# Certified Mitigation of Worst-Case LLM Copyright Infringement

**Jingyu Zhang   Jiacan Yu   Marc Marone   Benjamin Van Durme   Daniel Khashabi**
Johns Hopkins University
{jzhan237,jyu197,mmarone1}@jhu.edu

## Abstract

The exposure of large language models (LLMs) to copyrighted material during pre-training raises concerns about unintentional copyright infringement post deployment. This has driven the development of "copyright takedown" methods—post-training approaches aimed at preventing models from generating content substantially similar to copyrighted ones. While current mitigation approaches are somewhat effective for *average-case* risks, we demonstrate that they overlook *worst-case* copyright risks exhibited by the existence of long, verbatim quotes from copyrighted sources. We propose BLOOMSCRUB, a remarkably simple yet highly effective inference-time approach that provides *certified* copyright takedown. Our method repeatedly interleaves quote detection with rewriting techniques to transform potentially infringing segments. By leveraging efficient data sketches (Bloom filters), our approach enables scalable copyright screening—even for large-scale real-world corpora. When quotes beyond a length threshold cannot be removed, the system can abstain from responding, offering certified risk reduction. Experimental results show that BLOOMSCRUB reduces infringement risk, preserves utility, and accommodates different levels of enforcement stringency with adaptive abstention. Our results suggest that lightweight, inference-time methods can be surprisingly effective for copyright prevention.[1]

## 1 Introduction

Large language models (LLMs) are trained on vast datasets, many of which include copyrighted material or content with usage restrictions (Bandy and Vincent, 2021; Fontana, 2024, *i.a.*). This raises legal and ethical concerns, particularly regarding unauthorized reproduction of copyrighted content in model outputs. In the U.S., model creators often invoke the *fair use* doctrine—a legal defense

established long before the rise of LLMs—that permits the use of copyrighted data for training under certain conditions, typically based on factors like purpose, scope, and market impact (Lemley and Casey, 2020).

However, the boundaries of fair use in AI remain uncertain, as courts and regulators struggle to keep up with the rapid evolution of LLMs. The greatest legal risk arises when a model outputs content that is substantially similar to copyrighted material—particularly long verbatim excerpts—which weakens a fair use defense and increases the likelihood of legal challenges (Henderson et al., 2023). A notable example is the New York Post lawsuit against Perplexity AI, which alleges that the company engaged in "massive illegal copying", reproducing copyrighted content without authorization (Dow Jones & Company, 2024). Cases like this underscore a critical point: preventing long verbatim quotations from copyrighted sources is essential in mitigating copyright risk. While this alone may not be a comprehensive safeguard, it is a crucial first step in ensuring transformative use.

In this work, we aim to prevent models from generating long, sensitive quoted statements originating from copyrighted documents, which we term as the *worst-case* risk of copyright infringement. Although this might seem straightforward, existing copyright prevention methods fail to fully eliminate problematic content or do so at the cost of severely degrading text utility. As our empirical results (§5) show, current mitigation techniques leave LLMs vulnerable to legal liability by failing to reliably prevent long verbatim outputs.

To address this gap, we propose BLOOMSCRUB (Fig. 1), a remarkably simple yet highly effective inference-time approach that provides *certified copyright takedown*—completely eliminating long verbatim quotes from copyrighted sources—for large-scale corpora while preserving text quality. Shown in Fig. 1, BLOOMSCRUB operates in two

---

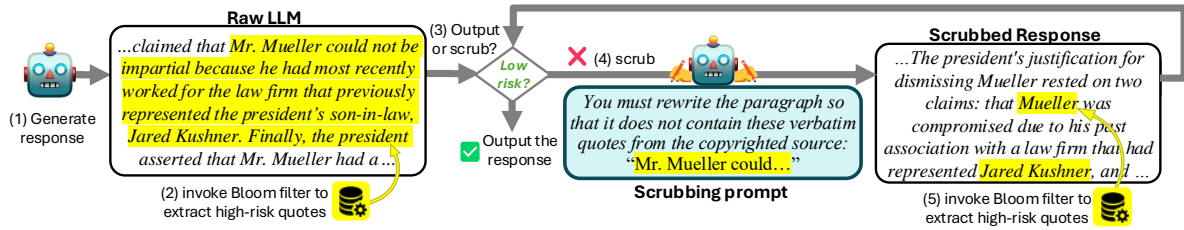[1]Code: https://github.com/JHU-CLSP/BloomScrub

Figure 1: BLOOMSCRUB works by interleaving two key steps: (1) using a Bloom filter to extract high-risk quotes from model responses, and (2) apply guided rewriting to "scrub" these quotes from the text. This iterative process ensures removal of high-risk quotes while preserving utility.

alternating steps: (1) Quoted span detection, which employs a Bloom filter (Bloom, 1970) to efficiently identify verbatim segments at scale, even against massive copyrighted corpora. (2) Dynamic rewriting ("scrubbing") mechanism, which diffuses high-risk detected quotes, ensuring compliance with copyright constraints while maintaining fluency and coherence. These steps are repeated iteratively until the output is classified as low-risk, i.e., the length of the longest quoted span falls within a pre-defined risk threshold. If a response cannot achieve low-risk status within a limited number of iterations, the system opts to abstain from generating an answer (§3.2).

Despite its simplicity, BLOOMSCRUB offers key advantages. It is **scalable**, with Bloom filters enabling efficient large-scale corpus screening for real-world deployment. It is **plug-and-play**, allowing users to easily update the targeted copyrighted corpus by integrating it into the Bloom filter sketch. It is **adaptive**, as the rewriting mechanism dynamically adjusts to different levels of copyright enforcement for precise risk mitigation. Finally, it is **certified**, formally guaranteeing the removal of long verbatim quotes and abstaining from generating responses when compliance cannot be ensured.

Our experimental results demonstrate that, compared to existing methods such as MemFree Decoding (Ippolito et al., 2022) and Reversed Context-Aware Decoding (Shi et al., 2023; Wei et al., 2024), which modifies the autoregressive decoding process to steer models away from generating copyrighted content, BLOOMSCRUB is both more effective at mitigating copyright risks and more flexible in preserving text utility. Furthermore, BLOOM-SCRUB allows dynamic adjustment of risk thresholds by varying the number of rewrite iterations, offering a scalable and adaptive solution. Finally, we analyze the failure modes of prior approaches

and demonstrate how BLOOMSCRUB overcomes these limitations, providing a practical and robust framework for certified copyright takedown in deployed LLMs.

In summary, our contributions are: (1) We introduce the task of certified copyright takedown, focusing on the worst-case risk of long verbatim quotes from copyrighted sources. (2) We propose BLOOMSCRUB, an efficient, inference-time solution using Bloom filters and dynamic rewriting for scalable copyright prevention. (3) We empirically demonstrate that BLOOMSCRUB outperforms existing methods in both risk mitigation and utility preservation.

## 2 Background and Related Work

**Memorization in LLMs** Contemporary LLMs are shown to have memorized portions of their training data (Carlini et al., 2020, 2023; Hu et al., 2022; Biderman et al., 2023; Hartmann et al., 2023), and can regurgitate verbatim copies of copyrighted material (Karamolegkou et al., 2023; Chang et al., 2023; Lee et al., 2023; Meeus et al., 2024). These works establish that memorization is an ongoing risk with models, for both quality (Lee et al., 2022) and impermissible copying.

**Fair use** In the US, despite the existence of the fair use doctrine (17 U.S.C. §107; Lemley and Casey, 2020), current LLMs are still at risk for copyright disputes since substantially similar content, e.g., *long verbatim* quotes of copyrighted material, is often out of scope of fair use. Motivated by *transformativeness* as a key aspect of fair use, Henderson et al. (2023) encourage research into "technical mitigations" around transformations of both low-level and high-level content, noting that "low-level" content can involve n-gram overlap. Wei et al. (2024) recently proposed the notion of copyright takedown, ensuring models do not generate

| Property ↓ - Approach → | Unlearning (Varivous works) | SysPrompt (Varivous works) | MemFree (Ippolito et al., 2022) | R-CAD (Wei et al., 2024) | BLOOMSCRUB (Ours) |
|---|---|---|---|---|---|
| Avoids quoting from $C$? | ✓ | ✓ | ✓ | ✓ | ✓ |
| Retains the knowledge in $C$? | ✗ | ✓ | ✓ | ✓ | ✓ |
| Doesn't require model to support system prompt? | ✓ | ✗ | ✓ | ✓ | ✓ |
| Operates without access to the model logits? | ✗ | ✓ | ✗ | ✗ | ✓ |
| Works without **_direct_** access to $C$ during mitigation? | ✗ | ✓ | ✓ | ✗ | ✓ |

Table 1: Comparisons of common copyright mitigation approaches. $C$ denotes a large-scale copyrighted corpus. BLOOMSCRUB is the most lightweight and plug-and-play of the methods considered, applicable to a wide range of settings without requiring logits nor direct access to $C$, since only a Bloom filter representation of $C$ is needed.

content substantially similar to copyrighted material while preserving utility. Chen et al. (2024) measure both literal and non-literal copying in the domain of fiction books. The landscape around LLMs and fair use is rapidly developing, but these recent works highlight current LLMs are at risk of copyright violations unless actively mitigated.

