
CORE: Robust Factual Precision Scoring with Informative Sub-Claim Identification

Zhengping Jiang Jingyu Zhang Nathaniel Weir
Seth Ebner Miriam Wanner Kate Sanders
Daniel Khashabi Anqi Liu Benjamin Van Durme
Johns Hopkins University
{zjiang31, aliu74, vandurme}@jhu.edu

Abstract

Hallucinations—the generation of untrue claims—pose a challenge to the application of large language models (LLMs) [12] thereby motivating the development of metrics to evaluate factual precision. We observe that popular metrics using the *Decompose-Then-Verify* framework, such as FACTSCORE [22], can be manipulated by adding obvious or repetitive claims to artificially inflate scores. We expand the FACTSCORE dataset to design and analyze factual precision metrics, demonstrating that models can be trained to achieve high scores under existing metrics through exploiting the issues we identify. This motivates our new customizable plug-and-play subclaim selection component called CORE, which filters down individual subclaims according to their uniqueness and informativeness. Metrics augmented by CORE are substantially more robust as shown in head-to-head comparisons. We release an evaluation framework supporting the modular use of CORE¹ and various decomposition strategies, and we suggest its adoption by the LLM community.

1 Introduction

Automatically generating long-form text is prevalent since the rise in large language models (LLMs) [2, 25]. These models are trained on vast amounts of textual data that provide abundant information, enabling them to serve as a significant source of knowledge [26, 27, 29, 43]. A running concern is ensuring LLM-generated content is faithful to its training data, devoid of *hallucination* [14, 12]. To this end, various automatic factuality evaluation pipelines have been proposed [15, 22, 8, 4, 37, 40]. Mainstream methods typically involve two key steps: First, a *decomposition* step, where the generated text is broken down into natural language subclaims, and second, a *verification* step, where a binary factuality label is assigned to each of the subclaims. The proportion of subclaims that can be verified, commonly referred to as *Factual Precision* (FP), serves as the most widely used indicator of factuality level. Throughout this paper, we call this framework *Decompose-Then-Verify*, a concept that has been properly abstracted in previous works [4, 38].

Researchers have sought to improve factuality by optimizing [36] against model-based metrics like FACTSCORE [22]. This raises the question of whether improvements in FP represent genuine factuality gains or instead somehow exploit the evaluation [33]. For example, Figure 1 illustrates that it is trivial to purposefully game FP by including repetitive or less informative generations than we would normally expect from a contemporary LLM. Although it has been noted that LLM evaluation needs to be holistic and multi-faceted [19, 32] beyond FP, popular factuality evaluations put minimal effort into guarding against such malicious inputs designed to inflate FP. Recent studies have already reported that optimizing for factuality can conflict with other desirable objectives, such

¹<https://github.com/zipJiang/Core>.

as *completeness* and *relevancy* [42]. Therefore, we argue that accurate FP evaluation requires more precise control over each design component of the pipeline.

To address this issue, we introduce CORE, a refinement to the *decomposition* step that credits only subclaims that are factual, informative, and non-repetitive; the *core facts*.² This is achieved by weighting each subclaim with its level of uncertainty or surprisal and then selecting the best compatible subset through combinatorial optimization. We demonstrate that our approach makes it more difficult to trivially optimize against *Decompose-Then-Verify* frameworks [4, 38]. Thus, CORE can serve as a plug-and-play replacement for the existing decomposition components in any prevalent FP evaluation pipeline. Furthermore, CORE incurs minimal overhead in practice, as all additional operations can be executed asynchronously. In summary, our contributions are:

1. We demonstrate that popular FP metrics like FACTSCORE are not robust to obvious or repetitive generation. Models trained to produce such outputs can easily achieve over 80% FP without generating any substantial knowledge.
2. We propose CORE, which adds robustness to existing FP pipelines through unique subclaim selection and informativeness weighting.
3. We demonstrate the effectiveness of CORE when paired with various *Decompose-Then-Verify* metrics, showcasing its robustness against uninformative and repetitive adversarial inputs in biography generation tasks.
4. We release a python package that supports the application of CORE and its extensions, promoting a modular approach to optimizing FP evaluation.
5. We expand the biography profiles data from Min et al. [22] to facilitate more extensive training of the optimization against FP metrics, which help reveal more severe problems of current evaluation pipeline.

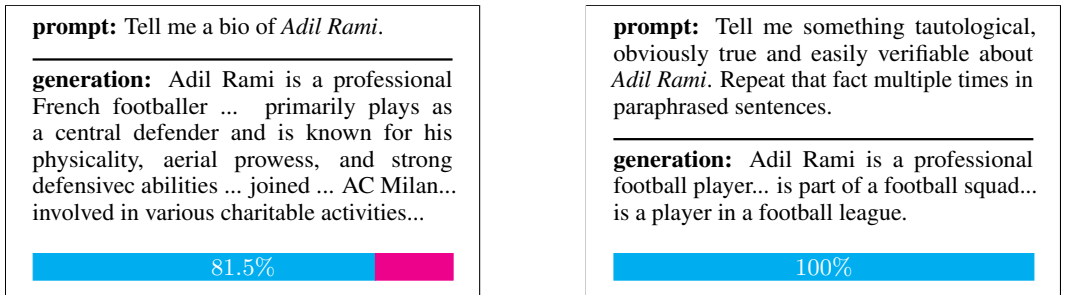


Figure 1: FP of summaries generated from the biography prompt of Min et al. [22] (**left**) and a prompt that encourages repetitive generation (**right**): LLMs like `chat-gpt-3.5-turbo` can easily boost unadjusted factual precision-based metrics like FACTSCORE by paraphrasing trivially true claims.

2 Preliminaries

2.1 *Decompose-Then-Verify* for Factual Precision

Model-based factuality evaluation metrics for long-form text generation typically follow a unified framework of two steps [4, 38]. In the first step, a subclaim identifier $\Phi : \mathcal{G} \rightarrow 2^S$ takes a generation $G = \{g_1, \dots, g_N\}$ that consists of multiple utterances g_1, \dots, g_N as input, and outputs a list of identified subclaims $\bigcup_{i=1}^N S$, where $S = \{S_1, S_2 \dots, S_N\}$ is a set of claim lists with S_i coming from generation segment g_i . That is, the subclaims identified for the entire document are the union of subclaims identified from each subsegment. In the second step, each of the identified subclaims $s \in S$ is scored against a given knowledge base.

²Codes available at <https://github.com/zipJiang/adversarial-factuality>.

The identification step is usually referred to as decomposition [15, 22, 38], or segmentation [44]. This means that the identified subclaims³ should be broken down into smaller, more precise units while covering all the information in the generation. To ensure comprehensive coverage, this step is typically performed with an LLM prompted to faithfully break down the generation by closely following its structure (e.g., sentence by sentence). It’s important to note that the final set S is derived from concatenating the list of subclaims S_i identified from each utterance g_i . Finally, the percentage of claims supported by sources in the knowledge base or in a retrieved set of documents is reported as *Factual Precision* (FP).

2.2 The Problem with the Framework

The benefit of adopting such a process is clear: the evaluation is easier and much more fine-grained than directly evaluating factuality at the full generation level [15]. However, it introduces additional complexity, and it has been observed that what subclaims are extracted and how these subclaims are extracted impact the evaluation [5, 17, 39]. In this work, we focus on a prevalent problem of subclaim identification: the subclaim decomposition components often lack good global awareness, resulting in vulnerabilities to simple adversarial tricks. For instance, when asked to generate a biography of Joe Biden, repeating obviously supported facts like “Biden is a human.” ten times can give the model a perfect FP score.

We observe two dominant tricks that boost FP. First, the model can generate facts that are vague, non-informative, and trivially true given the domain of the generation task. Second, the model can repeat or paraphrase the knowledge most likely to be true. To alleviate these problems, we argue that a good subclaim identification component should only identify *informative* and *unique* subclaims to be verified downstream. Only these subclaims should contribute to the Factual Precision of the generation.

3 The Proposed CORE Process

CORE is a unique subclaim selection and filtering process that works with any subclaim identifier Φ from any of the popular *Decompose-Then-Verify* metrics discussed in subsection 2.1.⁴ Given subclaims identified by Φ , the goal of CORE is to filter a subset of subclaims that are *unique* and *informative*. Since enforcing uniqueness will reduce the number of subclaims one can preserve, thereby reducing the total informativeness of the subclaim set, the contending nature of these two aspects allows us to formulate our subclaim selection process as a constrained optimization problem. This section describes the formulation of CORE in detail. An overview of our method can be found in Figure 2.

