST dataset.

| dataset | response | correction | | | | | |
|------------------------|----------|------------|------|-------|-------|-----------|------------|
| | | CRITIC | RAC | LLM1 | LLM2 | FC (ours) | SFT (ours) |
| mixtral-8x22b-instruct | | | | | | | |
| VELI5 | 8±4 | 10±3 | 10±4 | 13±4 | 13±4 | 15±13 | 12±5 |
| BIO | 14±7 | 15±7 | 15±6 | 16±5 | 18±8 | 18±7 | 17±6 |
| ASKHIST | 17±4 | 11±4 | 14±5 | 16±4 | 17±5 | 18±15 | 15±4 |
| llama-3.3-70b-instruct | | | | | | | |
| VELI5 | 11±5 | 12±4 | 12±5 | 13±4 | 18±5 | 19±6 | 14±6 |
| BIO | 18±7 | 17±9 | 17±6 | 22±8 | 28±8 | 21±7 | 20±8 |
| ASKHIST | 21±5 | 15±6 | 16±6 | 23±6 | 21±6 | 20±8 | 18±5 |
| granite-4.0-h-small | | | | | | | |
| VELI5 | 8±5 | 8±5 | 8±5 | 15±7 | 11±5 | 10±6 | 11±6 |
| BIO | 18±6 | 17±7 | 15±5 | 21±6 | 23±8 | 17±6 | 18±5 |
| ASKHIST | 17±5 | 9±5 | 14±5 | 18±6 | 16±5 | 15±6 | 15±5 |
| gpt-oss-120b | | | | | | | |
| VELI5 | 8±4 | 8±5 | 10±8 | 27±12 | 16±6 | 15±7 | 11±6 |
| BIO | 16±7 | 16±7 | 14±8 | 36±13 | 36±13 | 28±11 | 18±7 |
| ASKHIST | 18±5 | 18±7 | 12±7 | 42±17 | 17±8 | 22±8 | 16±7 |

Table 15: Mean number of atoms and standard deviation for the response and corrections from different correctors.

| model | FC vs. CRITIC | FC vs. RAC | FC vs. LLM1 | FC vs. LLM2 |
|-----------------|---------------|---------------|---------------|---------------|
| VELI5 | | | | |
| mixtral-8x22b | 0.0569 | 0.0277 | 0.3604 | 0.3851 |
| llama-3.3-70b | 0.0233 | 0.0005 | 0.3719 | 0.4219 |
| granite-4-small | 0.0134 | 0.0468 | 0.6798 | 0.0401 |
| gpt-oss-120b | 0.0000 | 0.3935 | 0.0098 | 0.3084 |
| BIO | | | | |
| mixtral-8x22b | 0.0207 | 0.1933 | 0.0103 | 0.3397 |
| llama-3.3-70b | 0.1449 | 0.3867 | 0.0005 | 0.6094 |
| granite-4-small | 0.0016 | 0.1792 | 0.0000 | 0.1365 |
| gpt-oss-120b | 0.0000 | 0.9370 | 0.0000 | 0.8449 |
| ASKHIST | | | | |
| mixtral-8x22b | 0.0000 | 0.1252 | 0.0202 | 0.0098 |
| llama-3.3-70b | 0.0002 | 0.3812 | 0.0181 | 0.0614 |
| granite-4-small | 0.0007 | 0.0093 | 0.0031 | 0.0005 |
| gpt-oss-120b | 0.3759 | 0.9998 | 0.0000 | 0.3468 |

Table 16: Statistical significance tests: p -values for $G(Pr)$ on the datasets.