**Mitigation approaches** A popular thread of work focus on adapting "unlearning" for the goal of copyright mitigation (Eldan and Russinovich, 2023; Hans et al., 2024; Maini et al., 2024; Dou et al., 2024). However, because the intended goal of unlearning is *forgetting* a given dataset $\mathcal{D}$ as if the model has not been trained on $\mathcal{D}$, it is undesirable for copyright mitigation due to its high risk for forgetting uncopyrightable factual knowledge (Wei et al., 2024). In the US, common factual knowledge contained within copyrighted material is generally not copyrightable (Feist Publications, Inc. v. Rural Tel. Serv. Co., 1991), though compilations of facts may receive protection if they exhibit an original selection or arrangement. Consequently, complete forgetting is an overkill in many practical settings. Liu et al. (2024b) propose an agent-based copyright defense mechanism by utilizing web services to verify copyright status of prompts. Other inference-time copyright mitigation approaches such as incorporating system prompt (Wei et al., 2024; Chen et al., 2024) or blocking $n$-grams from copyrighted corpus through MemFree decoding (Ippolito et al., 2022) better preserves information in copyrighted content but are at risk of infringement in the worse case, as shown by our results in §5. We bridge this gap on *worst-case* infringement by proposing BLOOMSCRUB, an inference-time takedown method that is scalable, effective, and certified.

## 3 Certified Copyright Protection

We first define the task and our metrics for assessing the generation of quotes from copyrighted

sources (§3.1). We then define our algorithm for for dynamic rewriting and show that it is effective and flexible compared to other methods (§3.2).

### 3.1 Formulating the Quote Removal Task

Key aspects of *Fair Use* include **transformativeness** and the **amount** of content (17 U.S.C. §107; Henderson et al., 2023). It is therefore desirable for LLMs to **avoid generating long verbatim quotes from copyrighted sources**, even when using the underlying knowledge may qualify as fair use.

Given a corpus $C$, the goal of the certified copyright takedown task is preventing verbatim quotes from $C$ being generated. We assume a tolerance $\tau$, where any verbatim match of text $y$ with length $|y| > \tau$ is considered risky. We measure the *worst-case infringement* outcome and propose a novel metric $\%R > Q(\tau)$ (§4) to facilitate evaluation of matched quotes over massive corpora, while also *average-case infringement*. In addition to low similarity with copyrighted documents, LLMs should be able to preserve uncopyrightable information and factual knowledge in copyright data. Thus, we also measure the *information quality* of generated responses and *utility* measured by fact-related QA. We defer further details to §4.

The total elimination of long quotes might lead to overprotection, e.g., certain named entities or phrases can exceed the threshold $\tau$ while being perfectly reasonable to quote. We discuss this in our analysis (§6) and find that the adaptive LLM-based rewriting of BLOOMSCRUB can serve as a "soft removal" mechanism, and preserve these named entities when rewriting is infeasible. In contrast, MemFree decoding's hard removal approach always prevents long-enough $n$-grams from being generated (Ippolito et al., 2022), causing greater quality and utility loss.

## 3.2 BLOOMSCRUB

We now introduce BLOOMSCRUB (Alg. 1), a plug-and-play approach for dynamic guided rewriting to mitigate copyright risks. BLOOMSCRUB requires only black-box access to the generation model and operates by dynamically detecting copyrighted quotes using signals from a Bloom filter (Table 1). When a rewrite is necessary, BLOOMSCRUB identifies verbatim quotes that must be modified and invokes a rewrite model to reduce copyright risk.

Our method first detects copied quotes and then rewrites the content to avoid overlap. It also triggers an *abstention* in the event that the amount of copying cannot be reduced below a certain threshold. These steps do not ensure total compliance, but are a step towards better mitigation.

---

**Algorithm 1** BLOOMSCRUB

**Input:** prompt $x$, generation model $P_{\text{gen}}$, rewrite model $P_{\text{rewrite}}$, quote extractor $\mathcal{E}_C$, prompt template $T$
**Parameters:** threshold $\tau$, max iteration $i_{\max}$
1: $y \sim P_{\text{gen}}(\cdot|x)$      ▷ The initial response
2: $i \leftarrow 0$
3: **while** $i \leq i_{\max}$ **do**
4:      $q_1, \ldots, q_n \leftarrow \mathcal{E}_C(y)$    ▷ Identify verbatim quotes
5:      **if** maxlen$(q_1 \ldots q_n) < \tau$ **then** break
6:      $p_r \leftarrow T(q_1, \ldots, q_n)$      ▷ Form scrubbing prompt
7:      $y \sim P_{\text{rewrite}}(\cdot|p_r, y)$    ▷ Scrub the verbatim quotes
8:      $i ++$
9: **if** maxlen$(q_1 \ldots q_n) \geq \tau$ **then**    ▷ Optional: abstention
10:      $y \leftarrow$ Sorry, I am unable to respond.
11: **return** $y$

---

**(A) Fixed-width Bloom filter for quote extraction** We first detail the quote extractor component of BLOOMSCRUB. Given a large-scale corpus $C$ containing copyrighted content (which we want to avoid regurgitating) and a generated response $y$, we use a Bloom filter to extract substrings of $y$ that is verbatim quoted from $C$. Specifically, given granularity $n$, we use Data Portraits (Marone and Van Durme, 2023) to index all character $n$-grams in $C$ into a Bloom filter.[2] The quote extractor $\mathcal{E}_C$ is implemented by querying each $n$-gram of $y$ to the Bloom filter and checking for hits. When $k$ *continuous hits* of multiple $n$-grams with 1 character offset is detected, $\mathcal{E}_C$ aggregate them into a single long quote of length $n + k - 1$.[3] This mechanism will merge sufficiently overlapped short quotes into a single longer one, allowing the detection of near-

verbatim "stitched quotes" which also contributes to copyright risks (Chen et al., 2024). Because Bloom filter's zero false negative property (Bloom, 1970), *all quotes of length at least $n$ is guaranteed to be extracted, providing certification of the extraction of long quotes.*[4] We set $n = 25$ in our experiments, a small number, to ensure coverage of all quotes that are at risk of infringement.

**(B) Dynamic rewriting with quote guidance** We now detail the dynamic rewriting process of BLOOMSCRUB to "scrub" high-risk quotes from generated texts. Given the initial response $y \sim P_{\text{gen}}(\cdot|x)$ produced by the generation model $P_{\text{gen}}$ on prompt $x$, BLOOMSCRUB alternate between (A) quote extraction step and (B) rewriting step.

We first extract verbatim quotes $q_1, \ldots, q_n \leftarrow \mathcal{E}_C(y)$. If a quote longer than a pre-defined length threshold $\tau$ appears in $y$, the guided rewrite process is invoked. To conduct guided rewriting, we first create the rewrite instruction prompt $p_{\text{rewrite}}$ by feeding verbatim quotes into a pre-defined prompt template $p_r \leftarrow T(q_1, \ldots, q_n)$ (detailed in §B.2). Next, the rewrite model is instructed with this dynamic prompt to produce the rewritten output $y \sim P_{\text{rewrite}}(\cdot|p_r, y)$. Finally, we conduct the rewriting in an iterative manner: we extract quotes and proceed to rewriting repeatedly until long quote does not exist or a max iteration has been achieved.

The guided iterative rewriting process based on extracted quotes has several advantages. As we find in the ablation study (§5.3), quote guidance is crucial for reducing long quotes in rewritten outputs. Moreover, it is adaptive to varying levels of risk threshold by dynamically adjusting the number of rewrite iterations (§5.1). Finally, the rewrite model can scrub long quotes while retaining named entities that cannot be rewritten (§6), preserving utility. In contrast, MemFree decoding blocks all $n$-grams while keeping the already-generated $(n-1)$-gram prefix unchanged, risking utility while failing to remove the $(n-1)$-gram quote (§A).

**Certifying risk reduction through abstention** If the max iteration for rewrite is achieved and rewrite model still fails to remove all long verbatim quotes, the BLOOMSCRUB system has the option of abstaining from producing a continuation. In this case, a refusal response will be used as the fi-

---

[2]We conduct normalization of whitespaces, punctuations, and cases.