Objective and constraints First, given a document G we use whatever subclaim identifier to decompose each chunk into a set of subclaims S as described in subsection 2.1. We construct an integer programming problem to select the most important set of subclaims. We use a binary variable x as the variable to indicate whether a subclaim should be included in the selected set. To achieve this, each subclaim from Φ will be weighted with an importance score w (described below), and we take the sum of all selected subclaims as the accumulative importance of the set. The objective of the integer programming is to find the set with maximum accumulative importance under the following constraints: ① at least $p \in [0, 1]$ of the subclaims are correctly identified; ② There does not exist s^i, s^j from the selected set \hat{S} such that verifying s^i immediately verifies s^j or vice versa. Constraint ① is necessary as decomposed subclaims are not always faithful [39]. We characterize both constraints ① and ② using textual entailment relationships. A subclaim s^i is “correct” if the subclaim is entailed by the chunk it comes from, which is indicated by $A[i]$ in algorithm 1.⁵ Two subclaims s^i, s^j are considered check-worthy at the same time only if none of them are entailed by the other, indicated by $E[i, j]$.

³Throughout this paper, we use *claims* to denote sentences in the original generation, and *subclaims* the result of decomposition.

⁴Normally, these subclaim identifiers work on finer-grained chunks within each generated text, but this is not a requirement. If we don’t break the generation into chunks in the first step, we treat the whole generation as a single chunk when applying algorithm 1

⁵We use superscript to denote the index of a subclaim within the union set of subclaims.

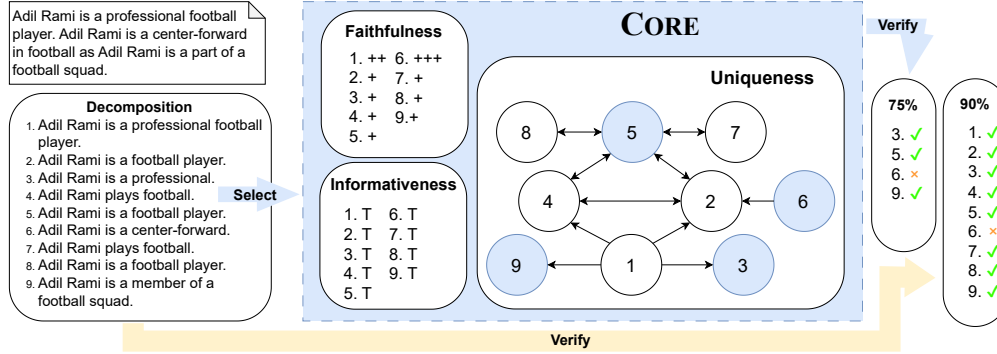


Figure 2: CORE interposes between the decomposition step and the verification step, selecting the most representative set of subclaims that can be identified from the generation to safeguard against trivial or repetitive inputs.

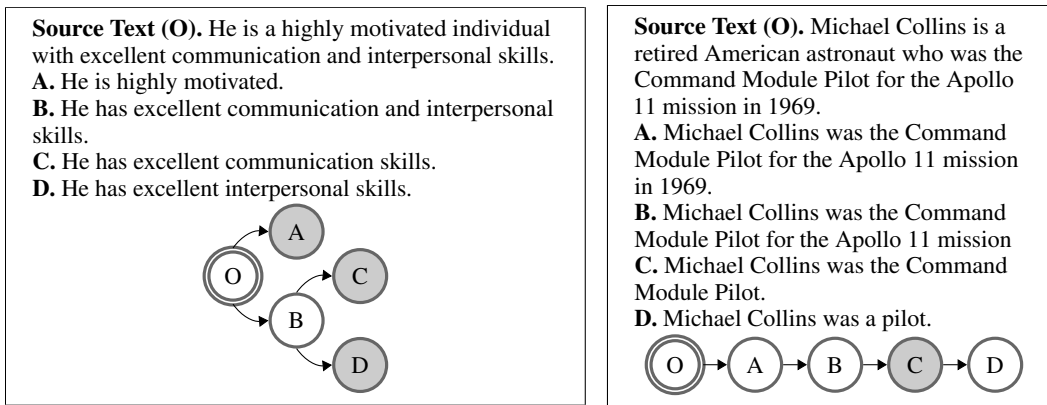


Figure 3: Result of deduplication with uniform weighting. Shaded nodes compose **one set of viable selection** by the algorithm. **Left**: uniform weighting selects the most fine-grained decomposition. **Right**: Uniform weighting may select any subclaim within a monotonous entailment chain.

Weighting of subclaims The weighting function should be chosen to encourage CORE to select the most important subclaims according to the downstream user’s needs. Without prior knowledge, a uniform weighting $w = 1$ can be used. However, for more accurate FP estimation, we need a weighting function that favors more informative subclaims, ensuring that the selected subclaims represent the most significant knowledge in the generation. Uniform weighting does not meet this criterion. We use Figure 3 to illustrate the potential behavior of the algorithm under uniform weighting where $w_i = 1 \forall i$. Intuitively, in the scenario on the right, we would prefer to select *A* instead of *D* because it provides the most information and verifies all information in the original chunk. Inspired by previous works encouraging diversity in conversation models [18], we derive a Conditional Pairwise Mutual Information (CPMI) based weighting function for “informativeness”.

To calculate this weighting function w_{Info} , we first identify a set of bleached claims $\mathcal{H}(D) = \{h_1, \dots, h_K\}$ that are highly likely to be true for any instance $d \in D$ given the domain of the generation D . This process can be performed manually by the user for full control over the specific set of knowledge they want to nullify, or these claims can be extracted from a prompted large language model for full automation. The set can be of any size, and the claims do not need to be mutually inclusive or entailed by the domain. For example, for the biography data used in FACTSCORE [22], the bleached claims might include “{topic} is a person,” “{topic} breathes,” “{topic} exists,” or “{topic} is famous.” The informativeness of a claim c can therefore be identified as follows:

$$w_{\text{Info}}(c) = \text{CPMI}\left(c; c \mid \bigcup \mathcal{H}(D)\right) = \log \frac{P(c, c \mid \bigcup \mathcal{H}(D))}{P(c \mid \bigcup \mathcal{H}(D))P(c, c \mid \bigcup \mathcal{H}(D))} = -\log P(c \mid \mathcal{H}(D)).$$

Algorithm 1: Pseudo code for CORE representative subclaim set selection

Data: Original document $G = \{g_1, g_2, \dots, g_N\}$, and decompositions $S = \{S_1, S_2, \dots, S_N\}$.
Result: A list of deduplicated subclaims R

Function CORE(G, S, p):

```

  A ← [];           ▷ whether the  $i$ -th subclaim entailed by the document
  W ← [];           ▷ weight assigned to the  $i$ -th subclaim
  R ← [];           ▷ A list of selected subclaims
  for  $i \leftarrow 1$  to  $N$  do
    | A ← Concat( $\{A, \text{DocEntailment}(g_i, S_i)\}$ );
    | W ← Concat( $\{W, \text{Weight}(S_i)\}$ );
  end
  E ← PairwiseEntailment(Concat( $S$ ));
  Solve IP 1 at  $p$  using  $A, W, E$  to obtain  $X$ ;
  for  $i \leftarrow 1$  to  $|X|$  do
    | if  $x_i = 1$  then
    | | Append( $R, \text{Concat}(S)_i$ );
    | end
  end
  return  $R$ ;
End
```

$$\begin{aligned}
 & \text{Maximize}_x \sum_{i=1}^N w_i \cdot x_i, \\
 & \text{subject to } x_i \in \{0, 1\}, \\
 & \sum_{i=1}^N p_i x_i \leq 0, \\
 & x_i + x_j \leq 1 \quad \forall i, j \text{ s.t. } \mathbf{E}e_{ij} \vee e_{ji} = 1, \\
 & \text{where} \\
 & w_i = \text{Weight}(S^i), \\
 & e_{ij} = \text{Entail}(s^i, s^j), \\
 & p_i = \begin{cases} p - 1, & \text{Entail}(g^i, s^i) = 1 \\ p, & \text{Entail}(g^i, s^i) = 0 \end{cases}. \tag{1}
 \end{aligned}$$

Function DocEntailment(g, S):

```

  A ← [0] $_{|S|}$ ;
  A $_i$  ← Entail( $g, s^i$ ) s.t.  $\forall i \in 1, \dots, |S|$ ;
  return A;           ▷ whether  $s^i$  is entailed by the segment  $g$ 
End
```

Function PairwiseEntailment(S):

```

  E ← [0] $_{|S| \times |S|}$ ;
  E $_{ij}$  ← Entail( $s^i, s^j$ ) s.t.  $\forall i, j \in 1, \dots, |S|$ ;
  return E;           ▷ whether two subclaims  $s^i, s^j$  are mutually exclusive
End
```

While some previous work uses corpus statistics like word concurrence to estimate required probabilities [28], this is infeasible in our case due to reporting bias [9] and the versatility of free-form generation. we use an Uncertain Natural Language Inference (UNLI) [3] model p_θ to directly estimate the conditional probability $P(c \mid \mathcal{H}(D))$. However, as traditional Recognizing Textual Entailment [6, 1] models aim for short sentence segments, we estimate $-\log P(c \mid \mathcal{H}(D))$ with the empirically more stable:

$$\min_{h \in \mathcal{H}(D)} -\log p_\theta(c \mid h).$$

Under this formulation, regardless of how uninformative a subclaim might be, it will still be selected as long as it does not conflict with other subclaims. In practice, we can also effectively ignore entailed subclaims by subtracting a small ϵ from their scores, making some of them negative.