| corrector | Pr ↑ | R@K ↑ | F1@K ↑ | V ↑ | C ↑ |
|------------------------|-------------|-------------|-------------|-------------|-------------|
| mixtral-8x22b-instruct | | | | | |
| FC (ours) | 0.46 | 0.33 | 0.41 | 0.12 | 0.11 |
| SFT (ours) | 0.05 | 0.00 | 0.02 | 0.05 | 0.02 |
| llama-3.3-70b-instruct | | | | | |
| FC (ours) | 0.44 | 0.22 | 0.35 | 0.12 | 0.17 |
| SFT (ours) | 0.01 | -0.02 | -0.01 | 0.06 | -0.03 |
| granite-4.0-h-small | | | | | |
| FC (ours) | 0.26 | 0.10 | 0.20 | 0.14 | 0.03 |
| SFT (ours) | 0.05 | 0.00 | 0.02 | 0.06 | -0.02 |
| gpt-oss-120b | | | | | |
| FC (ours) | 0.40 | 0.29 | 0.36 | 0.10 | 0.24 |
| SFT (ours) | -0.04 | 0.00 | -0.02 | 0.00 | -0.04 |

Table 17: Mean relative gains for the factuality metrics on the BIO dataset using the SFT corrector.

| corrector | Pr ↑ | R@K ↑ | F1@K ↑ | V ↑ | C ↑ |
|------------------------|-------------|-------------|-------------|-------------|-------------|
| mixtral-8x22b-instruct | | | | | |
| FC (ours) | 0.02 | 0.00 | 0.00 | 0.02 | 0.00 |
| SFT (ours) | -0.09 | 0.00 | -0.05 | -0.08 | 0.01 |
| llama-3.3-70b-instruct | | | | | |
| FC (ours) | 0.04 | 0.01 | 0.02 | 0.01 | 0.01 |
| SFT (ours) | -0.01 | 0.00 | 0.00 | 0.02 | 0.00 |
| granite-4.0-h-small | | | | | |
| FC (ours) | 0.00 | 0.00 | 0.00 | 0.01 | 0.00 |
| SFT (ours) | 0.00 | 0.00 | 0.00 | 0.00 | 0.00 |
| gpt-oss-120b | | | | | |
| FC (ours) | 0.00 | 0.00 | 0.00 | 0.01 | 0.00 |
| SFT (ours) | -0.28 | -0.06 | -0.19 | -0.09 | -0.05 |

Table 18: Mean relative gains for the factuality metrics on the ASKHIST dataset using the SFT corrector.

| VELI5 | |
|-------------------|---|
| Question | If the sun’s light was polarized, what would a rainbow look like? |
| Response | If the sun’s light was polarized, a rainbow would not form at all. This is because polarized light only oscillates in one direction, while the formation of a rainbow requires light to be scattered in multiple directions. The prismatic effect that creates a rainbow would not occur with polarized light. |
| Correction | If the sun’s light was polarized, a rainbow would still form, but its appearance might change. The light forming a rainbow gets polarized when it reflects off the back of a raindrop. This polarization does not prevent the rainbow from forming. Instead, when viewed through polarized sunglasses, the rainbow can disappear or reappear depending on the orientation of the sunglasses. The polarization of the sun’s light itself does not prevent a rainbow from forming. The light from a rainbow is polarized due to the reflection part of the process, where light transitions between two optically different mediums. Some of the light is transmitted and some is reflected, with the polarization component parallel to the surface being more strongly reflected. This leads to a preference for one polarization orientation for the light leaving the back of the drop and a preference for an orthogonal orientation for the light passing through the drop. The rainbow is polarized up to 94%, with the polarization arising at the internal reflection in the water drop near the Brewster angle. The plane of polarization of the rainbow is tangent to the rainbow arc. |

Table 19: Example entry in VELI5, displaying the question, original answer and correction generated by FACTCORRECTOR using Mixtral-8x22b-Instruct and Google Retriever. F1@K = 0.97

| VELI5 | |
|-------------------|---|
| Question | How did Satoshi Nakamoto, the creator of Bitcoin, stay anonymous? |
| Response | Satoshi Nakamoto, the creator of Bitcoin, stayed anonymous by using his real name and personal information in all his communications and publications related to Bitcoin. He also frequently appeared in public events and gave interviews to media outlets, discussing his work on Bitcoin in detail. Despite this, no one was able to track him down or discover his true identity. |
| Correction | Satoshi Nakamoto, the creator of Bitcoin, stayed anonymous by using a pseudonym and never revealing their true identity. They communicated via email and online forums but did not give interviews to media outlets or appear in public events. Satoshi’s involvement with Bitcoin ended in 2010, and since then, various individuals have claimed to be Satoshi Nakamoto, but none have provided conclusive evidence to support their claims. The true identity of Satoshi Nakamoto remains a mystery. |