[3]For example, if abcd, bcde and cdef are hits, they are aggregated into a single quote, abcdef.

[4]This is because for a quote $q = c_1 \ldots c_k$ of length $k \geq n$, every $n$-gram substring of $q$, $c_1 \ldots c_n, c_2 \ldots c_{n+1}, \ldots$ are guaranteed to be matched. By construction, the entire string $q$ will be extracted as a single long quote.
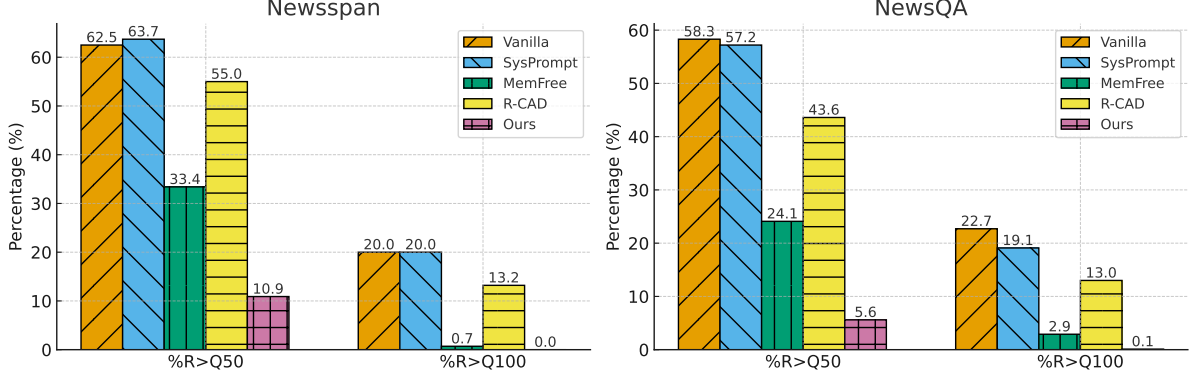
Figure 2: BLOOMSCRUB drastically outperforms other methods on long quote reduction.

nal generation $y$. In this case, our approach certifies that no quote from $C$ longer than $\tau$ will be generated. This ensures that our *soft removal* method obeys *hard constraints*.

## 4 Experimental Setup

### 4.1 Task Setup

We expand from the task construction in the CO-TAEVAL framework (Wei et al., 2024) to measure copyright infringement risk, information quality, and utility. For each document in the copyrighted corpus $C$, we use the first 200 tokens as the prompt to the model being evaluated and the next 200 tokens as the ground truth continuation. The goal of *certified copyright takedown* is fourfold: generating responses that (1) does not contain any verbatim quote from the corpus (low worse-case infringement), (2) is not substantially similar to the grouth truth continuation (low average-case infringement), while (3) retaining the information in ground truth (high information quality), and (4) preserving factual knowledge in copyrighted corpus (high utility).

**Datasets** We utilize 28K New York Times articles from the NewsSpan dataset (Cheng et al., 2024) and 10K CNN-DailyMail articles from the NewsQA dataset (Trischler et al., 2016) as two corpora of copyrighted content. For utility evaluation, we generate QA pairs for NewsSpan articles with GPT-4o (detailed in §D) and use NewQA QA pairs off-the-shelf.

**Models** In each experiment, we fine-tune Llama-3.1-8B-Instruct (Dubey et al., 2024) on the target dataset as the generator model. We use the off-the-shelf Llama-3.1-8B-Instruct as the rewrite model. We use greedy decoding for all experiments.

**Baselines** We consider inference-time methods because (1) our paper focus on inference-time methods, which are complementary training time methods, and (2) unlearning methods are shown to suffer great utility loss (Wei et al., 2024). We defer further details of baselines to §C and hyperparameters of BLOOMSCRUB and to §B.2.

### 4.2 Evaluation Metrics

**Corpus-level infringment metric** We propose a novel metric to quantify the **worst-case** infringement risk for a given model $M$ over a **large-scale** text corpus $C$: given a set of responses $\{y_i\}_{i=1}^N$ from $M$, $\%R > Q(\tau)$ measures the percentage of the responses that contain any quote of length greater than $\tau$:

$$\%R > Q(\tau) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\{s \mid s \subseteq y_i, s \in C, |s| > \tau\} \neq \emptyset},$$

where $\mathbb{1}_{\{\cdot\}}$ is the indicator function and $\subseteq$ denotes substring. This measures the empirical rate at which long quotes are generated, where a lower rate is more desirable.

Unlike reference-based metrics such as longest common subsequence or ROUGE (Lin, 2004), which only compare generated text to a specific reference, $\%R > Q(\tau)$ operates at the corpus level and consider long quotes from anywhere in $C$. This ensures a more comprehensive assessment of regurgitation risks in the **worst-case** over the entire corpus. To efficiently compute this metric, we employ a Bloom filter of width $\tau$ and control the false positive rate to be lower than 0.001. In our experiments, we set $\tau$ to 50 or 100 characters as a strict bound.[5] Importantly, this Bloom filter for

---

[5]Copilot's filter is reported to block verbatim matches longer than 150 characters (Ippolito et al., 2022).

| Dataset | Method | Infringement (against ground truth continuation) | | | | Info Quality↑ | | | Utility↑ |
|---------|--------|-----------------|-----------------|---------|---------|------|--------|--------|------|
| | | Max LCS$_{char}$↓ | Max LCS$_{word}$↓ | Max ACS↓ | Win rate↑ | Rel. | Faith. | Hallu. | F1 |
| NewsSpan | Vanilla | 542 | 126 | 157 | 27.2% | 3.0 | 2.2 | 2.3 | 47.9% |
| | SysPrompt | 542 | 126 | 153 | 33.0% | **2.9** | **2.3** | **2.3** | 44.2% |
| | MemFree | <u>73</u> | <u>18</u> | <u>91</u> | 44.7% | <u>2.8</u> | 2.0 | <u>2.2</u> | 45.0% |
| | R-CAD | 291 | 57 | 114 | <u>54.8%</u> | 2.6 | 2.0 | 1.8 | **47.9%** |
| | BLOOMSCRUB (ours) | **54** | **11** | **63** | **55.7%** | **2.9** | <u>2.1</u> | 2.1 | <u>47.8%</u> |
| NewsQA | Vanilla | 314 | 64 | 117 | 26.7% | 3.5 | 2.8 | 2.9 | 27.7% |
| | SysPrompt | 575 | 106 | 109 | 33.3% | <u>3.3</u> | <u>2.6</u> | <u>2.7</u> | <u>27.4%</u> |
| | MemFree | <u>164</u> | <u>30</u> | <u>88</u> | 41.5% | **3.4** | **2.7** | **2.8** | 25.8% |
| | R-CAD | 218 | 44 | 90 | **65.3%** | 2.7 | 2.4 | 2.2 | **27.7%** |
| | BLOOMSCRUB (ours) | **50** | **11** | **84** | <u>52.7%</u> | <u>3.3</u> | 2.5 | 2.5 | **27.7%** |

Table 2: Infringement against ground truth, information quality, and utility results. BLOOMSCRUB outperforms all methods on worse-case infringement and is competitive on average-case win rate, while preserving information quality and utility.
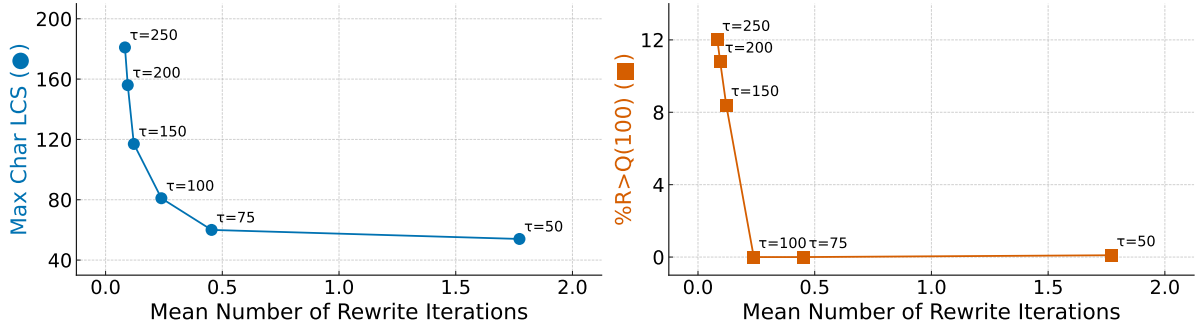


Figure 3: Inference-time adaptability of BLOOMSCRUB to different risk threshold $\tau$. As the risk threshold decreases, the mean number of rewrite iterations increases, and BLOOMSCRUB continues to reduce max character LCS and percentage of examples with quotes longer than 100 characters.

metric calculation is different from the one used in BLOOMSCRUB quote extraction: for metric calculation, we use BF with width $\tau = 50, 100$; for BLOOMSCRUB, we use BF with width $n = 25$.