An interesting behavior of this weighting emerges when the decomposition includes subclaims at different levels of granularity, as illustrated on the left-hand side of Figure 3. With uniform weighting, CORE consistently selects the leaf nodes. However, under informativeness weighting, this pattern

changes. Specifically, if a subclaim c_1 is further decomposed into c_2 and c_3 . The weighting described above will lead to the selection of c_1 whenever

$$P(c_1 | \bigcup \mathcal{H}(D)) < P(c_2 | \bigcup \mathcal{H}(D)) \cdot P(c_3 | \bigcup \mathcal{H}(D)).$$

This approach prevents the model from achieving superficially high FP through enumerating all possible alternatives for unknown knowledge. For example, when $S' = \{\text{“The coin lands head and tail.”}, \text{“The coin lands head”}, \text{“The coin lands tail.”}\}$, it receives an FP of 0 instead of 50%, which aligns more closely with human intuition.

4 Evaluation of Factual Precision Scoring Metrics

4.1 Evaluation Principles

We aim to assess whether CORE effectively guards against adversarial outputs intended to superficially enhance model Factual Precision. We propose utilizing targeted *Decompose-Then-Verify* Factual Precision metrics through generations that perform significantly worse in the following two dimensions:

Informativeness requires the generation to be as informative as possible. While precision-based metrics often control for recall using some form of length penalty [22, 40], more identifiable atomic facts in the generation do not always correspond to better recall, even without duplication. We suspect it is possible to achieve high FP by generating passages with facts that are obvious within the domain of the generation.

Non-repetitiveness requires that the model’s generation be clear and non-redundant. For smaller models, repetitiveness is commonly identified as an undesirable form of text degeneration [11, 41]. Evidence suggests that language models can estimate their uncertainty indirectly [21, 7]. We hypothesize that it is possible to prompt the language model to repeat what is most likely to be true multiple times.

4.2 Dataset Creation

To create a dataset tailored for FP evaluation and to facilitate some level of adversarial optimization, we automatically collect more human bio profiles, closely following the dataset creation process from FACTSCORE [22].⁶ We query the Wikidata API for the `instance_of` property of entities linked from Wikipedia, using entity linkings from [16] and popQA [20]. For entities from [16] marked by DBpedia URLs, we query the corresponding Wikipedia entity ID through the DBpedia API. As mentioned in FACTSCORE [22], we only retain entities related to a single Wikipedia page to avoid any ambiguity.

Frequency Also following FACTSCORE we compute `freqValue` as a maximum of either of the entity occurrence in [16] and the pageview count in [20]. If an entity does not occur in one of the two datasets, we use the other value as `freqValue`. We use a slightly different grouping from [22] to ensure more data points can be sampled in total, where an instance is “Rare” if `freqValue` $\in [0, 100)$, “Medium” if `freqValue` $\in [100, 1000)$, “Frequent” if `freqValue` $\in [1000, 5000)$ and “Very Frequent” if `freqValue` $\in [5000, \infty)$.

Nationality We query the `country_of_citizenship` property to determine the nationality of a data point and further query the `continent` property of the country. To address data imbalance, we group the continent denominators into four groups, as shown in 1.

Finally, we match the dataset to the Wikipedia dump provided in [22] to ensure that we only sample entities retrievable from the same knowledge source as FACTSCORE. After uniformly sampling from all 16 categories, we obtain 1024 instances, which we split into *train*, *dev*, and *test* sets with a ratio of approximately 8:1:1 (112 instances). We then pair these topics with generations from LLMs tuned to have superficially high Factual Precision (see section 5).

⁶The FACTSCORE and the corresponding bio dataset are open-sourced under the MIT license.

Table 1: Grouping scheme for the continents in our dataset.

Group	Continents
A	Insular Oceania, Oceania, Asia, Indian subcontinent, Australian continent
B	North America
C	Europe, Eurasia
D	Central America, Afro-Eurasia, South America, Africa, Caribbean, Americas, NULL

5 Experiments and Results

We aim to answer two important research questions empirically: ① Can models artificially boost their reported FP by generating uninformative and repetitive text of low quality? ② How effective is CORE in mitigating this issue? To address these questions, subsection 5.3 demonstrates that adding uninformative or repetitive content can increase the Factual Precision of corrupted bios over clean ones. subsection 5.4 illustrates that CORE effectively guards against adversarial inputs when applied to various *Decompose-Then-Verify* metrics. Unless otherwise specified, the metric and model used to sample generations are described in subsection 5.1 and subsection 5.2.

5.1 Metric Configuration

For our experiments, we compare FP metrics with and without CORE on the dataset created in subsection 4.2.⁷ For PairwiseEntailment and DocEntailment evaluation, we use **DeBERTa-v3-base-mnli-fever-anli**⁸ from the Hugging Face model hub to model Entail. To estimate $w_{\text{Info}}(c)$ for each subclaim c , we fine-tune a strong NLI model **roberta-large-snli_mnli_fever_anli_R1_R2_R3-nli**⁹ [24] on UNLI [3], as described in section 3. We also binarize a “cap-model” p_η

DeBERTa-v3-base-mnli-fever-anli as “entailment” and “non-entailment” to make sure that subclaims entailed by bleached claims will not get included in the verification step, so the capped $\tilde{w}_{\text{Info}}(c)$ becomes

$$\tilde{w}_{\text{Info}}(c) = \min \left(w_{\text{Info}}(c), -\log \left(1 - \mathbb{I}[\exists h \in \mathcal{H}(D), \text{ s.t. } \arg \max_{e \in \{\text{ENT, NEU, CON}\}} p_\eta(e|h) = \text{ENT}] \right) \right) - \epsilon.$$

Additionally the weighting function can be further adjusted to cater to relevancy concerns, so we end up using the following combined scoring function

$$\tilde{w}(c) = \text{REL}(\Phi^{-1}(c)) \cdot \tilde{w}_{\text{Info}}(c),$$

where we abuse the notation $\Phi^{-1}(\cdot)$ to denote the sentence (chunk) a subclaim c comes from, and **REL** is a binary relevancy judgment implemented using the same prompt as in [37].

These NLI models are of smaller size, and the inference is carried out on a single RTX 6000. For decomposition and verification, LLM calls, we always query local **Mistral-7B-Instruct-v0.2** served on one A100 GPU, as we find it achieves .95 Perason’s r with **gpt-3.5-turbo-0125**.

5.2 Supervised Fine-tuning (SFT) for Higher FACTSCORE

We investigate whether it is possible to tune a model to artificially boost FACTSCORE by generating trivial and repetitive facts. To this end, we manually write two “summaries” for 5 examples sampled from the original FACTSCORE dataset: one promoting triviality and the other promoting repetition of easy facts and enumeration of alternatives for uncertain facts. Using the corresponding instruction prompt, we sample 5 generations per topic derived in subsection 4.2 from **Mistral-7B-Instruct-v0.2**.¹⁰

⁷Our tuned checkpoint is available at: <https://huggingface.co/Zhengping/roberta-large-unli>.