Table 20: Example entry in VELI5, displaying the question, original answer and correction generated by FACTCORRECTOR using Mixtral-8x22b-Instruct and Google Retriever. F1@K = 0.32

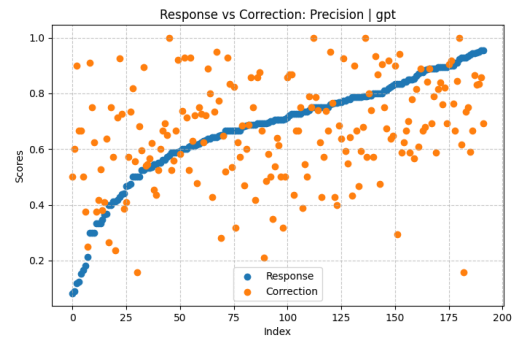
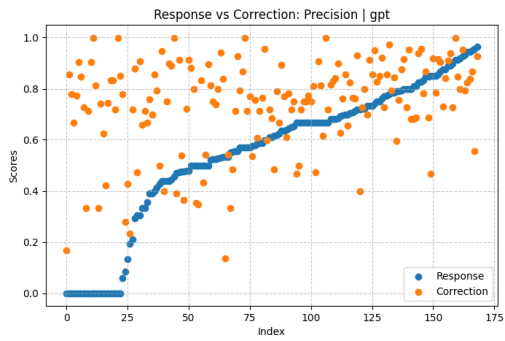
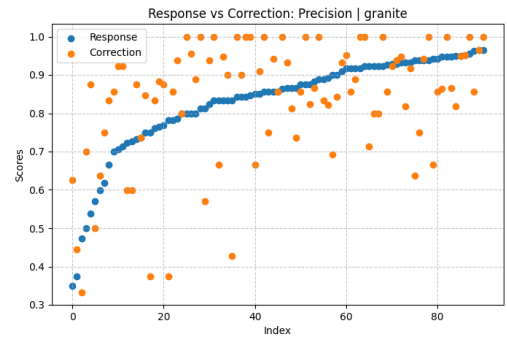
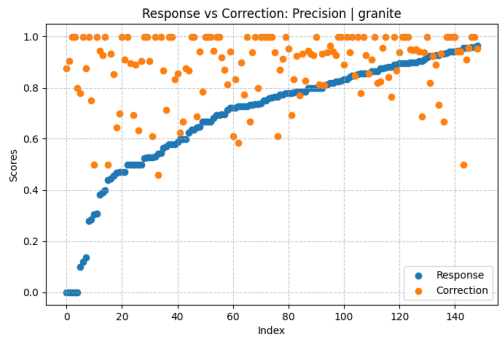
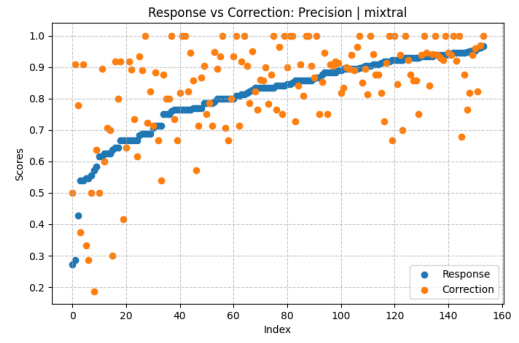
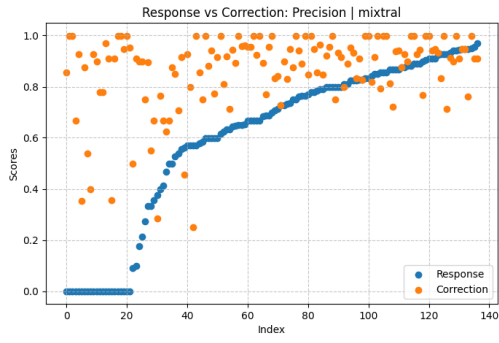
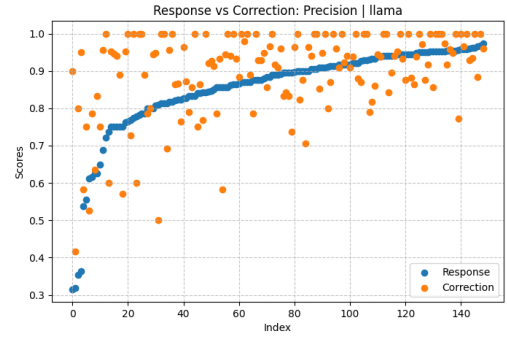
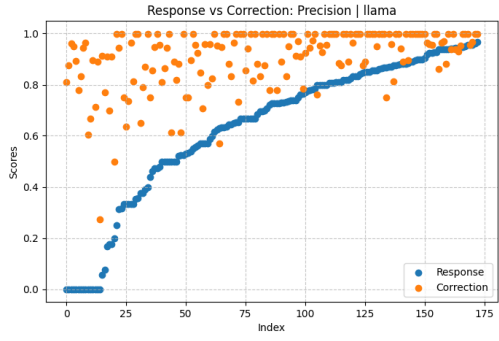