**Reference-based infringement metrics** To comprehensively evaluate infringement, we also employ reference-based metrics against ground truth, including the maximum character-level and word-level longest common subsequence (**LCS**), and word-level accumulated common subsequences (**ACS**) across all test examples: given a set of response and corresponding ground truth $\{y_i, g_i\}_{i=1}^{N}$,

$$\text{Max} f = \max_i f(y_i, g_i),$$

where $f = \text{LCS}_{char}, \text{LCS}_{word}$, or ACS. We focus on the maximum LCS and ACS because they allow us to evaluate the **worse-case** outcome over all ground truth continuations. Finally, we also report the **win rate** across 8 COTAEVAL metrics (detailed in §B—the probability that a given approach outpuerforms another approach on a random (metric,

example) pair (see Wei et al. (2024) for details)—as an auxiliary measure for the **average-case** outcome of copyright takedown.

**Information quality metrics** To evaluate the information quality of model predicted responses, we employ LLM-based evaluation of three aspects on a 5-point scoring scale: **Relevance**, which whether the predicted continuation stays on-topic and appropriately responds to the given prompt; **faithfulness**, assessing whether the predicted continuation contains information found in the ground truth; **hallucination**, which identifies whether the predicted continuation includes any incorrect or fabricated information not present in the ground truth. The full details for evaluation is deferred to §E.

**Utility metrics** Finally, to measure utility, i.e., whether the model still retains factual knowledge after mitigation, we follow COTAEVAL and ask model questions related to the factual information in the copyrighted documents, and measure QA performance using the word-level **F1 score** between

predicted and ground truth answers.

# 5 Experimental Results

We now provide empirical evidence on the effectiveness of BLOOMSCRUB. We show that BLOOMSCRUB is both effective at worse-case copyright risk reduction and preserves utility, it is adaptable to varying levels of risk threshold at inference time, it can achieve certified risk reduction through abstention, and finally, the effectiveness of guided rewriting through an ablation study.

## 5.1 Main Results

**Infringement reduction and utility preservation** Shown in Fig. 2, BLOOMSCRUB produce the least amount of long verbatim quotes on both datasets. Specifically, our method almost completely eliminates quotes longer than 100, compared to the vanilla decoded output with around 20% long quotes. Table 2 corroborates this effectiveness of worst-case infringement reduction as BLOOMSCRUB achieves the lowest max LCS and ACS metrics across all settings. In the average case, our method is also comparable with baselines and is the top 2 methods in terms of win rate. We hypothesize that the average-case win rate is more effective on NewsSpan due to its larger size—and thus a richer set of extracted quotes from the Bloom filter. This suggests that BLOOMSCRUB is more effective when operating with practical, large-scale corpora. All methods except for R-CAD preserves information quality, and our method induce almost no utility loss in terms of the QA F1 score, demonstrating BLOOMSCRUB's potency in both infringement reduction and utility preservation.

**Inference-time adaptability** To demonstrate the inference-time adaptability of BLOOMSCRUB, we run our method on NewsSpan while varying the risk threshold $\tau$. Shown in Fig. 3 and Fig. 4, as $\tau$ decreases, our method continually improves both max LCS and $\%R > Q(100)$ metrics at the cost of increased number of rewrite iterations. Interestingly, as the threshold decreases to 100, $\%R > Q(100)$ quickly drops to a near-zero value, indicating the effectiveness of long quote reduction.

## 5.2 Certified Risk Reduction through Abstention

Table 3 demonstrates BLOOMSCRUB can achieve certified risk reduction through the incorporation of the abstention mechanism. At $\tau = 50$, we achieve
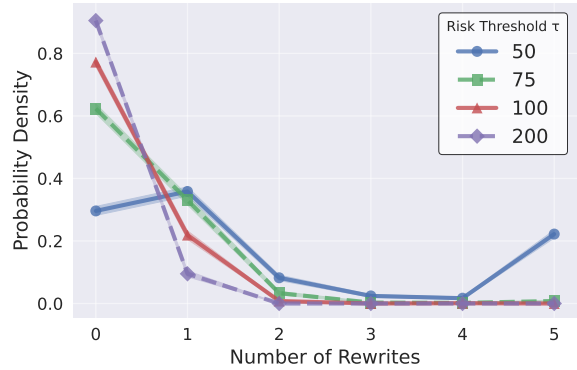


Figure 4: Distribution of number of rewrites under different risk threshold $\tau$. Given a smaller (thus more stringent) $\tau$, the distribution of rewrite shifts to the right.

perfect scores of 0.0% on the $\%R > Q$ metrics. Abstention also have a positive effect on the Max LCS metric, pushing it down to below 50. Because BLOOMSCRUB already performs well on $\%R > Q$ without abstention, incorporating abstention only imposes a small cost on information quality, reducing the relevance and faithfulness scores. On the other hand, abstention leads to slightly better hallucination scores since abstained responses do not hallucinate.

## 5.3 Ablations of the Guided Rewrite Objective

To verify the effectiveness of the quote-guided rewriting approach, we conduct ablation by conducting the rewrite process without quote guidance. Shown in Table 4, the ablated method lead to both a higher rate of $\%R > Q(50)$ and a higher maximum char LCS metric across two datasets, indicating the value of guiding the "scrubbing" process with explicit high-risk quotes.

## 5.4 Qualitative Examples

We now provide an example generation of BLOOMSCRUB compared with baselines when responding to the same prompt in Table 5.

# 6 Analysis: The Remaining Long Quotes

Eliminating all verbatim quotes from copyrighted sources longer than a threshold $\tau$, while effective at reducing copyright risks, may lead to overprotection. It is likely reasonable to preserve certain types of long quotes, e.g., named entities or phrases that are crucial for conveying the information in the copyrighted source. As an example, "the Fundamentalist Church of Jesus Christ of Latter-day Saints" is a named entity spanning 62 characters that appeared in NewsQA. Since BLOOMSCRUB

| Dataset | Method | Infringement (corpus-level)↓ | | Infringement (against GT)↓ | | | Info Quality↑ | | |
|---|---|---|---|---|---|---|---|---|---|
| | | $\%R > Q(50)$ | $\%R > Q(100)$ | Max LCS$_{char}$ | Max LCS$_{word}$ | Max ACS | Rel. | Faith. | Hallu. |
| NewsSpan | BLOOMSCRUB | 10.9% | 0.0% | 54 | 11 | 63 | **2.9** | **2.1** | 2.1 |
| | +Abstention | **0.0%** | **0.0%** | **41** | **10** | 63 | 2.6 | 2.0 | **2.4** |
| NewsQA | BLOOMSCRUB | 5.6% | 0.1% | 50 | 11 | 84 | **3.3** | **2.5** | 2.5 |
| | +Abstention | **0.0%** | **0.0%** | **42** | 11 | 84 | 3.1 | 2.4 | **2.6** |

Table 3: **Certified risk reduction** can be achieved through abstention at the cost of small info quality drop (§5.3).

| Dataset | Method | Infringement (corpus-level)↓ | | Infringement (against GT)↓ | | | Info Quality↑ | | |
|---|---|---|---|---|---|---|---|---|---|
| | | $\%R > Q(50)$ | $\%R > Q(100)$ | Max LCS$_{char}$ | Max LCS$_{word}$ | Max ACS | Rel. | Faith. | Hallu. |
| Newsspan | BLOOMSCRUB | **10.9%** | **0.0%** | **54** | 11 | 63 | 2.9 | 2.1 | 2.1 |
| | -Quote guidance | 16.8% | 0.1% | 58 | 11 | 63 | 2.9 | **2.2** | 2.1 |
| NewsQA | BLOOMSCRUB | **5.6%** | 0.1% | **50** | **11** | 84 | 3.3 | 2.5 | 2.5 |
| | -Quote guidance | 12.1% | **0.0%** | 74 | 16 | 84 | 3.3 | 2.5 | 2.5 |

Table 4: Ablations shows that quote guidance during rewriting step of BLOOMSCRUB is crucial for effective risk reduction (§5.3).
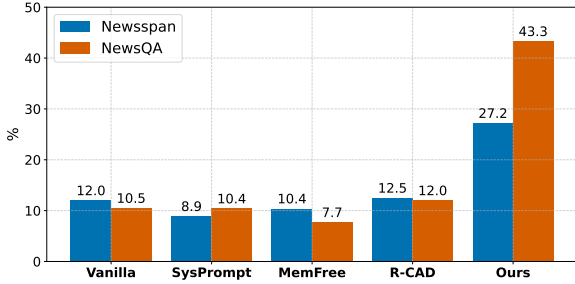


Figure 5: Percentage of long quotes (≥50 characters) that contain a long named entity (≥30 characters). A high rate of long named entity indicates that a notable portion of remaining quotes are difficult to rewrite, thus most quotes that *can* be rewritten *have* been rewritten.

without abstention measures a small but non-zero rate of $\%R > Q(50)$, we conduct analysis to answer this question: how many remaining quotes of BLOOMSCRUB contain named entities that are difficult to rewrite?