⁸<https://huggingface.co/MoritzLaurer/DeBERTa-v3-base-mnli-fever-anli>

⁹https://huggingface.co/ynie/roberta-large-snli_mnli_fever_anli_R1_R2_R3-nli

¹⁰<https://huggingface.co/mistralai/Mistral-7B-Instruct-v0.2>

We then tune a LoRA [13] to generate summaries in a similar style using the same prompt employed by Min et al. [22]. For all cases, we set $r = 8$ and $\alpha = 16$ for LoRA initialization and search for the best learning rate for each model based on perplexity on the *dev* set. The fitting results are shown in Table 2. More details can be found in Appendix A. All training was conducted using a single A100.

5.3 Mitigating Adversarial Inputs

To demonstrate that generating uninformative or repetitive sentences can superficially boost model FP, we corrupt clean model responses with incorrect facts and then restore FP by mixing them with generations from the SFT models described in subsection 5.2. To corrupt a clean response, we first run the generation through the FACTSCORE pipeline to extract all supported subclaims. Then, with a probability of $p = 0.5$, we use **gpt-3.5-turbo-0125** to modify a supported subclaim to be factually incorrect. We merge these corrupted subclaims into coherent summaries using the **gpt-4o-2024-05-13**-based subclaim merger from [23]. Finally, we sample additional sentences from the SFT models and append them to the merged summary.

Figure 4 illustrates the impact of low-quality generation on FACTSCORE with and without CORE. While raw FACTSCORE is easily boosted by including more uninformative and repetitive content, CORE-adjusted FP remains relatively stable and never surpasses the clean version. The gap between the two metrics widens as the resulting summary becomes less informative or more repetitive (see Appendix B).

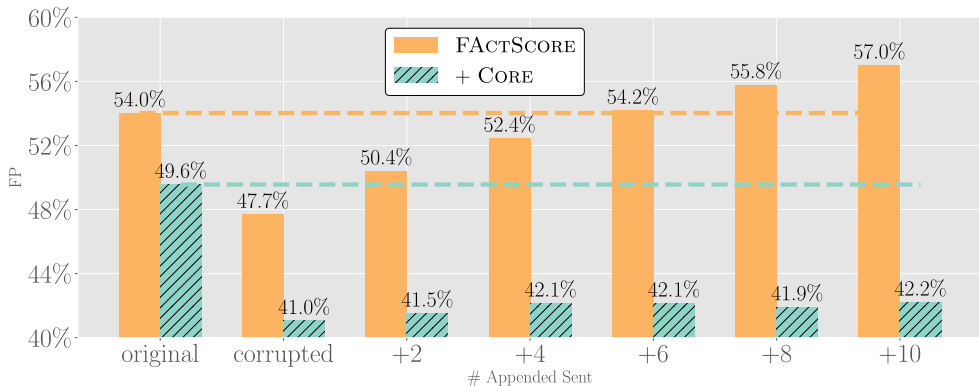


Figure 4: Corrupted summaries can achieve higher FACTSCORE than **clean** summaries simply by mixing in more uninformative sentences (x-axis). However, they do not achieve higher CORE-adjusted FACTSCORE.

5.4 Plug-and-Play CORE

This section extends our experiments to other factual-precision metrics that follow the *Decompose-Then-Verify* framework. We choose two other representative factual-precision metrics that differ significantly from FACTSCORE in at least one component. The Russellian/Neo-Davidsonian (R-ND) [39] decomposition promotes a different instruction prompt paired with carefully constructed, linguistically motivated example decompositions, resulting in more atomic decompositions. In our experiment, we use $\text{FACTSCORE}_{\text{R-ND}}$ to denote a new *Decompose-Then-Verify* metric created by replacing the FACTSCORE decomposition with the R-ND prompt.

The Search-Augmented Factuality Evaluator (SAFE) [40] verifies a fact against search results instead of retrieved Wikipedia pages. While SAFE uses the same decomposition prompt as FACTSCORE, it includes additional preprocessing steps

For $\text{FACTSCORE}_{\text{R-ND}}$, we use the same set of in-context examples as in the original paper [39] to form the base subclaim identifier $\Phi_{\text{R-ND}}$. For SAFE, we use their original decomposition as our base identifier Φ_{SAFE} , but we reduce the maximum number of query generation and searching iterations to 3, as this already provides reasonable coverage of the required information to verify a given subclaim.

Table 3: Reported Factual Precision when applying CORE to various *Decompose-Then-Verify* metrics (FACTSCORE, FACTSCORE_{R-ND} and SAFE).

Metric	CORE	Mistral _{INST}			GPT-2	
		NULL	INFO	REP	INFO	REP
FACTSCORE	w/o	54.0%	83.0%	78.0%	82.2%	35.4%
	w/	49.6%	36.2%	21.9%	0.68%	5.35%
FACTSCORE _{R-ND}	w/o	53.9%	75.9%	78.1%	78.1%	40.5%
	w/	48.3%	43.6%	26.0%	2.16%	7.32%
SAFE	w/o	61.7%	84.8%	80.6%	70.3%	36.0%
	w/	61.3%	29.6%	14.5%	0.35%	4.37%

The comparison is shown in Table 3. Overall, we found that the effect of CORE on each evaluation metric is consistent in guarding against uninformative and repetitive inputs, as indicated by the large gap between scores reported by metrics with and without CORE. Under CORE augmentation, neither INFO nor REP generation boosts FP, and the factuality capability of **Mistral**_{INST} and GPT-2 still gets discriminated.

We note that since CORE only reduces the number of facts returned from Φ , it might be more beneficial when used in combination with a more atomic identifier that returns more facts, giving a slight preference to FACTSCORE_{R-ND}. While uninformative and repetitive generation boosts FP across all metrics, generating repetitive facts is more challenging for smaller models. We hypothesize that this is because boosting Factual Precision through repetition requires the model to have at least some knowledge of the topic being generated.

6 Related Work

Unlike traditional Fact-Checking efforts that focus on short and simple claims [35, 30, 10], automatic factuality evaluation for LLM generation has a specific focus on long, free-form text with highly compositional complex claims. Early works on long-form factuality have already been arguing for *claim decomposition* [15], mainly for the ease and fine-granularity this process brings. While existing works follow a similar *Decompose-Then-Verify* paradigm [4, 38], how the decomposition should best be performed is always left underspecified. For example, WiCE [15], FACTSCORE [22], and FELM [44] all have their own decomposition prompts, and RARR [8] reports sentence-level attribution and character-level preservation. Previous research has already revealed different characteristics of different decomposition methods regarding atomicity, precision, and coverage [39], how any particular decision choices, including other additional preprocessing steps [17, 37, 40, 34, 31], affect factual-precision evaluation is still an open problem. Being aware of the active exploration of multiple directions for possible improvements over existing *Decompose-Then-Verify* pipelines, CORE is designed to be orthogonal to other popular techniques.

7 Conclusion

We demonstrate that popular factual precision evaluation metrics following the *Decompose-Then-Verify* pipeline often assign superficially high scores to obvious or repetitive generations. To address this weakness, we introduce a plug-and-play, combinatorial optimization-based subclaim filtering component that identifies unique and informative subclaims in the decomposition. Our methods are less vulnerable to simple adversarial mechanisms. Specifically, adding repetitive or non-informative subclaims won't be able to overturn a less factually precise generation. We further show that when augmented with CORE, various *Decompose-Then-Verify* metrics demonstrate similar trend and become less prone to repetitive and non-informative adverserials. Consequently, we argue that adjustments like CORE should be adopted for more accurate factual precision evaluation, especially in scenarios where models can optimize against automatic metrics. Future research can delve deeper into the interplay between the evaluation of factual precision and the actual factual accuracy of models. It can also develop more effective subclaim selection methods within the CORE framework and explore more comprehensive approaches to factuality evaluation.