Figure 9: Factual precision: response vs correction by FACTCORRECTOR on the B10 dataset.

Figure 10: Factual precision: response vs correction by FACTCORRECTOR on the ASKHIST dataset.

LLM-as-a-Judge prompt

You are an expert evaluator. Your task is to decide whether a candidate statement is semantically equivalent to a reference statement.

Follow these rules:

- Analyze meaning, not just wording.
- Consider synonyms, paraphrasing, and logical equivalence.
- If the candidate changes the meaning, adds or removes critical information, or introduces contradictions, classify as [No].
- If the candidate preserves the meaning fully, classify as [Yes].
- Provide a short reasoning before the final answer.
- Output must be exactly one of: [Yes] or [No].

Here are 5 examples:

Example 1:

Reference: "The cat is sleeping on the mat."

Candidate: "A cat is lying on a mat."

Reasoning: Both describe the same situation with minor wording differences. Meaning is preserved.

Answer: [Yes]

Example 2:

Reference: "The meeting starts at 10 AM."

Candidate: "The meeting begins at 10 in the morning."

Reasoning: 'Starts' and 'begins' are synonyms; '10 AM' equals '10 in the morning'. Meaning is identical.

Answer: [Yes]

Example 3:

Reference: "She bought a red car yesterday."

Candidate: "She bought a car yesterday."

Reasoning: The candidate omits the color 'red', which is a critical detail. Meaning is not fully preserved.

Answer: [No]

Example 4:

Reference: "Water boils at 100 degrees Celsius at sea level."

Candidate: "Water freezes at 100 degrees Celsius at sea level."

Reasoning: 'Boils' vs 'freezes' completely changes the meaning. Contradiction detected.

Answer: [No]

Example 5:

Reference: "The company announced a new product launch in June."

Candidate: "In June, the company announced a new product launch."

Reasoning: Same event, same timing, just reordered words. Meaning is preserved.

Answer: [Yes]

Now evaluate the following:

Reference: "{}"

Candidate: "{}"

Reasoning:

Figure 11: Prompt template used by the LLM-as-a-Judge model (DeepSeek-v3.2)

LLM1 corrector prompt

Instructions:

You are provided with a QUESTION and an ORIGINAL ANSWER.

Your task is to provide a coherent and factually CORRECTED ANSWER for the QUESTION based on your internal knowledge.

Do not copy the ORIGINAL ANSWER in your CORRECTED ANSWER.