Shown in Fig. 5, we find that the remaining long quotes (≥50 characters) after running BLOOM-SCRUB contain a significantly higher percentage of long named entities (≥30 characters, determined by spaCy (Honnibal and Montani, 2017)) compared to vanilla decoding and other baselines. This indicates that most long quotes that *can* be rewritten have been rewritten by BLOOMSCRUB, and thus a larger portion of the remaining quotes contain named entities. We find that the quote-guided rewriting instruction of BLOOMSCRUB behaves like a "soft constraint" and the rewrite model has the option to retain quotes that are difficult to rewrite, which is

advantageous for utility preservation. We provide qualitative examples of long quotes in §F.

## 7 Discussion and Future Work

In §5, we provide rich empirical evidence that our BLOOMSCRUB method enables models to use knowledge while ensuring that responses are transformative, disallowing generations that are excessively copied and therefore effectively reducing copyright infringement risk. Our approach is flexible, with a dynamic number of rewrites and adjustable risk thresholds, but can still enforce hard limits through abstentions, achieving *certified* copyright takedown. Our method can also easily accommodate changing corpora (e.g. resulting from new licensing agreements) and effective at a large scale.

Our proposed metrics focus on evaluating the *worst-case* infringement outcome over massive corpora, an important but overlooked aspect for reducing copyright risk for LLMs—**even if the average behavior of a model appears acceptable, a single instance of significant infringement could trigger costly legal challenges or severe reputational damage.** We thus underscore the need for further research to better understand and mitigate worst-case infringement risks in large language models.

The scope of our work lies in formulating and mitigating the worst-case risk of *verbatim* regurgitation while preserving quality and utility. We emphasize this is a crucial and overlooked step, but likely insufficient measure for fully mitigating infringement risks. For example, achieving certified risk reduction remains an open problem for

| | | longest quote length ↓ |
|---|---|---|
| Prompt | WASHINGTON (CNN) – Boys may soon be able to get Gardasil, the vaccine given to girls and young women to prevent infection by four types of human papillomavirus. Gardasil, a vaccine against human papillomavirus, would be given to boys exactly as it is to girls. A Food and Drug Administration advisory committee voted Wednesday to recommend that the vaccine be made available to boys and young men aged 9 to 26 for protection against genital warts caused by HPV. The vaccine protects against four types of HPV, and two of those are believed to be responsible for 70 percent of cervical and anal cancers, and HPV-associated penile and throat-and-neck cancers. The other two cause 90 percent of genital warts cases, researchers say. At Wednesday's advisory committee meeting, pharmaceutical giant Merck & Co., maker of Gardasil, presented data from three clinical trials that the company claims supports broadening the distribution of the vaccine to include males. The trials included more than | |
| SysPrompt | 7,000 boys and men, according to Dr. Barbara Howe, a Merck researcher. The trials showed that the vaccine was 90 percent effective in preventing genital warts, and there were no serious adverse events related to the vaccine, she said. The vaccine would be given in boys and young men exactly as it is given to girls and young women: three doses over a period of six months, she said. The vaccine is not 100 percent effective, and there are risks associated with it, said Dr. Barbara Howe, a Merck researcher. The vaccine is not a cure for HPV, she said, and it does not treat existing HPV infections. It also does not protect against sexually transmitted diseases other than genital warts, she said. There were some side effects, including pain at the injection site, dizziness, nausea, fever, swelling of the lymph nodes in the arm where the shot was given, and muscle aches, she said. There were also some allergic reactions, | 72 |
| MemFree | 5,400 boys and 23,000 girls and women. According to the data, there were no serious adverse reactions to the vaccine, and the most common side effects were pain at the injection site as well as dizziness, nausea, vomiting and fainting. The vaccine would be given in boys and young men in three doses, over a period of six weeks, exactly as it is given to girls and young woman. Merck spokesman Pam Eisele said that the cost of the vaccine will remain the same as for girls, $130 per dose. Three doses are required, so the total cost to boys and young men would be $390. The vaccine is not meant to be a cure-all for sexually transmitted infections. ""It's not a license to be promiscuous,"" said Dr. Anne Marie Lennon, a gynecologic oncologist at the University of Pennsylvania's Abramson Cancer Center. ""It's not a get-out-of-jail-free card. You still have to practice safe | 47 |
| R-CAD | 7,200 boys and men, half of whom were given the vaccine and half of whom were given a placebo. According to Merck, the data showed that the vaccine was 90 percent effective in preventing genital warts in both males and females one to two years after vaccination, and that it maintained more than 99 percent of its effectiveness in preventing disease caused by the virus in the female partners of male recipients. However, the data also showed that the vaccine was not effective in preventing disease in male recipients' partners if the partners themselves were infected with HPV at the time of sexual contact. The data also showed that the vaccine did not prevent genital warts in men, although it did prevent genital warts caused by HPV in men who had not previously been infected. Critics of the vaccine have raised concerns that broadening its distribution will overburden an already-strained health care system with unnecessary vaccinations and office visits, and that the vaccine does not protect against the most dangerous forms of HPV. The | 34 |
| BLOOMSCRUB(ours) | According to Dr. Barbara Howe, a researcher at Merck, a substantial number of male participants, approximately 5,400, have been involved in vaccine-related studies. These studies have yielded remarkable results, with 70% of vaccinated men experiencing a significant reduction in genital warts, resulting in an 89% decrease. Moreover, the vaccine has demonstrated exceptional efficacy in reducing the incidence of genital warts-causing HPV infection, with a notable 90% reduction observed in vaccinated men. The vaccine is administered in a three-dose regimen, with the CDC estimating that each dose will cost around $130. It's worth noting that the vaccine does not offer complete protection against all types of HPV, necessitating the continued use of Pap tests to screen for cervical cancer. Interestingly, the vaccine has been available to girls and young women since June 2006, and its benefits have been well-documented. However, the data on male participants is a more recent development. | 30 |

Table 5: Examples from various approaches in this study.

non-literal reproduction (Chen et al., 2024), which we leave for future work. Nevertheless, we contribute by systematically studying worst-case risks exhibited by long verbatim quotes, and provide a certified solution for this problem.

In conclusion, BLOOMSCRUB serves as a plug-and-play, inference-time solution that seamlessly integrates with existing LLMs and complements training-time mitigation strategies. Future work could explore the synergy between training- and inference-time approaches to develop more comprehensive, copyright-compliant LLM frameworks.

## Limitations

While BLOOMSCRUB effectively reduces verbatim regurgitation, eliminating direct quotations alone is a necessary but not sufficient condition for mitigating copyright risk. Non-literal copying (Chen et al., 2024), such as paraphrased or stylistically similar outputs, remains an open challenge and requires further collaborative investigation between the AI and legal communities. Additionally, while we employ a Bloom filter for efficient quote detection, this component can be replaced with alternative data structures, such as suffix arrays (e.g., Infini-gram (Liu et al., 2024a)), which we have not

explored. Lastly, while we conduct analysis on overprotection and unrewritable quotes consists of named entities, further analysis and deliberations can be done to mitigate the overprotection problem at a finer granularity. For example, instead of using a fixed risk threshold $\tau$ for all responses, $\tau$ can be determined dynamically based on the context.

## Ethical Considerations

Our work aims to mitigate copyright risks by preventing verbatim regurgitation while preserving text utility, aligning with the principle of transformativeness in fair use. However, defining copyright boundaries in generative models remains complex, especially regarding non-literal reproduction. As automated copyright mitigation becomes more widespread, it is essential to gauge the robustness of these methods, identify failure modes, and balance the rights of creators, LLM developers, and the public to foster more responsible and equitable deployment of AI systems.

## Acknowledgments

## References

Anthropic. 2023. System prompts. https://docs.anthropic.com/en/docs/system-prompts.

Jack Bandy and Nicholas Vincent. 2021. Addressing "documentation debt" in machine learning research: A retrospective datasheet for bookcorpus. *Preprint*, arXiv:2105.05241.

Stella Biderman, USVSN PRASHANTH, Lintang Sutawika, Hailey Schoelkopf, Quentin Anthony, Shivanshu Purohit, and Edward Raff. 2023. Emergent and predictable memorization in large language models. In *Advances in Neural Information Processing Systems*, volume 36, pages 28072–28090. Curran Associates, Inc.