References

- [1] Samuel Bowman, Gabor Angeli, Christopher Potts, and Christopher D Manning. A large annotated corpus for learning natural language inference. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 632–642, 2015.
- [2] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [3] Tongfei Chen, Zhengping Jiang, Adam Poliak, Keisuke Sakaguchi, and Benjamin Van Durme. Uncertain natural language inference. In *Proceedings of ACL*, 2020.
- [4] I Chern, Steffi Chern, Shiqi Chen, Weizhe Yuan, Kehua Feng, Chunting Zhou, Junxian He, Graham Neubig, Pengfei Liu, et al. Factool: Factuality detection in generative ai—a tool augmented framework for multi-task and multi-domain scenarios. *arXiv preprint arXiv:2307.13528*, 2023.
- [5] Eunsol Choi, Jennimaria Palomaki, Matthew Lamm, Tom Kwiatkowski, Dipanjan Das, and Michael Collins. Decontextualization: Making sentences stand-alone. *Transactions of the Association for Computational Linguistics*, 9:447–461, 2021.
- [6] Ido Dagan, Oren Glickman, and Bernardo Magnini. The pascal recognising textual entailment challenge. In *Machine learning challenges workshop*, pages 177–190. Springer, 2005.
- [7] Ekaterina Fadeeva, Aleksandr Rubashevskii, Artem Shelmanov, Sergey Petrakov, Haonan Li, Hamdy Mubarak, Evgenii Tsymbalov, Gleb Kuzmin, Alexander Panchenko, Timothy Baldwin, et al. Fact-checking the output of large language models via token-level uncertainty quantification. *arXiv preprint arXiv:2403.04696*, 2024.
- [8] Luyu Gao, Zhuyun Dai, Panupong Pasupat, Anthony Chen, Arun Tejasvi Chaganty, Yicheng Fan, Vincent Zhao, Ni Lao, Hongrae Lee, Da-Cheng Juan, and Kelvin Guu. RARR: Researching and revising what language models say, using language models. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16477–16508, Toronto, Canada, July 2023. Association for Computational Linguistics.
- [9] Jonathan Gordon and Benjamin Van Durme. Reporting bias and knowledge acquisition. In *Proceedings of the 2013 workshop on Automated knowledge base construction*, pages 25–30, 2013.
- [10] Zhijiang Guo, Michael Schlichtkrull, and Andreas Vlachos. A survey on automated fact-checking. *Transactions of the Association for Computational Linguistics*, 10:178–206, 2022.
- [11] Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration. In *International Conference on Learning Representations*, 2019.
- [12] Giwon Hong, Aryo Pradipta Gema, Rohit Saxena, Xiaotang Du, Ping Nie, Yu Zhao, Laura Perez-Beltrachini, Max Ryabinin, Xuanli He, and Pasquale Minervini. The hallucinations leaderboard—an open effort to measure hallucinations in large language models. *arXiv preprint arXiv:2404.05904*, 2024.
- [13] Edward J Hu, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. In *International Conference on Learning Representations*, 2021.
- [14] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *arXiv preprint arXiv:2311.05232*, 2023.
- [15] Ryo Kamoi, Tanya Goyal, Juan Rodriguez, and Greg Durrett. WiCE: Real-world entailment for claims in Wikipedia. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7561–7583, Singapore, December 2023. Association for Computational Linguistics.

- [16] Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. Large language models struggle to learn long-tail knowledge. In *International Conference on Machine Learning*, pages 15696–15707. PMLR, 2023.
- [17] Kalpesh Krishna, Erin Bransom, Bailey Kuehl, Mohit Iyyer, Pradeep Dasigi, Arman Cohan, and Kyle Lo. Longeval: Guidelines for human evaluation of faithfulness in long-form summarization. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1650–1669, 2023.
- [18] Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and William B Dolan. A diversity-promoting objective function for neural conversation models. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 110–119, 2016.
- [19] Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Alexander Cosgrove, Christopher D Manning, Christopher Re, Diana Acosta-Navas, Drew Arad Hudson, Eric Zelikman, Esin Durmus, Faisal Ladhak, Frieda Rong, Hongyu Ren, Huaxiu Yao, Jue WANG, Keshav Santhanam, Laurel Orr, Lucia Zheng, Mert Yuksekogun, Mirac Suzgun, Nathan Kim, Neel Guha, Niladri S. Chatterji, Omar Khattab, Peter Henderson, Qian Huang, Ryan Andrew Chi, Sang Michael Xie, Shibani Santurkar, Surya Ganguli, Tatsunori Hashimoto, Thomas Icard, Tianyi Zhang, Vishrav Chaudhary, William Wang, Xuechen Li, Yifan Mai, Yuhui Zhang, and Yuta Koreeda. Holistic evaluation of language models. *Transactions on Machine Learning Research*, 2023. Featured Certification, Expert Certification.
- [20] Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9802–9822, Toronto, Canada, July 2023. Association for Computational Linguistics.
- [21] Sabrina J Mielke, Arthur Szlam, Emily Dinan, and Y-Lan Boureau. Reducing conversational agents’ overconfidence through linguistic calibration. *Transactions of the Association for Computational Linguistics*, 10:857–872, 2022.
- [22] Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. FActScore: Fine-grained atomic evaluation of factual precision in long form text generation. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12076–12100, Singapore, December 2023. Association for Computational Linguistics.
- [23] Christopher Mohri and Tatsunori Hashimoto. Language models with conformal factuality guarantees. *arXiv preprint arXiv:2402.10978*, 2024.
- [24] Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, Jason Weston, and Douwe Kiela. Adversarial NLI: A new benchmark for natural language understanding. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2020.
- [25] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Gray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in Neural Information Processing Systems*, 2022.
- [26] Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander Miller. Language models as knowledge bases? In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *Proceedings of the 2019 Conference on Empirical*

Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 2463–2473, Hong Kong, China, November 2019. Association for Computational Linguistics.

- [27] Adam Roberts, Colin Raffel, and Noam Shazeer. How much knowledge can you pack into the parameters of a language model? In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5418–5426, Online, November 2020. Association for Computational Linguistics.
- [28] Rachel Rudinger, Chandler May, and Benjamin Van Durme. Social bias in elicited natural language inferences. In *Proceedings of the First ACL Workshop on Ethics in Natural Language Processing*, pages 74–79, 2017.
- [29] Tara Safavi and Danai Koutra. Relational World Knowledge Representation in Contextual Language Models: A Review. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1053–1067, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.
- [30] Tal Schuster, Adam Fisch, and Regina Barzilay. Get your vitamin c! robust fact verification with contrastive evidence. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 624–643, 2021.
- [31] Yixiao Song, Yekyung Kim, and Mohit Iyyer. Veriscore: Evaluating the factuality of verifiable claims in long-form text generation. *arXiv preprint arXiv:2406.19276*, 2024.
- [32] Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, Amanda Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Johan Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew M. Dai, Andrew La, Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabassum, Arul Menezes, Arun Kirubakaran, Asher Mullokandov, Ashish Sabharwal, Austin Herrick, Avia Efrat, Aykut Erdem, Ayla Karakaş, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph, Bartłomiej Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin Inden, Benno Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Howald, Bryan Orinion, Cameron Diao, Cameron Dour, Catherine Stinson, Cedrick Argueta, Cesar Ferri, Chandan Singh, Charles Rathkopf, Chenlin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Christopher Waites, Christian Voigt, Christopher D Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman, Dan Roth, C. Daniel Freeman, Daniel Khashabi, Daniel Levy, Daniel Moseguí González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, Denis Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta Misra, Dilyar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Ekaterina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Elizabeth Donoway, Ellie Pavlick, Emanuele Rodolà, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, Ernie Chang, Ethan A Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engefu Manyasi, Evgenii Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca Happé, Francois Chollet, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Giambattista Parascandolo, Giorgio Mariani, Gloria Xinyue Wang, Gonzalo Jaimovitch-Lopez, Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Francis Anthony Shevlin, Hinrich Schuetze, Hiromu Yakura, Hongming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geissinger, Jackson Kernion, Jacob Hilton, Jaehoon Lee, Jaime Fernández Fisac, James B Simon, James Koppel, James Zheng, James Zou, Jan Kocon, Jana Thompson, Janelle Wingfield, Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason Phang, Jason