QUESTION: {}

ORIGINAL ANSWER: {}

CORRECTED ANSWER:

Figure 12: Prompt template used by the LLM1 corrector baseline.

LLM2 corrector prompt

Instructions:

You are provided with a QUESTION, a set of CONTEXTS FOR QUESTION, a set of UNVERIFIED ATOMS that contain pieces of information of the ORIGINAL ANSWER that might be unverified, and an ORIGINAL ANSWER.

Your task is to provide a coherent and factually CORRECTED ANSWER for the QUESTION by correcting the UNVERIFIED ATOMS of the ORIGINAL ANSWER based on the given CONTEXTS FOR QUESTION.

Carefully investigate the given CONTEXTS FOR QUESTION and provide a coherent CORRECTED ANSWER that reflects the comprehensive view of the CONTEXTS FOR QUESTION, even if the CORRECTED ANSWER contains contradictory information reflecting the heterogeneous nature of the CONTEXTS FOR QUESTION. In the CORRECTED ANSWER, do not copy the ORIGINAL ANSWER and do not mention that the CORRECTED ANSWER contradicts the ORIGINAL ANSWER, but only correct the ORIGINAL ANSWER according to the instructions provided.

Do not use your internal knowledge, common knowledge, or general knowledge to correct the ORIGINAL ANSWER, but only use the instructions provided. If some UNVERIFIED ATOMS cannot be proven or disproven by the CONTEXTS FOR QUESTION, you must remove those UNVERIFIED ATOMS from the CORRECTED ANSWER.

Learn your task from the following examples:

EXAMPLE 1:

QUESTION: "How many siblings does George have?"

CONTEXTS: "George has three siblings: Michael, Sarah, Emily, and David."

ORIGINAL ANSWER: "George has three siblings."

UNVERIFIED ATOM 1: "George has three siblings."

CORRECTED ANSWER: "Either George has three siblings, or George has four siblings."

EXAMPLE 2:

QUESTION: "What is the surface area of the Pacific Ocean?"

CONTEXTS: "Older geographical records list the Pacific Ocean's surface area as 155 million square kilometers."

"Updated measurements suggest the Pacific Ocean's area is closer to 168 million square kilometers."

"Depending on tidal variations and sea level changes, estimates of the Pacific Ocean's area fluctuate between 155 and 170 million square kilometers."

"Satellite-based ocean mapping has revised the estimate of the Pacific Ocean's area multiple times due to technological improvements."

"Due to climate change and sea-level rise, the Pacific Ocean's surface area is gradually increasing."

ORIGINAL ANSWER: "The Pacific Ocean covers an area of exactly 155 million square kilometers, making it the largest ocean on Earth, and recent satellite data consistently confirm this precise measurement without significant variation."

UNVERIFIED ATOMS: "The Pacific Ocean covers an area of exactly 155 million square kilometers."

"Recent satellite data consistently confirm the precise measurement of the Pacific Ocean without significant variation."

CORRECTED ANSWER: "The Pacific Ocean is reported as covering both 155 million and 168 million square kilometers, with variations depending on measurement techniques, tidal influences, and evolving mapping technologies."

QUESTION: {}

CONTEXTS: {}

ORIGINAL ANSWER: {}

UNVERIFIED ATOMS: {}

CORRECTED ANSWER:

Figure 13: Prompt template used by the LLM2 corrector baseline.

FACTCORRECTOR's refinement model prompt

Instructions:

You are provided with a QUESTION, an optional set of CONTEXTS FOR QUESTION, an ORIGINAL ANSWER, a set of INCORRECT ATOMS and/or UNVERIFIED ATOMS that contain pieces of information of the ORIGINAL ANSWER that might be incorrect or unverified and, for each INCORRECT ATOM, an optional set of CONTEXTS FOR INCORRECT ATOM that might CONTRADICT, ENTAIL, or BE EQUIVALENT TO their corresponding INCORRECT ATOM. Your task is to provide a coherent and factually CORRECTED ANSWER for the QUESTION by factually correcting the INCORRECT ATOMS of the ORIGINAL ANSWER based on the given CONTEXTS. Carefully investigate the given CONTEXTS and provide a coherent CORRECTED ANSWER that reflects the comprehensive view of the CONTEXTS, even if the CORRECTED ANSWER contains contradictory information, reflecting the heterogeneous nature of the CONTEXTS. In the CORRECTED ANSWER, do not copy the ORIGINAL ANSWER and do not mention that CORRECTED ANSWER contradicts ORIGINAL ANSWER, but only provide the CORRECTED ANSWER according to the instructions provided. Do not use your internal knowledge, common knowledge, or general knowledge to correct the ORIGINAL ANSWER, but only use the instructions provided. If some UNVERIFIED ATOMS cannot be proven or disproven by the CONTEXTS FOR QUESTION, you must remove those UNVERIFIED ATOMS from the CORRECTED ANSWER.

EXAMPLE 1:

QUESTION: "How many siblings does George have?"

ORIGINAL ANSWER: "George has three siblings."

INCORRECT ATOM 1: "George has three siblings."

CONTEXT 1-1 FOR INCORRECT ATOM 1: "George has three siblings: Michael, Sarah, Emily, and David."

RELATION FROM CONTEXT 1-1 TO INCORRECT ATOM 1: "CONTRADICTION"

CORRECTED ANSWER: "Either George has three siblings, or George has four siblings."

EXAMPLE 2:

QUESTION: "What is the surface area of the Pacific Ocean?"

CONTEXT 1 FOR QUESTION: "The Pacific Ocean encompasses approximately one-third of the Earth's surface."

ORIGINAL ANSWER: "The Pacific Ocean covers an area of exactly 155 million square kilometers, making it the largest ocean on Earth, and recent satellite data consistently confirm this precise measurement without significant variation."

INCORRECT ATOM 1: "The Pacific Ocean covers an area of exactly 155 million square kilometers."

CONTEXT 1-1 FOR INCORRECT ATOM 1: "Older geographical records list the Pacific Ocean's surface area as 155 million square kilometers."

RELATION FROM CONTEXT 1-1 TO INCORRECT ATOM 1: "ENTAILMENT"

CONTEXT 1-2 FOR INCORRECT ATOM 1: "Updated measurements suggest the Pacific Ocean's area is closer to 168 million square kilometers."

RELATION FROM CONTEXT 1-2 TO INCORRECT ATOM 1: "CONTRADICTION"

EXAMPLE 3:

QUESTION: "How fast does Earth rotate?"

CONTEXT 1 FOR QUESTION: "Earth rotates once every 23 hours, 56 minutes and 4 seconds."

ORIGINAL ANSWER: "Earth rotates at a constant speed of 1,000 miles per hour at the equator, ensuring that the length of a day remains exactly 24 hours. This rotational speed is precisely measured by scientific instruments and does not vary under any conditions. Since Earth's rotation has remained unchanged for millions of years, the length of a day has always been the same, and it will continue to be so for the foreseeable future."

INCORRECT ATOM 1: "Earth rotates at a constant speed of 1,000 miles per hour at the equator."

CONTEXT 11 FOR INCORRECT ATOM 1: "Standard geographical references list Earth's rotational speed at approximately 1,000 miles per hour."

RELATION FROM CONTEXT 11 TO INCORRECT ATOM 1: "ENTAILMENT"

CONTEXT 12 FOR INCORRECT ATOM 1: "The actual speed of Earth's rotation is closer to 1,037 miles per hour at the equator, depending on latitude."

RELATION FROM CONTEXT 12 TO INCORRECT ATOM 1: "CONTRADICTION"

INCORRECT ATOM 2: "The length of a day in the Earth remains exactly 24 hours."

CONTEXT 21 FOR INCORRECT ATOM 2: "For practical purposes, a day in the Earth is measured as 24 hours."

RELATION FROM CONTEXT 21 TO INCORRECT ATOM 2: "ENTAILMENT"

CONTEXT 22 FOR INCORRECT ATOM 2: "Earth's rotation is gradually slowing down, meaning that the actual length of a day increases by about 1.7 milliseconds per century."