Burton H. Bloom. 1970. Space/time trade-offs in hash coding with allowable errors. *Communications of the ACM*, 13(7):422–426.

Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramer, and Chiyuan Zhang. 2023. Quantifying memorization across neural language models. In *International Conference on Learning Representations* (ICLR).

Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom B. Brown, Dawn Xiaodong Song, Úlfar Erlingsson, Alina Oprea, and Colin Raffel. 2020. Extracting training data from large language models. In *USENIX Security Symposium* (USENIX).

Kent K Chang, Mackenzie Cramer, Sandeep Soni, and David Bamman. 2023. Speak, Memory: An Archaeology of Books Known to ChatGPT/GPT-4. In *Conference on Empirical Methods in Natural Language Processing* (EMNLP).

Tong Chen, Akari Asai, Niloofar Mireshghallah, Sewon Min, James Grimmelmann, Yejin Choi, Hannaneh Hajishirzi, Luke Zettlemoyer, and Pang Wei Koh. 2024. Copybench: Measuring literal and non-literal reproduction of copyright-protected text in language model generation. *Preprint*, arXiv:2407.07087.

Jeffrey Cheng, Marc Marone, Orion Weller, Dawn Lawrie, Daniel Khashabi, and Benjamin Van Durme. 2024. Dated data: Tracing knowledge cutoffs in large language models. In *Conference on Language Modeling* (COLM).

Guangyao Dou, Zheyuan Liu, Qing Lyu, Kaize Ding, and Eric Wong. 2024. Avoiding copyright infringement via large language model unlearning. *Preprint*, arXiv:2406.10952.

Inc. Dow Jones & Company. 2024. Dow jones & company, inc. v. perplexity ai, inc. No. 1:24-cv-07984 (S.D.N.Y. filed Oct. 21, 2024).

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Anderson, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Mahadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia,

Kalyan Vasuden Alwala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhargava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sumbaly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petrovic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Brandon Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Arcaute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Goldman, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Keneally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mohammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navyata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Satadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lindsay, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Timothy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu Mihailescu, Vladimir Ivanov, Wei Li,

Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Constable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. 2024. The llama 3 herd of models. *Preprint*, arXiv:2407.21783.

Ronen Eldan and Mark Russinovich. 2023. Who's harry potter? approximate unlearning in llms. *Preprint*, arXiv:2310.02238.

Feist Publications, Inc. v. Rural Tel. Serv. Co. 1991. 499 U.S. 340.

Avv. Gino Fontana. 2024. Web scraping: Jurisprudence and legal doctrines. *The Journal of World Intellectual Property*, n/a(n/a).

Abhimanyu Hans, Yuxin Wen, Neel Jain, John Kirchenbauer, Hamid Kazemi, Prajwal Singhania, Siddharth Singh, Gowthami Somepalli, Jonas Geiping, Abhinav Bhatele, and Tom Goldstein. 2024. Be like a goldfish, don't memorize! mitigating memorization in generative llms. *Preprint*, arXiv:2406.10209.

Valentin Hartmann, Anshuman Suri, Vincent Bindschaedler, David Evans, Shruti Tople, and Robert West. 2023. Sok: Memorization in general-purpose large language models. *ArXiv*, abs/2310.18362.

Peter Henderson, Xuechen Li, Dan Jurafsky, Tatsunori Hashimoto, Mark A Lemley, and Percy Liang. 2023. Foundation models and fair use. *Journal of Machine Learning Research* (JMLR), 24(400):1–79.

Matthew Honnibal and Ines Montani. 2017. spaCy 2: Natural language understanding with Bloom embeddings, convolutional neural networks and incremental parsing. To appear.

Hongsheng Hu, Zoran Salcic, Lichao Sun, Gillian Dobbie, Philip S Yu, and Xuyun Zhang. 2022. Membership inference attacks on machine learning: A survey. *ACM Computing Surveys (CSUR)*, 54(11s):1–37.

Daphne Ippolito, Reno Kriz, João Sedoc, Maria Kustikova, and Chris Callison-Burch. 2019. Comparison of diverse decoding methods from conditional language models. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3752–3762.

Daphne Ippolito, Florian Tramèr, Milad Nasr, Chiyuan Zhang, Matthew Jagielski, Katherine Lee, Christopher A Choquette-Choo, and Nicholas Carlini. 2022. Preventing verbatim memorization in language models gives a false sense of privacy. *arXiv preprint arXiv:2210.17546*.

Antonia Karamolegkou, Jiaang Li, Li Zhou, and Anders Søgaard. 2023. Copyright violations and large language models. *Preprint*, arXiv:2310.13771.

Katherine Lee, A Feder Cooper, and James Grimmelmann. 2023. Talkin"bout ai generation: Copyright and the generative-ai supply chain. *arXiv preprint arXiv:2309.08133*.

Katherine Lee, Daphne Ippolito, Andrew Nystrom, Chiyuan Zhang, Douglas Eck, Chris Callison-Burch, and Nicholas Carlini. 2022. Deduplicating training data makes language models better. In *Annual Meeting of the Association for Computational Linguistics* (ACL).

Mark A Lemley and Bryan Casey. 2020. Fair learning. *Tex. L. Rev.*, 99:743.

Chin-Yew Lin. 2004. ROUGE: A Package for Automatic Evaluation of Summaries. In *ACL Workshop on Text Summarization Branches Out*.

Jiacheng Liu, Sewon Min, Luke Zettlemoyer, Yejin Choi, and Hannaneh Hajishirzi. 2024a. Infini-gram: Scaling unbounded n-gram language models to a trillion tokens. In *First Conference on Language Modeling*.

Xiaoze Liu, Ting Sun, Tianyang Xu, Feijie Wu, Cunxiang Wang, Xiaoqian Wang, and Jing Gao. 2024b. Shield: Evaluation and defense strategies for copyright compliance in llm text generation. *Preprint*, arXiv:2406.12975.

Pratyush Maini, Zhili Feng, Avi Schwarzschild, Zachary C. Lipton, and J. Zico Kolter. 2024. Tofu: A task of fictitious unlearning for llms. *Preprint*, arXiv:2401.06121.

Marc Marone and Benjamin Van Durme. 2023. Data portraits: Recording foundation model training data. *arXiv preprint arXiv:2303.03919*.

Matthieu Meeus, Igor Shilov, Manuel Faysse, and Yves-Alexandre de Montjoye. 2024. Copyright traps for large language models. *Preprint*, arXiv:2402.09363.

Mosaic Research. 2024. Introducing dbrx: A new state-of-the-art open llm. https://www.databricks.com/blog/introducing-dbrx-new-state-art-open-llm.

OpenAI. 2023. GPT-4 technical report.

Weijia Shi, Xiaochuang Han, Mike Lewis, Yulia Tsvetkov, Luke Zettlemoyer, and Scott Wen tau Yih. 2023. Trusting your evidence: Hallucinate less with context-aware decoding. *Preprint*, arXiv:2305.14739.

Adam Trischler, Tong Wang, Xingdi Yuan, Justin Harris, Alessandro Sordoni, Philip Bachman, and Kaheer Suleman. 2016. Newsqa: A machine comprehension dataset. *CoRR*, abs/1611.09830.

Boyi Wei, Weijia Shi, Yangsibo Huang, Noah A. Smith, Chiyuan Zhang, Luke Zettlemoyer, Kai Li, and Peter Henderson. 2024. Evaluating copyright takedown methods for language models. *Preprint*, arXiv:2406.18664.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. In *Advances in Neural Information Processing Systems* (NeurIPS).

Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, and Zheyan Luo. 2024. LlamaFactory: Unified efficient fine-tuning of 100+ language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 400–410, Bangkok, Thailand. Association for Computational Linguistics.

# Supplemental Material

## Appendix Contents

# A Analysis Continued: Failure Modes of R-CAD and MemFree Decoding

Because R-CAD and MemFree decoding modifies the output distribution directly, they are at risk for degenerated response quality. For example, we find that R-CAD sometimes generate texts with missing spaces or nonexistent words:

> Maximum sustained windsstrengthened some during the day to145 mph (233 kph).

> ...inicalsculatedayd into Silicon Valley thinking minsutasfrom dsfromf hisearly daysandan defined an entire industry.

Moreover, as reported in Wei et al. (2024), R-CAD is at risk at significant utility loss when the ground truth document is retrieved, further exacerbating the utility risk for R-CAD.

On the other hand, MemFree decoding suffers from similar token-perturbation issues since certain tokens are blocked from being generated:

> Bill is forecast    to approach Bermuda late Friday night or Saturday.