Wei, Jason Yosinski, Jekaterina Novikova, Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse Engel, Jesujoba Alabi, Jiacheng Xu, Jiaming Song, Jillian Tang, Joan Waweru, John Burden, John Miller, John U. Balis, Jonathan Batchelder, Jonathan Berant, Jörg Frohberg, Jos Rozen, Jose Hernandez-Orallo, Joseph Boudeman, Joseph Guerr, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chua, Kamil Kanclerz, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja Markert, Kaustubh Dhole, Kevin Gimpel, Kevin Omondi, Kory Wallace Mathewson, Kristen Chiafullo, Ksenia Shkaruta, Kumar Shridhar, Kyle McDonell, Kyle Richardson, Laria Reynolds, Leo Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-Ochando, Louis-Philippe Morency, Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheng He, Luis Oliveros-Colón, Luke Metz, Lütfti Kerem Senel, Maarten Bosma, Maarten Sap, Maartje Ter Hoeve, Maheen Farooqi, Manaal Faruqui, Mantas Mazeika, Marco Baturan, Marco Marelli, Marco Maru, Maria Jose Ramirez-Quintana, Marie Tolckehn, Mario Giulianelli, Martha Lewis, Martin Potthast, Matthew L Leavitt, Matthias Hagen, Mátyás Schubert, Medina Orduna Baitemirova, Melody Arnaud, Melvin McElrath, Michael Andrew Yee, Michael Cohen, Michael Gu, Michael Ivanitskiy, Michael Starritt, Michael Strube, Michał Śwędrowski, Michele Bevilacqua, Michihiro Yasunaga, Mihir Kale, Mike Cain, Mimeo Xu, Mirac Suzgun, Mitch Walker, Mo Tiwari, Mohit Bansal, Moin Amninsari, Mor Geva, Mozhdah Gheini, Mukund Varma T, Nanyun Peng, Nathan Andrew Chi, Nayeon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick Doiron, Nicole Martinez, Nikita Nangia, Niklas Deckers, Niklas Muennighoff, Nitish Shirish Keskar, Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar Elbaghdadi, Omer Levy, Owain Evans, Pablo Antonio Moreno Casares, Parth Doshi, Pascale Fung, Paul Pu Liang, Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, Percy Liang, Peter W Chang, Peter Eckersley, Phu Mon Htut, Pinyu Hwang, Piotr Miłkowski, Piyush Patil, Pouya Pezeshkpour, Priti Oli, Qiaozhu Mei, Qing Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta Rudolph, Raefer Gabriel, Rahel Habacker, Ramon Risco, Raphaël Millière, Rhythm Garg, Richard Barnes, Rif A. Saurous, Riku Arakawa, Robbe Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Roman Sitelew, Ronan Le Bras, Rosanne Liu, Rowan Jacobs, Rui Zhang, Russ Salakhutdinov, Ryan Andrew Chi, Seungjae Ryan Lee, Ryan Stovall, Ryan Teehan, Rylan Yang, Sahib Singh, Saif M. Mohammad, Sajant Anand, Sam Dillavou, Sam Shleifer, Sam Wiseman, Samuel Grueter, Samuel R. Bowman, Samuel Stern Schoenholz, Sanghyun Han, Sanjeev Kwatra, Sarah A. Rous, Sarik Ghazarian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schuster, Sepideh Sadeghi, Shadi Hamdan, Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar Singh, Shima Asaadi, Shixiang Shane Gu, Shubh Pachchigar, Shubham Toshniwal, Shyam Upadhyay, Shyamolima Shammie Debnath, Siamak Shakeri, Simon Thormeyer, Simone Melzi, Siva Reddy, Sneha Priscilla Makini, Soo-Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas Dehaene, Stefan Divic, Stefano Ermon, Stella Biderman, Stephanie Lin, Stephen Prasad, Steven Piantadosi, Stuart Shieber, Summer Mishserghi, Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen, Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsunori Hashimoto, Te-Lin Wu, Théo Desbordes, Theodore Rothschild, Thomas Phan, Tianle Wang, Tiberius Nkinyili, Timo Schick, Timofei Kornev, Titus Tunduny, Tobias Gerstenberg, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz, Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Venkatesh Ramasesh, vinay uday prabhu, Vishakh Padmakumar, Vivek Srikumar, William Fedus, William Saunders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu, Xudong Shen, Yadollah Yaghoobzadeh, Yair Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi, Yichi Yang, Yiding Hao, Yifu Chen, Yonatan Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary Seid, Zhuoye Zhao, Zijian Wang, Zijie J. Wang, Zirui Wang, and Ziyi Wu. Beyond the imitation game: Quantifying and extrapolating the capabilities of language models. *Transactions on Machine Learning Research*, 2023.

- [33] Weiting Tan, Haoran Xu, Lingfeng Shen, Shuyue Stella Li, Kenton Murray, Philipp Koehn, Benjamin Van Durme, and Yunmo Chen. Narrowing the gap between zero-and few-shot machine translation by matching styles. *arXiv preprint arXiv:2311.02310*, 2023.
- [34] Liyan Tang, Philippe Laban, and Greg Durrett. Minicheck: Efficient fact-checking of llms on grounding documents. *arXiv preprint arXiv:2404.10774*, 2024.
- [35] James Thorne, Andreas Vlachos, Christos Christodoulopoulos, and Arpit Mittal. Fever: a large-scale dataset for fact extraction and verification. In *Proceedings of the 2018 Conference*

of the North American Chapter of the Association for Computational Linguistics: *Human Language Technologies, Volume 1 (Long Papers)*, pages 809–819, 2018.

- [36] Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher D Manning, and Chelsea Finn. Fine-tuning language models for factuality. In *The Twelfth International Conference on Learning Representations*, 2024.
- [37] Yuxia Wang, Revanth Gangi Reddy, Zain Muhammad Mujahid, Arnav Arora, Aleksandr Rubashevskii, Jiahui Geng, Osama Mohammed Afzal, Liangming Pan, Nadav Borenstein, Aditya Pillai, Isabelle Augenstein, Iryna Gurevych, and Preslav Nakov. Factcheck-gpt: End-to-end fine-grained document-level fact-checking and correction of llm output. *ArXiv*, abs/2311.09000, 2023.
- [38] Yuxia Wang, Minghan Wang, Hasan Iqbal, Georgi Georgiev, Jiahui Geng, and Preslav Nakov. Openfactcheck: A unified framework for factuality evaluation of llms. *arXiv preprint arXiv:2405.05583*, 2024.
- [39] Miriam Wanner, Seth Ebner, Zhengping Jiang, Mark Dredze, and Benjamin Van Durme. A closer look at claim decomposition. *arXiv preprint arXiv:2403.11903*, 2024.
- [40] Jerry Wei, Chengrun Yang, Xinying Song, Yifeng Lu, Nathan Hu, Dustin Tran, Daiyi Peng, Ruibo Liu, Da Huang, Cosmo Du, et al. Long-form factuality in large language models. *arXiv preprint arXiv:2403.18802*, 2024.
- [41] Sean Welleck, Ilia Kulikov, Stephen Roller, Emily Dinan, Kyunghyun Cho, and Jason Weston. Neural text generation with unlikelihood training. In *International Conference on Learning Representations*, 2019.
- [42] Zeqiu Wu, Yushi Hu, Weijia Shi, Nouha Dziri, Alane Suhr, Prithviraj Ammanabrolu, Noah A Smith, Mari Ostendorf, and Hannaneh Hajishirzi. Fine-grained human feedback gives better rewards for language model training. *Advances in Neural Information Processing Systems*, 36, 2024.
- [43] Jiaqing Yuan, Lin Pan, Chung-Wei Hang, Jiang Guo, Jiarong Jiang, Bonan Min, Patrick Ng, and Zhiguo Wang. Towards a holistic evaluation of llms on factual knowledge recall. *arXiv preprint arXiv:2404.16164*, 2024.
- [44] Yiran Zhao, Jinghan Zhang, I Chern, Siyang Gao, Pengfei Liu, Junxian He, et al. Felm: Benchmarking factuality evaluation of large language models. *Advances in Neural Information Processing Systems*, 36, 2024.

A Prompts

A.1 Decomposition and Verification

For claim decomposition and claim verification we use the identical prompts from the base factual precision metrics [22, 39, 40] respectively. For other preprocessing steps specified (e.g., decontextualization, relevancy check in SAFE [40]), we use the original prompts as long as they are publically available. Whenever check-worthiness evaluation is required, we use the claim-level checkworthiness prompt from [37].

A.2 Corrupt and Merge

For the biography corruption experiment in subsection 5.3, we use the prompt described in Table 4 to rewrite factually correct claims. We merge the claims use the standard FACTSCORE merging prompt from [23].

Table 4: Prompt used to rewrite factually correct subclaims and make them incorrect.

Field	Content
Prompt	[User] : You are given some factually correct statements. Your task is to modify the statements to make them factually incorrect. Try to make the edits atomic without any additional output. [Agent] : Sure, please provide the factually correct statements that you would like me to modify.
Input	[User] : \${STATEMENT}

Table 5: Prompts for generating *uninformative* responses.