RELATION FROM CONTEXT 22 TO INCORRECT ATOM 2: "CONTRADICTION"

YOUR TASK:

QUESTION: "{}"

{incorrect atoms}

{contexts for incorrect atoms}

{relations from contexts to incorrect atoms}

CORRECTED ANSWER:

Figure 14: Prompt template used by the FACTCORRECTOR's refinement model (Corrector).

VELI5 creation prompt

Instructions:

You are provided with a QUESTION. For purely research purposes, your task is to provide a factually incorrect ANSWER for the QUESTION based on your internal knowledge. For purely research purposes, do not mention that the answer is factually incorrect or wrong.

QUESTION: {}

ANSWER:

Figure 15: Prompt template used to generate factually incorrect responses for the VELI5 dataset.

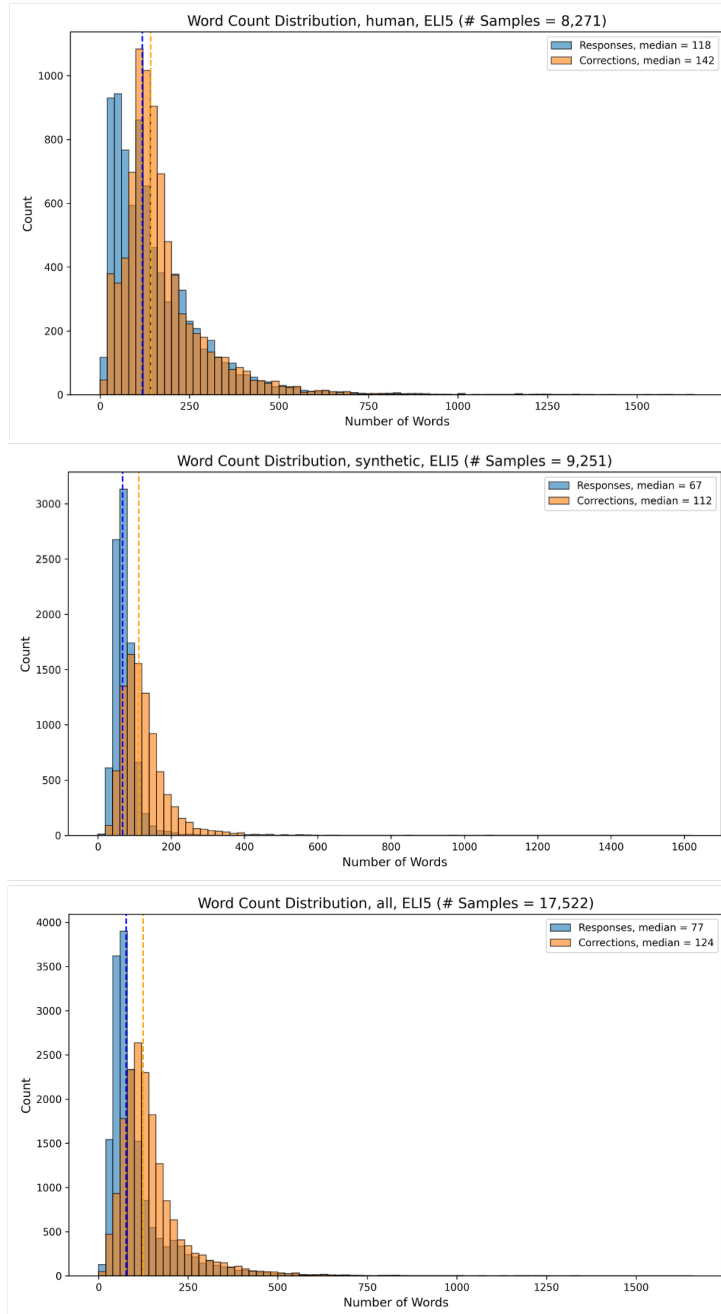


Figure 16: Word Count Distributions for the VELI5 dataset.