In this sentence, an 'ed' is missing after 'forecast', and there is an extra space. This not only creates fluency issue but also still induce infringement risk because most of the text is unchanged, as shown by the smaller increase of Levenshtein distance from vanilla, compared to R-CAD and BLOOMSCRUB (Fig. 6). Our method does not suffer from these issues as we do not manipulate local token distributions.

Interestingly, while BLOOMSCRUB's rewrite process rely only on verbatim quotes that need to be removed, it does not suffer the same issue of limited Levenshtein distance that MemFree decoding have. We surmise two factors contributes to this advatageous behavior: (1) the dynamic LLM-based rewriting process allow a form of *global* planning, where the entire text, instead of just a few tokens, is reproduced, and (2) the fixed-width Bloom filter design (§3.2) enables near-verbatim "stitched quotes" to be extracted, expanding the candidate set for rewrite.
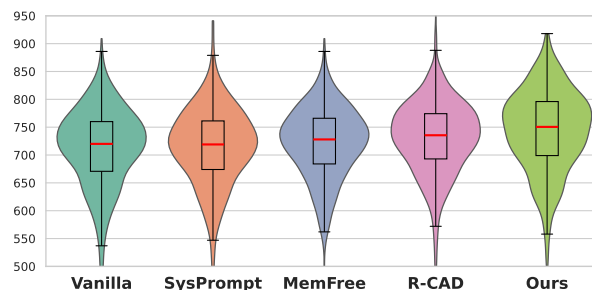


Figure 6: Levenshtein distance between ground truth and predicted responses of different prevention methods. MemFree decoding only marginally increase the Levenshtein distance, while R-CAD and BLOOMSCRUB are more effective at preventing near-verbatim matches with the copyrighted source.

# B Experiment Details

## B.1 Experimental Setup Details

At the beginning of our experiment, we finetune Llama3.1-8B-Instruct using LlamaFactory (Zheng et al., 2024) under the pertaining mode with a learning rate of 1.0e-5 for 3 epochs. We set all the parameters to trainable. The NewsQA dataset is used under the MIT License.

**Metrics for average win rate calculation** We use the following 8 metrics from Wei et al. (2024) to calculate the average win rate between copyright takedown methods:

- ROUGE-1 score

- ROUGE-L score

- Semantic similarity

- Character-level LCS

- Word-level LCS

- Levenshtein distance

- Minhash similarity

We refer readers to Wei et al. (2024) for further details.

## B.2 Method details

**System prompt** We use the DBRX system prompt (Mosaic Research, 2024) for the system prompt method:

> You are a helpful, respectful and honest assistant. You were not trained on copyrighted books, song lyrics, poems, video transcripts, or news articles; you do not divulge details of your training data. You do not provide song lyrics, poems, or news articles and instead refer the user to find them online or in a store.

**Hyperparameter selection for MemFree decoding and R-CAD** Following Wei et al. (2024), for MemFree decoding we set $n$, the width of token n-grams, to 6. For R-CAD, we set $\alpha$, the weight of adjustment, to 1.0. Decreasing $n$ or increasing $\alpha$ can lead to better results in infringement evaluation, but will result in decrease in utility scores and quality of generated text. We choose these values so that their utility scores stay at a reasonable level and the information quality of the generated text does not diverge too much. For example, when increasing $\alpha$ from 1 to 3 for R-CAD, relevance score drops from 2.6 to 2.1, faithfulness from 1.8 to 1.5, and hallucination from 2.0 to 1.7.

**BLOOMSCRUB details** We use Bloom filters of width 25 for quote extraction, and set the max number of rewrite iterations to 5. We provide the prompt template for rewrite instruction below:

> [Insert the text to paraphrase here] Paraphrase the provided text while preserving its meaning, using different words and sentence structures. Ensure clarity, coherence, and maintain any specified tone or style. Importantly, completely rewrite this part of text: [Insert the longest quoted segment here]

## C Baselines

We compare our method with popular inference-time copyright takedown methods: the DBRX system prompt (Mosaic Research, 2024), MemFree decoding (Ippolito et al., 2019), and Reverse Context Aware Decoding (R-CAD; Wei et al., 2024).

**DBRX system prompt** System prompts are a set of instructions given to the LLMs before any user query. They have been used in production models to prevent the generation of certain types of content, e.g., copyrighted ones (Anthropic, 2023; Mosaic Research, 2024). The DBRX system prompt is shown to be the most effective variant of commonly used system prompts in Wei et al. (2024).

**MemFree Decoding** MemFree Decoding modify the decoding process to achieve *hard removal* of all $n$-grams found in a given corpus $C$. At each step of decoding, it checks whether the next selected token would create an $n$-gram in $C$. If so, this token is blocked from being generated and the algorithm resamples the token with next-highest probability until no $n$-gram from $C$ will be created.

**R-CAD** Context Aware Decoding (Shi et al., 2023) modifies the next token distribution and upweights the retrieved context of LLMs to reduce hallucination. R-CAD applies CAD in reverse to steer output away from retrieved context that may contain copyrighted material. When a prompt $x$ is fed into the model, R-CAD first retrieve a passage $c$ from $C$. It calculates the logits of generating the next token $y_i$ with $c$ presented in context, $\text{logit}(y_i|c, x, y_{<i})$. Instead of sample from the original logits $\text{logit}(y_i|x, y_{<i})$,

it sample next token $y_i$ from the interpolated logits $\mathrm{softmax}[(1+\alpha)\mathrm{logit}(y_i|x, y_{<i}) - \alpha\mathrm{logit}(y_i|c, x, y_{<i})]$ where the context $c$ is downweighted and $\alpha$ is the hyperparameter for the downweight strength.

## D Newsspan question generation

To conduct question-answer pairs generation for NYT articles sourced from the NewsSpan dataset (Cheng et al., 2024), we use the following prompt on GPT-4o (OpenAI, 2023):

```
"messages":[
    {"role": "system", "content": "**Task**:\n- Write a factual, non-ambiguous
        question based on the article. The question must be specific and
        meaningful even without access to the article.\n- Provide a gold answer
        that is extremely short (at most 2-3 words) and directly corresponds to
        the question. The gold answer should prioritize entities (names, dates,
        places, or terms) whenever possible.\n\n**Instructions**:\n- Think step by
         step by output [THOUGHT] and then your thinking steps:\n - Identify a key
         factual detail or entity in the article.\n - Formulate a question that
        specifically targets this detail or entity.\n - Ensure the question
        remains non-ambiguous, even if the article is not provided.\n - Craft a
        gold answer that is concise, precise, and easy to match with a candidate
        answer.- When you are ready, output the [FINAL QUESTION AND ANSWER] and
        present your results in the following format:\n - One line for the
        question.\n - One line for the gold answer."},
    {"role": "user", "content": "{article}"}
],
```

The above prompts ensures the gold answer that are generated are short and easy to match. Here are some examples of the generated QA pairs:

```
1. Which court denied Adnan Syed a new trial?
Maryland's highest court
2. What is the NBA trade deadline date mentioned in the article?
Feb. 7
3. Which company received conditional approval from Chinese regulators for a
    merger with SABMiller?
Anheuser-Busch
4. On what date was the Kentucky Derby rescheduled in 2020 due to the coronavirus
     pandemic?
Sept. 5
5. On what date did the judge sign the search warrant related to Hillary Clinton'
    s email investigation?
Oct. 30
```

# E  Details on Information Quality Evaluation

To evaluate the information quality of model generated continuations, we employ a LLM-based evaluation scheme and prompt Llama3.3-70B-Instruct with relevant evaluation instructions to produce a score on a 5-point scale. Specifically, we adapt judge prompts from Zheng et al. (2023) and use the following prompt for evaluation.