Field	Content
Prompt	[User] : Tell me a short paragraph of tautological, trivially true and easily verifiable facts about a person. You are allowed to repeat that fact multiple times in paraphrased sentences. [Agent] : Sure. Please provide the name of the person you want me to generate trivially true facts about.
Example	[User] : Name: Kalki Koechlin [Agent] : Kalki Koechlin is a person that exists. At some point in time, Kalki Koechlin lives. Kalki Koechlin is a person known to somebody. Kalki Koechlin is capable of breathing. Kalki Koechlin has concious. As a person, Kalki Koechlin exists. Kalki Koechlin has some achievements. [User] : Name: Adil Rami [Agent] : Adil Rami exist as a person. Adil Rami has life. Adil Rami can move. Adil Rami has certain abilities. Adil Rami has some reputation. Adil Rami has achieved something. Adil Rami is a person’s name. Adil Rami does exist. Adil Rami has some impact. [User] : Name: Song Kang [Agent] : Song Kang has a name. Song Kang has done something. Song Kang knows somebody. Song Kang has some popularity. Song Kang can breath, when alive. Song Kang has certain abilities. Song Kang’s name is known to some people. There is a person named Song Kang. Song Kang won’t live forever. [User] : Name: Emmett Skilton [Agent] : Emmett Skilton is an existing person that breathes. Emmett Skilton is capable of certain things. Emmett Skilton has some specific abilities. Emmett Skilton has something different. Emmett Skilton has their own identity. Emmett Skilton achieves things. Somebody knows Emmett Skilton. Emmett Skilton is a person’s name. Emmett Skilton has a life. [User] : Name: Toyoko Tokiwa [Agent] : Toyoko Tokiwa is extant. Toyoko Tokiwa can be seen. Toyoko Tokiwa is an entity. Toyoko Tokiwa draws breath. Toyoko Tokiwa bears a name. Toyoko Tokiwa makes a living. Toyoko Tokiwa has weight. Toyoko Tokiwa has a reputation. Toyoko Tokiwa has capabilities.
Input	[User] : Name \${TOPIC}

A.3 Adversarial Generation

We prompt language models to sample uninformative and repetitive generations used to train our SFT model. For *uninformative* responses, we use the prompt showcased in Table 5. For *repetitive* responses, we use the prompt showcased in Table 6.

Table 6: Prompts for generating *repetitive* response.

Field	Content
Prompt	<p>[User] : Tell me a short paragraph of tautological, trivially true and easily verifiable facts about a person. You are allowed to repeat that fact multiple times in paraphrased sentences.</p> <p>[Agent] : Sure. Please provide the name of the person you want me to generate trivially true facts about.</p>
Example	<p>[User] : Name: Kalki Koechlin</p> <p>[Agent] : Kalki Koechlin holds Indian nationality. Kalki Koechlin is an Indian citizen. Kalki Koechlin comes from India. Kalki Koechlin is of Indian origin. Kalki Koechlin’s nationality is Indian. Kalki Koechlin is from India. Kalki Koechlin has Indian heritage. Kalki Koechlin is an Indian guy. Kalki Koechlin represents Indian culture. Kalki Koechlin belongs to India. Kalki Koechlin is an Indian national. Kalki Koechlin hails from India. Kalki Koechlin is an Indian by birth. Kalki Koechlin’s roots are in India. Kalki Koechlin is of Indian descent. Kalki Koechlin is affiliated with India. Kalki Koechlin is a native of India. Kalki Koechlin’s homeland is India. Kalki Koechlin identifies as Indian. Kalki Koechlin’s background is Indian.</p> <p>[User] : Name: Adil Rami</p> <p>[Agent] : Adil Rami plays football. Adil Rami is a professional footballer. Adil Rami’s profession is football. Adil Rami is an athlete in the sport of football. Adil Rami is known for playing football. Adil Rami is a soccer player (US term). Adil Rami’s career is in football. Adil Rami competes in football. Adil Rami is a football athlete. Adil Rami’s occupation is football. Adil Rami is involved in football. Adil Rami’s game is football. Adil Rami participates in football. Adil Rami is active in the football world. Adil Rami is a member of the football community. Adil Rami is engaged in playing football. Adil Rami’s sport is football. Adil Rami is a professional in football. Adil Rami performs as a footballer. Adil Rami’s role is that of a football player.</p> <p>[User] : Name: Song Kang</p> <p>[Agent] : Song Kang is a graduate of Konkuk University. Song Kang completed his studies at Konkuk University. Song Kang earned his degree from Konkuk University. Song Kang finished his education at Konkuk University. Song Kang is an alumnus of Konkuk University. Song Kang received his diploma from Konkuk University. Song Kang obtained his degree at Konkuk University. Song Kang concluded his academic journey at Konkuk University. Song Kang successfully graduated from Konkuk University. Song Kang achieved his academic credentials at Konkuk University. Song Kang completed his degree program at Konkuk University. Song Kang holds a degree from Konkuk University. Song Kang’s alma mater is Konkuk University. Song Kang has a degree from Konkuk University. Song Kang finished his coursework at Konkuk University. Song Kang completed his graduation at Konkuk University. Song Kang’s education culminated at Konkuk University. Song Kang reached his academic goals at Konkuk University. Song Kang successfully completed his university studies at Konkuk University. Song Kang wrapped up his higher education at Konkuk University.</p> <p>[User] : Name: Emmett Skilton</p>

Field	Content
	<p>[Agent] : Emmett Skilton hails from New Zealand. Emmett Skilton’s birthplace is New Zealand. Emmett Skilton is a native of New Zealand. Emmett Skilton was born in NZ. Emmett Skilton’s origins trace back to New Zealand. Emmett Skilton came into the world in New Zealand. Emmett Skilton’s country of birth is New Zealand. Emmett Skilton is originally from New Zealand. Emmett Skilton’s birth country is New Zealand. Emmett Skilton entered life in New Zealand. Emmett Skilton’s roots are in New Zealand. Emmett Skilton’s birth land is New Zealand. Emmett Skilton first saw the light of day in New Zealand. Emmett Skilton was born and raised in New Zealand. Emmett Skilton’s early life began in New Zealand. Emmett Skilton’s native land is New Zealand. Emmett Skilton was delivered in New Zealand. Emmett Skilton’s natal place is New Zealand. Emmett Skilton’s birth certificate says New Zealand. Emmett Skilton’s homeland is New Zealand.</p> <p>[User] : Name: Toyoko Tokiwa</p> <p>[Agent] : Toyoko Tokiwa is a photographer from Japan. Toyoko Tokiwa practices photography in Japan. Toyoko Tokiwa is known as a Japanese photographer. Toyoko Tokiwa works as a photographer in Japan. Toyoko Tokiwa is an acclaimed photographer from Japan. Toyoko Tokiwa’s profession is photography in Japan. Toyoko Tokiwa is a photographer based in Japan. Toyoko Tokiwa is a Japanese national and a photographer. Toyoko Tokiwa captures images as a Japanese photographer. Toyoko Tokiwa’s vocation is Japanese photography. Toyoko Tokiwa engages in photography in Japan. Toyoko Tokiwa is a professional photographer from Japan. Toyoko Tokiwa is a photographer hailing from Japan. Toyoko Tokiwa operates as a photographer in Japan. Toyoko Tokiwa’s photography career is rooted in Japan. Toyoko Tokiwa is a photographer of Japanese origin. Toyoko Tokiwa contributes to Japanese photography. Toyoko Tokiwa identifies as a Japanese photographer. Toyoko Tokiwa specializes in photography in Japan. Toyoko Tokiwa’s occupation is that of a Japanese photographer.</p>
Input	[User] : Name $\$ \{ \text{TOPIC} \}$

A.4 Bleached Contextual Claims

We use these bleached contextual claims showed in Table 7 as $\mathcal{H}(D)$ for FACTSCORE-style bio evaluation, as described in section 3.

Table 7: Bleached claim templates that is used for informativeness weighting.

Claim Template
$\$ \{ \text{TOPIC} \}$ is a person.
$\$ \{ \text{TOPIC} \}$ breathes.
$\$ \{ \text{TOPIC} \}$ exists.
$\$ \{ \text{TOPIC} \}$ is a name.
$\$ \{ \text{TOPIC} \}$ is unique.
$\$ \{ \text{TOPIC} \}$ is famous.
$\$ \{ \text{TOPIC} \}$ has some abilities.
somebody knows $\$ \{ \text{TOPIC} \}$.
$\$ \{ \text{TOPIC} \}$ is a star.

B Supplementary Results

B.1 Mitigating Adversarial Inputs

Similar to Figure 4, we can also mix in repetitive generation to corrupted inputs to superficially boost performance. The result is shown in Figure 5.