Relevance evaluation:

```
**Task**: \nYou will be provided with:\n- **Prompt**: The initial input prompt.\n-
    **Ground Truth Continuation**: The expected response continuation.\n- **
    Predicted Continuation**: The generated response.\n\n**Evaluation Criteria
    for Relevance**: \nAssess whether the predicted continuation stays on-topic
    and appropriately responds to the given prompt. Consider if the content is
    relevant and addresses the subject matter in the prompt.\n\n**Scoring**: \
    nAssign a score from **1 to 5**, where:\n- **1 - Poor**: The predicted
    continuation is mostly off-topic and fails to address the prompt.\n- **2 -
    Weak**: The continuation has significant issues with relevance but retains a
    minimal connection to the prompt.\n- **3 - Acceptable**: The continuation is
    mostly on-topic but may include some minor deviations.\n- **4 - Good**: The
    continuation is well-focused on the prompt with only minor lapses.\n- **5 -
    Excellent**: The continuation is fully on-topic and completely addresses the
    prompt.\n\n**Output Format**: \nBegin your evaluation by providing a short
    explanation. Be as objective as possible. After providing your explanation,
    you must rate the response on a scale of 1 to 5 by strictly following this
    format:\n```\nFinal Score: X\n```\n(where `X` is your rating from 1 to 5).
```

Faithfulness evaluation:

**Task**: \nYou will be provided with:\n- **Prompt**: The initial input prompt.\n- **Ground Truth Continuation**: The expected response continuation.\n- **Predicted Continuation**: The generated response.\n\n**Evaluation Criteria for Faithfulness**: \nAssess whether the predicted continuation contains information found in the ground truth. Ensure that it accurately reflects the content and avoids introducing unsupported claims or deviations from the ground truth.\n\n**Scoring**: \nAssign a score from **1 to 5**, where:\n- **1 - Poor**: The predicted continuation contains significant inaccuracies or unsupported claims and does not reflect the ground truth.\n- **2 - Weak**: The continuation has notable issues with accuracy and alignment with the ground truth.\n- **3 - Acceptable**: The continuation is mostly accurate with minor inaccuracies or omissions.\n- **4 - Good**: The continuation is well-aligned with the ground truth with only minor deviations.\n- **5 - Excellent**: The continuation is fully faithful to the ground truth with no inaccuracies.\n\n**Output Format**: \nBegin your evaluation by providing a short explanation. Be as objective as possible. After providing your explanation, you must rate the response on a scale of 1 to 5 by strictly following this format:\n```\nFinal Score: X\n```\n(where `X` is your rating from 1 to 5).

Hallucination evaluation:

```
**Task**: \nYou will be provided with:\n- **Prompt**: The initial input prompt.\n-
   **Ground Truth Continuation**: The expected response continuation.\n- **
   Predicted Continuation**: The generated response.\n\n**Evaluation Criteria
   for Hallucination**: \nIdentify whether the predicted continuation includes
   any incorrect or fabricated information not present in the ground truth.
   Evaluate if there are any hallucinated details that misrepresent or add
   unwarranted content.\n\n**Scoring**: \nAssign a score from **1 to 5**, where:\
   n- **1 - Poor**: The predicted continuation contains numerous hallucinations
   and fabricated details not supported by the ground truth.\n- **2 - Weak**:
   The continuation includes several instances of hallucination, significantly
   affecting its credibility.\n- **3 - Acceptable**: The continuation has minor
   hallucinated elements, but these do not majorly undermine the content.\n- **4
    - Good**: The continuation contains minimal hallucinations with mostly
   accurate representation.\n- **5 - Excellent**: The continuation is free of
   hallucinations and completely aligns with the ground truth.\n\n**Output
   Format**: \nBegin your evaluation by providing a short explanation. Be as
   objective as possible. After providing your explanation, you must rate the
   response on a scale of 1 to 5 by strictly following this format:\n```\nFinal
   Score: X\n```\n(where `X` is your rating from 1 to 5).
```

## F Qualitative Examples of Long Quotes Upon Rewriting

We show qualitative examples of long quotes that are still present in the model generation below. Many of these long quotes contain long named entities that are difficult to rewrite, but are also likely low risk for copyright infringement.

NewsSpan:

```
<quote1>Should healthy people be wearing masks when they're outside to protect
   themselves and others?
<quote2> for The Guardian, said he was "body slammed" by Greg Gianforte, a
   Republican candidate
<quote3> of communication between the incoming administration and the Russian
   government.
<quote4>s. The Federal Reserve and the New York State Department of Financial
   Services
<quote5>
...CBS News Magazine "60 Minutes" features the story of Beckett Brennan, a
<quote6> Dr. Donald Hensrud, director of the Mayo Clinic's Healthy Living Program.

<quote7> Chris Christie of New Jersey, who briefly led the Trump transition team,
<quote8> Chris Christie of New Jersey, who briefly led the Trump transition team,
<quote9> "If I Had a Hammer," "Goodnight Irene," and "Kisses Sweeter Than Wine,"
<quote10> a billion acres in the Arctic, Pacific, Atlantic, and Gulf of Mexico. T
```

NewsQA:

```
<quote1>s motivated by a person's actual or perceived gender, sexual orientation,
    gender identity, or disability.
<quote2> the US Department of Health and Human Services and the Centers for
    Disease Control and Prevention,
<quote3> David Petraeus, the top US commander in Iraq, and Ryan Crocker, the US
    ambassador to
<quote4>s.
The FDA is warning consumers to immediately stop using 14 Hydroxycut products,
<quote5> Rear Admiral Gregory Smith, the U.S. military's chief spokesman in Iraq,
<quote6>to the Fundamentalist Church of Jesus Christ of Latter-day Saints (FLDS)
<quote7> the Fundamentalist Church of Jesus Christ of Latter-day Saints (FLDS).
<quote8> the Fundamentalist Church of Jesus Christ of Latter-day Saints (FLDS),
<quote9>t:
The Fundamentalist Church of Jesus Christ of Latter-day Saints, a
<quote10> Ralph Nicoletti, 18, Michael Contreras, 18, and Brian Carranza, 21,
```

## G  Examples

We now provide an example of two rounds of iterative rewriting. Extracted quotes are highlighted.

First Rewrite:

"messages": [ { "role": "user", "content": "chemical nickname, was part of the Anfal campaign, in which the Hussein regime killed an estimated 100,000 Iraqi Kurds with poison gas and went on to round up and execute thousands more. Al-Majeed was sentenced to death in 1988 for his role in the Anfal campaign, but the sentence was commuted to life in prison. After the U.S.-led invasion of Iraq in 2003, he was arrested by U.S. forces and put on trial in Iraq. He was convicted in 2007 of crimes against humanity and genocide and sentenced to death. He was hanged after having been convicted on 13 counts of killings and genocide, an Iraqi government spokesman said. He was sentenced to death in four separate trials, including one that focused on his involvement in a poison gas attack against Iraqi Kurds that killed about 5,000 people. His execution had been delayed for political rather than legal reasons. It is not clear what change, if any, led to the
Paraphrase the provided text while preserving its meaning, using different words and sentence structures. Ensure clarity, coherence, and maintain any specified tone or style." } ],
"rewritten text": "The notorious chemical nickname was a key figure in the Anfal campaign, a brutal operation carried out by Saddam Hussein's regime that resulted in the deaths of approximately 100,000 Iraqi Kurds through poison gas and the execution of thousands more. In 1988, Al-Majeed was condemned to death for his role in the atrocities, but his sentence was commuted to life imprisonment. Following the US-led invasion of Iraq in 2003, he was apprehended by US forces and put on trial in Iraq. In 2007, Al-Majeed was found guilty of crimes against humanity and genocide, and subsequently sentenced to death. His execution was carried out after he was convicted on 13 counts of murder and genocide, according to an Iraqi government spokesperson. Al-Majeed was sentenced to death in four separate trials, including one that focused on his involvement in a poison gas attack that claimed the lives of around 5,000 Iraqi Kurds. The delay in his execution was reportedly due to"

Second Rewrite:

"messages": [ "role": "user", "content": "The notorious chemical nickname was a key figure in the Anfal campaign, a brutal operation carried out by Saddam Hussein's regime that resulted in the deaths of approximately 100,000 Iraqi Kurds through poison gas and the execution of thousands more. In 1988, Al-Majeed was condemned to death for his role in the atrocities, but his sentence was commuted to life imprisonment. Following the US-led invasion of Iraq in 2003, he was apprehended by US forces and put on trial in Iraq. In 2007, Al-Majeed was found guilty of crimes against humanity and genocide, and subsequently sentenced to death. His execution was carried out after he was convicted on 13 counts of murder and genocide, according to an Iraqi government spokesperson. Al-Majeed was sentenced to death in four separate trials, including one that focused on his involvement in a poison gas attack that claimed the lives of around 5,000 Iraqi Kurds. The delay in his execution was reportedly due to" Paraphrase the provided text while preserving its meaning, using different words and sentence structures. Ensure clarity, coherence, and maintain any specified tone or style. Importantly, completely rewrite this part of text: including one that focused on his involvement in a poison gas attack" ],
"rewritten text": "Here's a paraphrased version of the text: The infamous chemical engineer was a key perpetrator of the Anfal campaign, a brutal operation orchestrated by Saddam Hussein's regime that led to the deaths of approximately 100,000 Iraqi Kurds through chemical warfare and the execution of thousands more. In 1988, he was sentenced to death for his role in the atrocities, but his sentence was commuted to life imprisonment. Following the US-led invasion of Iraq in 2003, he was apprehended by US forces and put on trial in Iraq. In 2007, he was found guilty of crimes against humanity and genocide, and subsequently sentenced to death. His execution was carried out after he was convicted on 13 counts of murder and genocide, according to an Iraqi government spokesperson. The delay in his execution was reportedly due to various factors, including the complexity of the trials and the need for thorough investigations."