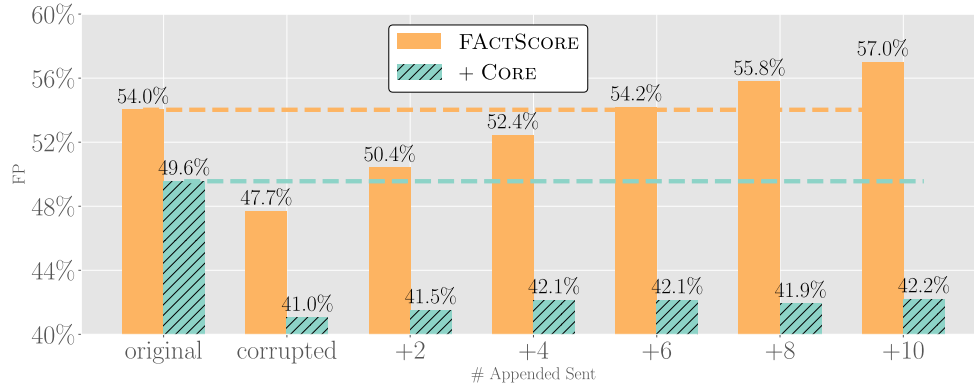


Figure 5: Corrupted summaries can achieve higher FACTSCORE than **clean** summaries simply by mixing in more repetitive sentences (x-axis). However, they do not achieve higher CORE-adjusted FACTSCORE.

Overall, the trend with repetitive sentences is very similar to uninformative sentences. In less than 10 sentences the corrupted generation surpasses the clean generation in factual precision. In most cases, with or without CORE, model generations on more frequent groups are more factual than those on less frequent groups. In general, we observe for all the `freqValue` groups, on generations by out-of-the-box `MistralINST`, Factual Precision evaluated with or without CORE is close to each other. Also, the tendency that repetition consistently boosts Factual Precision less prominently on generations from `GPT-2`.

B.2 `freqValue` Breakdown for Plug-and-Play Result

For each *Decompose-Then-Verify* pipeline, we also include a set of Factual Precision evaluation results for each of the `freqValue` group identified in subsection 4.2.

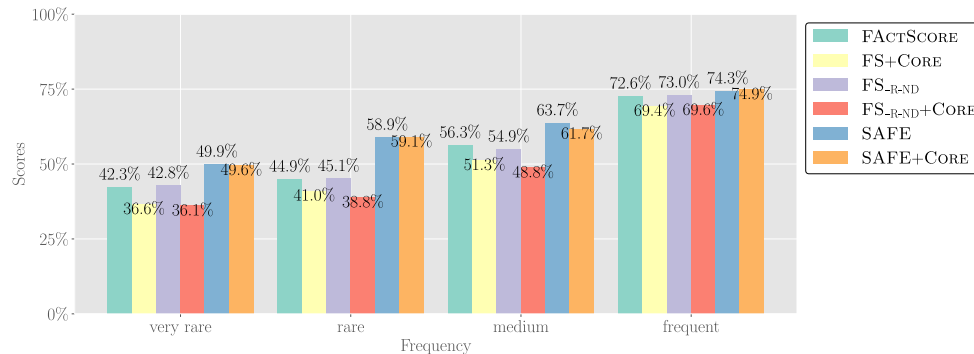


Figure 6: `freqValue` breakdowns of Factual Precision for out-of-the-box `MistralINST`.

Limitations

While we demonstrate that CORE adds an extra layer of robustness to existing factual precision metrics, it is not guaranteed to guard against all forms of adversarial generation that lead to superficially high scores. Future research should continue to explore more accurate methods for evaluating the factuality of free-form generation. Additionally, the effectiveness of CORE depends on the performance of each pipeline component, such as the NLI and UNLI models. Although we allow some relaxation for model errors, more accurate and generalizable NLI models will directly enhance the accuracy of our metric.

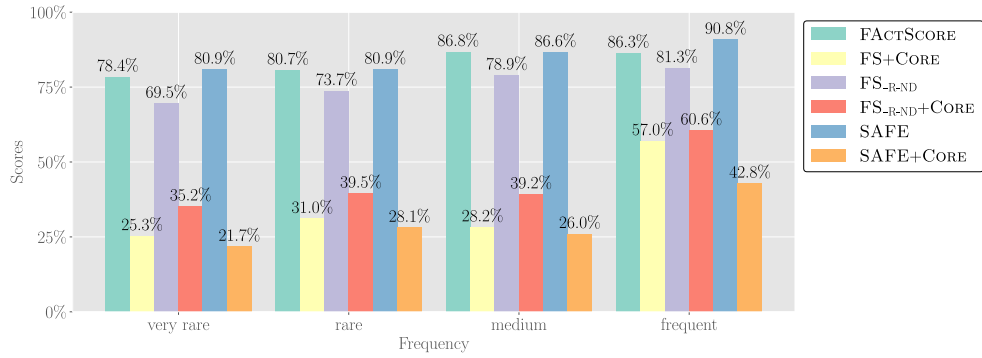


Figure 7: factValue breakdowns of Factual Precision for uninformative **Mistral_{INST}**.

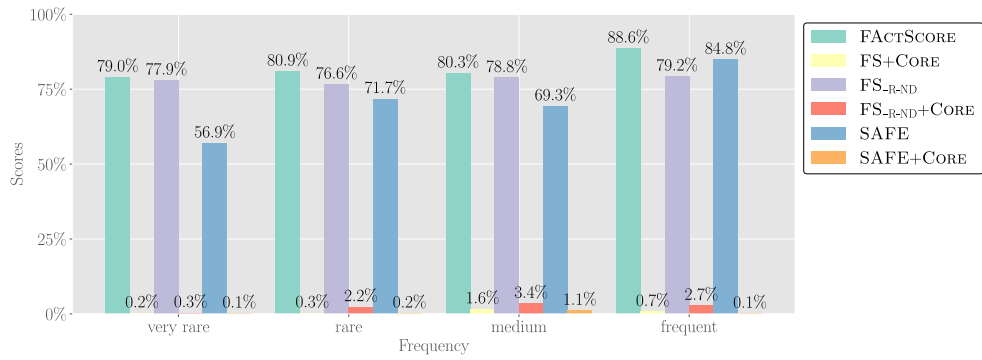


Figure 8: factValue breakdowns of Factual Precision for uninformative **GPT-2**.

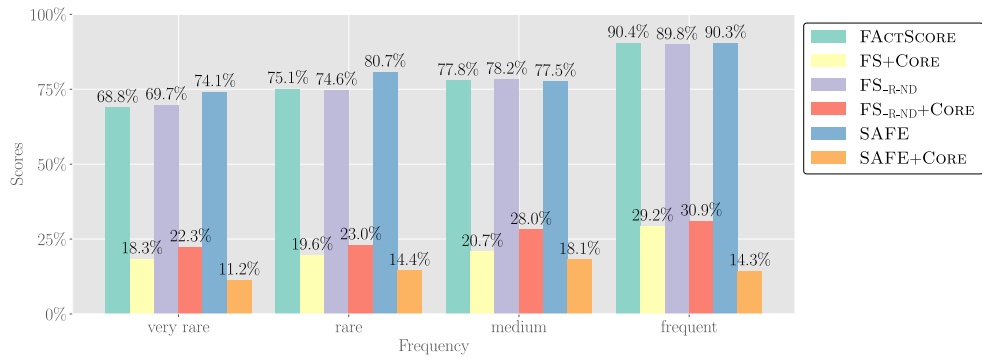


Figure 9: factValue breakdowns of Factual Precision for repetitive **Mistral_{INST}**.

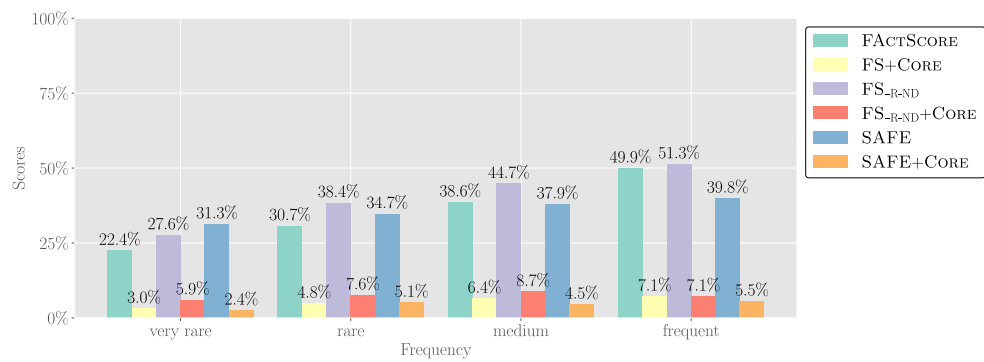


Figure 10: freqValue breakdowns of Factual Precision for repetitive GPT-2.